
AI Governance in Social Work: A Triple Mandate-Informed Accountability Model

Anonymous Authors¹

Abstract

AI systems are increasingly deployed across public-sector contexts, including social work, informing decisions about risk assessment, resource allocation, and service delivery. Technologies deployed in social work have to operate on its unique characteristics, such as involuntary client engagement, life-altering and often irreversible decisions, relational practice, professional discretion amid moral uncertainty, and working with structurally marginalized populations. Whether existing AI governance frameworks, developed largely for private-sector and routine public-sector contexts, can be effectively applied in social work setting remains under-examined. We identify six operating conditions that prevailing frameworks presuppose and show how each misaligns with social work practice in ways that compromise algorithmic accountability. Drawing on Staub-Bernasconi's triple mandate, we propose a minimum accountability layer organized around three mandate domains (client, organizational, and professional) across key accountability requirements: transparency, contestability, and redress. Our central argument is that effective AI governance must preserve the professional mandate as the independent ethical fulcrum mediating between organizational power and client rights. This work contributes to an emerging conversation about adapting AI governance knowledge to relational, high-stakes service contexts.

1. Introduction

AI systems are increasingly deployed across public-sector contexts, informing decisions about risk assessment, resource allocation, and service delivery (Afrouz & Lucas,

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

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

2023; Lee et al., 2024; Park et al., 2025). Alongside this expansion, a growing body of work on AI governance has proposed frameworks for ensuring such systems operate in ways that are fair, transparent, and accountable (Jobin et al., 2019; Corrêa et al., 2023; Wirtz et al., 2022). These frameworks reflect considerable effort to translate abstract ethical principles into practical mechanisms for oversight and accountability. As deployment expands into an increasingly diverse set of service domains, however, an open question is how well these frameworks, developed primarily with private-sector or routine public-administration contexts in mind, translate into domains with fundamentally different operating conditions.

Social work is one such domain. It operates across the United States through more than 800,000 positions, with state and local government among the largest employers (Bureau of Labor Statistics, 2025), and it is increasingly a site where AI systems inform consequential decisions: predictive risk modeling in child welfare, automated eligibility determination in benefits programs, and algorithmic triage in service referrals (Báez et al., 2026; Brown et al., 2019; Wang et al., 2023). Several features distinguish social work from other public services in ways that matter for governance design. Clients often engage involuntarily, compelled by law, institutional mandate, or the absence of alternatives (Rooney & Mirick, 2018; Redden et al., 2020). Decisions are typically high-stakes and often irreversible, such as removing a child from a family, approving essential benefits, or mandating treatment, with consequences that propagate across a lifetime (Gillingham, 2019; Taylor, 2021). Practice is fundamentally relational, with human relationships serving as the core mechanism through which change occurs rather than a channel for service delivery (Devlieghere et al., 2022; Forenza & Eckert, 2018). Professionals exercise significant discretion as street-level bureaucrats, navigating genuinely irresolvable moral dilemmas under intense public scrutiny (Berrick, 2018; Lipsky, 2010; Ahn et al., 2025). And the populations served are disproportionately affected by poverty, racial marginalization, and structural inequality (Fong, 2017; Roberts, 2014).

High-profile deployment failures suggest what is at stake when these conditions are not accounted for. Michigan's

MiDAS system wrongfully accused tens of thousands of people of unemployment fraud at an extraordinary error rate, stripping them of benefits without adequate notice or opportunity to respond (Charette, 2018). Australia's Robodebt scheme (Mao, 2023) and the Dutch childcare benefits affair (Amnesty International, 2021) produced comparable patterns: wrongful accusations, housing loss, family separation. These failures do not appear to be simply matters of flawed implementation; they suggest that accountability mechanisms designed for other contexts may not deliver equivalent protections when transplanted into relational, high-stakes service domains. AI is not a neutral technology introduced into a blank context: it encodes biases present in existing data, embeds normative assumptions into its models, and often operates with limited transparency (Carter & Dale, 2025; Eubanks, 2019; Keddell, 2019; Mittelstadt et al., 2016). Whether existing AI governance frameworks can effectively address these dynamics in domains like social work is a question that has received limited systematic attention (David, 2024; Wang et al., 2023).

In this paper, we examine how existing AI governance knowledge can be adapted to social work, a domain shaped by distinct ethical obligations, marginalized populations, and practitioner discretion. Our aim is not to reject AI in social work nor to propose a comprehensive alternative framework, but to offer a diagnostic and a starting point for adaptation. We first identify six operating conditions that existing AI governance frameworks tend to presuppose, and examine where and why these conditions break down in social work's distinctive context. Drawing on Staub-Bernasconi's (Staub-Bernasconi, 2007; 2009; 2016) triple mandate as a professional accountability structure, we then propose a set of reflective guiding questions organized around the three mandate domains (client, organizational, and professional) intersecting with accountability requirements of transparency, contestability, and redress. We intend this framework as a minimum accountability layer that can be incorporated into governance approaches applied to social work, and more broadly as an example of how sector-specific conditions can inform governance design.

2. AI Governance: Key Elements and Operating Conditions

2.1. Definition and key elements

Among many efforts to define AI governance, the most prominent is "a system of rules, practices, processes, and technological tools that are employed to ensure an organization's use of AI technologies aligns with the organization's strategies, objectives, and values; fulfills legal requirements; and meets principles of ethical AI followed by the organization" (Mäntymäki et al., 2022, p. 604). Operationalized through tools, norms, standards, and regulations (Butcher &

Beridze, 2019; Corrêa et al., 2023; Jobin et al., 2019), AI governance frameworks commonly converge around four central themes: technology, stakeholders and context, regulation, and processes (Birkstedt et al., 2023; Mäntymäki et al., 2022). Among these, accountability warrants particular attention, as AI governance frameworks are, at their core, mechanisms for ensuring that those who build and deploy AI systems can be held answerable for their outcomes (Raji et al., 2020).

2.2. Operating conditions of AI governance frameworks

While a growing body of literature examines AI governance in the public sector, including human services (Kuziemski & Misuraca, 2020; Robles & Mallinson, 2025; Wang et al., 2023), much of this work implicitly assumes conditions that prevails in private-sector services or routine public administration. We identify six such operating assumptions that recur across prevailing frameworks.

First, perhaps the most widely shared assumption is a technosolutionist orientation (Morozov, 2013) that treats AI deployment as intrinsically desirable, centering on optimizing implementation rather than evaluating whether it is warranted (Attard-Frost & Lyons, 2025; Corrêa et al., 2023; Wirtz et al., 2022). Even frameworks that question the extent of human-machine collaboration focus on the degree of automation (Wang et al., 2023), rather than whether algorithmic decision-making is appropriate for the specific domain. This orientation works well where AI benefits are well-supported and automating judgment does not alter the fundamental nature of the service, as in traffic signal optimization, weather forecasting, or routine document processing.

Second, AI governance frameworks often assume that AI risks can be classified into distinct, predefined categories and addressed through structured, sequential mitigation (Wirtz et al., 2020; 2022). Frameworks operationalize this through dimensional risk taxonomies, modular governance layers mapped to technical, ethical, and legal domains (Gasser & Almeida, 2017), and lifecycle-based models (Batool et al., 2025; Mäntymäki et al., 2022; Raji et al., 2020). This works where risk factors are sufficiently independent to be addressed in isolation, an assumption that holds for standardized tasks but becomes unstable where risks are entangled with human relationships and contextual judgment.

Third, improving predictive accuracy is widely treated as the path to better service outcomes (McCadden et al., 2025; Van Amsterdam et al., 2025). Embedded are three premises: outcomes of interest are accurately captured in available data, more accurate models lead to better decisions, and prediction does not itself shape the outcome being predicted (Kolt et al., 2025; Metcalf et al., 2021). These premises hold where data reliably measure what they claim to measure

and where predictions inform decisions without reshaping the conditions they describe. They do not hold where data capture institutional behavior rather than the underlying phenomenon of interest (Farrell, 2025; Hu, 2025; McNellan et al., 2022), or where labeling someone as high-risk triggers responses that worsen the very outcome predicted.

Fourth, human oversight is treated as a primary safeguard, expecting humans to detect errors, exercise independent judgment, and override flawed algorithmic decisions (Mišić et al., 2025; Office of Management and Budget, 2024). This requires that AI systems provide explanations operators can understand and act on, that operators have time, expertise, and authority to evaluate outputs independently rather than defer, and that decision-making pace allows deliberation. These conditions erode where operators face high caseloads and time pressure (Kuziemski & Misuraca, 2020), and where AI systems remain opaque with model internals protected by intellectual property restrictions (Camilleri, 2024; David, 2024).

Fifth, many existing AI governance frameworks presuppose at the organizational level that institutions have, or can readily develop, the infrastructure to execute governance requirements, including dedicated AI governance teams, compliance mechanisms, and clear accountability lines (Mäntymäki et al., 2023). The U.S. OMB memorandum (Office of Management and Budget, 2024), for instance, requires each federal agency to designate a Chief AI Officer and convene agency-wide AI governance bodies. Yet these requirements assume organizational capacity many agencies cannot readily meet, widening disparities across the service landscape and leaving smaller, under-resourced agencies without realistic compliance means (Attard-Frost & Lyons, 2025).

Finally, governance frameworks often assume that affected populations have meaningful agency over engagement with AI systems, such that individuals can consent to or refuse AI-driven services, opt out without penalty, and appeal through institutional channels (Corrêa et al., 2023; Fjeld et al., 2020; Jobin et al., 2019). These conditions hold only where participation is voluntary and alternatives are available.

The above assumptions are not inherently flawed; they reflect the contexts in which AI governance frameworks were developed and serve important functional purposes within those contexts. Yet even the frameworks designed for the general public service context often carry forward the assumptions drawn on the private sector and administrative dimensions of the public sector, without examining whether they hold equally across all public sector domains, including social work.

3. Where Governance Assumptions Meet Social Work's Distinctive Context

Unlike the private sector, the public sector is accountable to citizens and legislators rather than shareholders, obligated to protect rights rather than maximize profit, and required to be transparent rather than permitted to withhold information as proprietary (Kuziemski & Misuraca, 2020; Mišić et al., 2025). Yet the public sector is not monolithic. Social work is distinguished even within the public sector by the involuntary nature of client engagement, life-altering high-stakes decisions, and the structural marginalization of populations served. While no single characteristic is entirely unique to social work, their combination creates a context where the operating conditions assumed by standard AI governance frameworks require re-examination. Figure 1 summarizes the five distinctive conditions we discuss below and the seven points of misalignment they produce with standard governance assumptions.

3.1. Involuntary engagement

Among public sector services, social work is distinguished by the involuntary nature of client engagement. Families under child protection investigation, individuals subject to court-ordered treatment, and people dependent on welfare benefits do not choose to engage but are compelled by law, institutional mandate, or the absence of alternatives, producing relationships defined by profound asymmetries of power (Redden et al., 2020; Rooney & Mirick, 2018). Unlike other public services where voluntariness varies but rarely reaches this extreme, social work clients cannot opt out without penalty, cannot choose alternative providers, and face legal consequences for disengagement. When the state controls access to vital resources, consent to data collection is not a genuine choice, as refusal risks losing essential resources or triggering legal sanctions (Garrett, 2025; Varon & Peña, 2021). This becomes particularly consequential when AI systems are introduced: those who cannot leave the system and lack the ability to remove their data have no means to prevent, or even be informed of, how information collected in the context of seeking help is repurposed for algorithmic risk profiling, transforming the relationship from support into surveillance (Keddell, 2019).

3.2. High-stakes, life-altering consequences

In social work, decision stakes are often life-altering. Removing a child from a family, approving essential benefits, or mandating treatment carries profound consequences even when warranted (Gillingham, 2019; Taylor, 2021). A wrongful removal or a failure to prevent harm can have lasting effects across a lifetime. Michigan's MiDAS system, designed to automate unemployment insurance claims without human oversight, wrongfully accused tens of thousands of people of

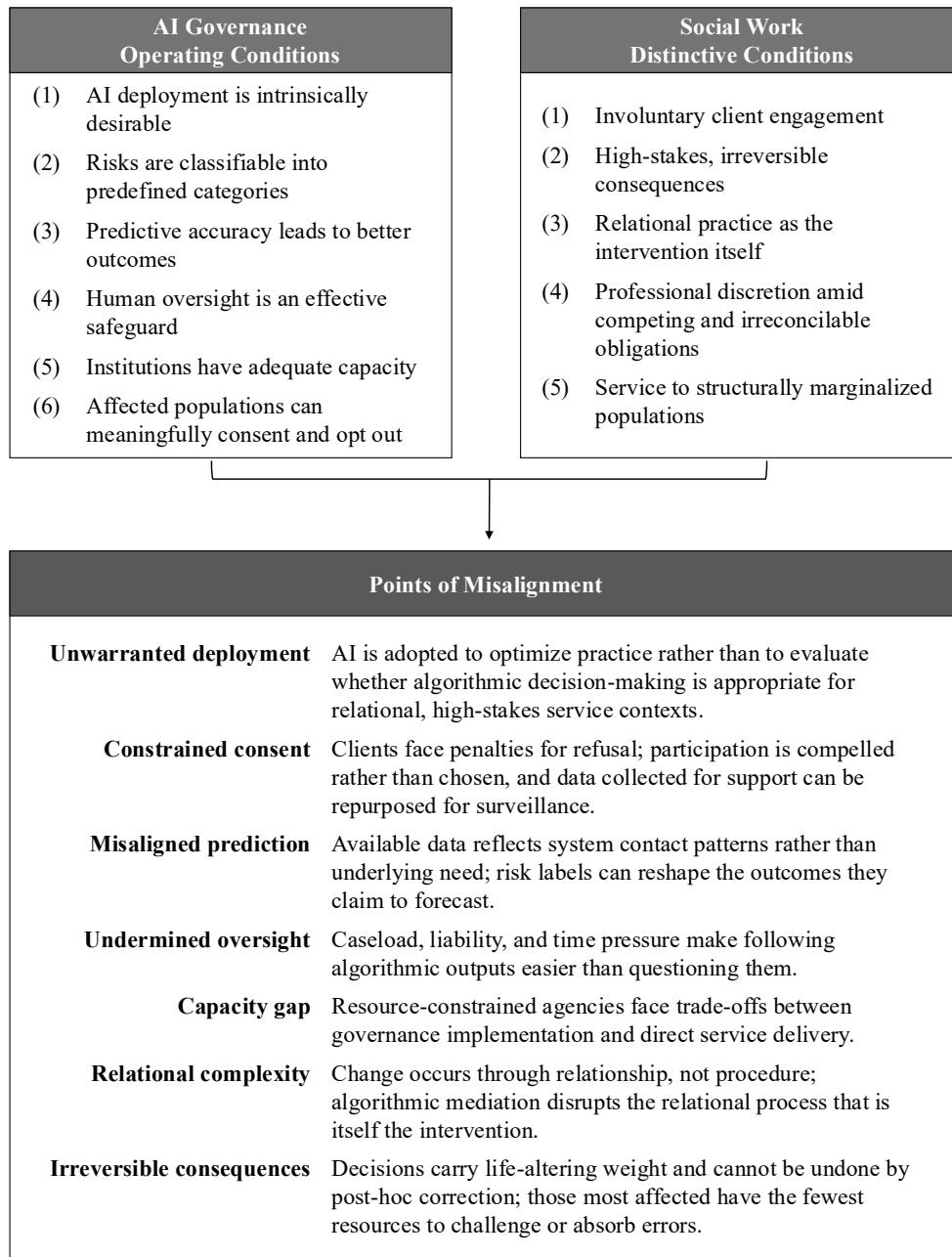


Figure 1. Where standard AI governance operating conditions (left) meet social work’s distinctive conditions (right). The bottom panel summarizes seven points of misalignment that emerge when frameworks developed for other contexts are applied to relational, high-stakes service domains.

fraud at an extraordinarily high error rate, stripping them of benefits and imposing severe financial penalties without adequate notice or opportunity to respond (Charette, 2018; Gray & Farrington, 2018). Australia's Robodebt scheme (Mao, 2023) and the Dutch childcare benefits affair (Amnesty International, 2021) similarly produced wrongful accusations, housing loss, and family separation. In each case, the harm was not simply that the system erred, but that it erred in a domain where decisions already carry extraordinary weight and where those affected had no meaningful capacity to challenge or avoid the consequences.

3.3. Relational practice as the intervention itself

While public services vary in how relational they are, social work sits at the far end of this spectrum. The profession is grounded in the principle that human relationships are not merely a vehicle for service delivery but the core mechanism through which change occurs (Devlieghere et al., 2022; Forenza & Eckert, 2018; Mishna et al., 2021). Outcomes depend not only on what decision is reached but on how the client experiences the process of reaching it, which requires engaging with individuals as whole persons, reading non-verbal cues, interpreting emotional states, and drawing on trust built over time (Garrett, 2025; Gillingham & Graham, 2017). Within each interaction, empathy, assessment, intervention, and relationship maintenance occur simultaneously and cannot be decomposed into discrete, documentable actions. These characteristics sit uneasily with AI governance approaches built around classifiable data, standardizable procedures, and measurable accuracy.

3.4. Professional discretion amid uncertainty and public scrutiny

Unlike some public services where decisions follow codifiable rules with limited interpretation (e.g., a postal worker routes a package, a parking enforcement officer issues a citation), social work requires practitioners to exercise significant discretion in every encounter (Berrick, 2018).¹ The problems social workers address, such as poverty, substance use, family crisis, and child safety, are complex, chronic, and rarely have a single correct resolution (Berrick, 2018; Lipsky, 2010). In child welfare, workers must decide whether to remove a child for safety, knowing that removal itself carries harm, under conditions where no solution may be clearly right at the point of decision (Berrick, 2018).

This discretion is exercised under intense public scrutiny. When decisions are associated with severe injury or death of a child, social workers may face not only institutional review and media criticism but criminal charges, as when child welfare workers in Los Angeles faced prosecution fol-

¹Anonymous self-citation omitted for review.

lowing the death of Gabriel Fernandez (Etehad & Winton, 2017). Such scrutiny has been associated with defensive practice: over-documenting, intervening more than necessary, or refraining from intervening when service is required (Chenot, 2011; Whittaker & Havard, 2016). AI governance frameworks built for conditions where decisions can be standardized and optimized for accuracy offer little guidance where "right" decisions may not exist, and the practitioner bears consequences for getting it wrong.

3.5. Serving populations already affected by structural inequality

Social work disproportionately works with populations affected by poverty, racial marginalization, and structural inequality, whose circumstances are often compounded by heightened state surveillance and the lasting consequences of prior system involvement (Fong, 2017; Hailu et al., 2022; Roberts, 2014). Serving these populations is central to the profession's mission (National Association of Social Workers, 2021), and their concentration of structural disadvantage intensifies every condition above: involuntary engagement carries greater weight when those compelled to participate already lack resources to advocate for themselves; irreversible decisions are more devastating when there is no safety net to absorb their impact; and relational trust is harder to build when clients have been repeatedly failed by the systems meant to help them.

Beyond amplifying these conditions, working with structurally marginalized populations introduces distinct challenges for data-driven governance. Data generated through social work practice reflects not only genuine need but also patterns shaped by historical bias and disproportionate surveillance, including from adjacent systems such as criminal justice, where low-income and racial minority communities are overrepresented (Gillingham, 2019; Neil & Zanger-Tishler, 2025). The populations most affected by algorithmic systems are also the least positioned to participate in the governance processes that shape them or to contest decisions made about their lives (Garkisch & Goldkind, 2024; Ahn, 2025).

4. A Guiding Model for AI Governance in Social Work

4.1. From misalignment to accountability

Ensuring accountability has been recognized as a central purpose of AI governance in organizational and technical contexts (Busuioc, 2021; Raji et al., 2020), yet what this requires in social work practice remains largely unexamined. When AI systems shape decisions that affect fundamental rights and family integrity, accountability raises questions specific to this context: who is responsible for those de-

275 cisions, to whom justification is owed, and through what
 276 mechanisms that responsibility can be enforced (Novelli
 277 et al., 2024). Accountability is a relational concept in which
 278 a decision-maker must explain and justify their conduct to
 279 those affected or to an oversight authority, which in turn
 280 can question, evaluate, and impose consequences (Bovens,
 281 2007). For a meaningful accountability process, three condi-
 282 tions must be met: the decision-maker's conduct must
 283 be transparent enough to be examined (information); the
 284 reasoning behind decisions must be open to questioning and
 285 justification (explanation); and decision-makers must face
 286 the possibility of consequences, whether sanctions, required
 287 remediation, or redress, when an oversight authority judges
 288 their conduct to have fallen short (Bovens, 2007; Busuioc,
 289 2021).

290 Where AI systems are integrated into social work, each
 291 accountability requirement is compromised in ways that
 292 follow directly from the misalignments identified above. In-
 293 formation is constrained because involuntary clients who
 294 cannot opt out and lack meaningful agency often have no ac-
 295 cess to how algorithmic outputs shape decisions about their
 296 lives. Explanation is undermined because system opacity,
 297 compounded by caseload pressure and liability fear, makes
 298 it difficult to produce or demand meaningful justification
 299 for algorithmic recommendations. Consequences are struc-
 300 turally compromised at two levels: oversight authorities
 301 rarely have sufficient access to how algorithmic systems
 302 reach their conclusions, and the populations most affected
 303 are the least positioned to seek redress through institutional
 304 appeals, a structural exclusion that extends to the fundamen-
 305 tal question of who has standing to contest decisions made
 306 about their lives (Garkisch & Goldkind, 2024; Ahn, 2025).

307 4.2. The triple mandate as a professional accountability 308 framework

309 Addressing these accountability failures requires under-
 310 standing how professional obligation in social work is struc-
 311 tured. The triple mandate—simultaneous obligations to the
 312 client, the employing organization or state, and the pro-
 313 fession itself (Staub-Bernasconi, 2007; 2009; 2016)—is a
 314 foundational feature of professional identity that remains
 315 essential in the AI era (Garkisch & Goldkind, 2024). What
 316 makes it significant is not the complications posed by hav-
 317 ing three dimensions, but the distinctive role of the third.
 318 Without the professional mandate, social work reduces to
 319 a double mandate where the practitioner becomes either an
 320 instrument of the state or an advocate for the client, with no
 321 independent ethical standing to mediate between the two. It
 322 is the professional mandate, grounded in ethical principles
 323 and disciplinary knowledge, that provides the normative
 324 basis to push back against both organizational directives and
 325 client preferences when professional judgment demands it
 326 (Staub-Bernasconi, 2009; 2016).
 327
 328
 329

In practice, these three obligations frequently pull in oppos-
 ing directions (Staub-Bernasconi, 2009; 2016), and AI shifts
 the balance of power within each tension. Between client
 and organization, organizational mandates demand compli-
 ance, efficiency, and risk management, while the client man-
 date demands self-determination, dignity, and minimal inter-
 vention. AI intensifies this tension by equipping organiza-
 tions with seemingly objective justification for intervention
 that involuntary clients have virtually no capacity to contest
 (Garkisch & Goldkind, 2024). Between client and practi-
 tioner, professional expertise and client self-determination
 can already conflict; AI compounds this by layering algori-
 thmic authority onto professional judgment, further reducing
 the client's space to dispute decisions (Busuioc, 2021). Be-
 tween practitioner and organization, professionals require
 discretion, contextual judgment, and time for relationships,
 while organizations seek standardization and efficiency. AI
 can create an asymmetric structure where following an al-
 gorithmic recommendation requires no justification, while
 overriding it demands documentation, explanation, and ex-
 posure to scrutiny (Busuioc, 2020).

Across all three tensions, the consistent pattern is the system-
 atic weakening of the professional mandate. AI strengthens
 organizational accountability through compliance metrics
 and efficiency measures while simultaneously eroding pro-
 fessional discretion, relational judgment, and ethical media-
 tion on which the professional mandate depends (Bovens,
 2007; Busuioc, 2021). As that mediating role is undermined,
 the triple mandate risks collapsing into a double mandate,
 leaving client and organization in direct, unmediated op-
 position where the power asymmetry between involuntary
 clients and the state becomes unmitigated. Preserving the
 profession's mediating role is therefore not merely a pro-
 fessional interest but a governance imperative that must be
 addressed deliberately in the design of AI governance for
 social work.

4.3. A triple mandate-informed guiding model

Our analysis established that the misalignment between
 AI governance operating conditions and social work prac-
 tice represents fundamental accountability failures (Busuioc,
 2021; Novelli et al., 2024), and that the triple mandate pro-
 vides the structural foundation through which these failures
 can be addressed (Staub-Bernasconi, 2009). While the triple
 mandate originally describes the competing obligations nav-
 igated by individual practitioners, we adapt it here to the
 governance level, arguing that AI governance in social work
 must preserve the conditions for fulfilling all three man-
 dates when algorithmic systems are integrated into practice
 (Garkisch & Goldkind, 2024).

Each mandate defines a core principle that AI governance
 must protect: client self-determination, dignity, and rights;

Table 1. Guiding Model for AI Governance in Social Work: Accountability Requirements Across the Triple Mandate.

Triple Mandate Domain	Guiding Questions	Practical Considerations
Client Mandate <i>Protecting client self-determination, dignity, and rights in AI-mediated services</i>	<p><i>Transparency:</i> Can clients receive understandable explanations of how algorithmic outputs influenced decisions about their lives?</p> <p><i>Contestability:</i> Is there an independent pathway for clients to question and contest AI-informed decisions?</p> <p><i>Redress:</i> Are there meaningful mechanisms for redress, including human-led appeals processes, when algorithmic errors cause harm?</p>	<p>Transparency may be technically infeasible for complex algorithms; meaningful appeals slow efficiency and increase costs; involuntary clients face power asymmetries that limit their capacity to exercise these rights even when formally available. Clients should be treated as rights-holders first, with their dignity placed above administrative efficiency.</p>
Organizational Mandate <i>Institutional accountability for responsible AI deployment and its impacts</i>	<p><i>Transparency:</i> Is the organization systematically monitoring AI system performance and its effects on the balance among the three mandates?</p> <p><i>Contestability:</i> Can the organization justify its AI adoption and operation in terms of preserving the triple mandate, not only technical performance or legal compliance?</p> <p><i>Redress:</i> Is the organization conducting regular audits for fairness and bias, and taking corrective action when structural inequalities or systematic harms are identified?</p>	<p>Monitoring and auditing require dedicated resources and technical expertise many social work agencies lack; procurement rarely includes social work-specific criteria; transparency about failures risks reputational damage.</p>
Professional Mandate <i>Protecting ethical judgment, discretion, and knowledge-based practice in AI-mediated decision-making</i>	<p><i>Transparency:</i> Are the assumptions, limitations, and value trade-offs embedded in the algorithm transparent to practitioners?</p> <p><i>Contestability:</i> Do practitioners have protected authority to override algorithmic recommendations, with AI positioned as a support tool rather than a replacement for professional judgment?</p> <p><i>Redress:</i> Is critical AI literacy supported through education, time, supervision, and organizational culture that values professional reasoning over algorithmic compliance?</p>	<p>Override authority is meaningless without supportive organizational culture; proprietary restrictions may limit algorithmic transparency; sustained investment in training is needed to prevent deskilling over time.</p>

institutional accountability for responsible AI deployment and its impacts; and professional ethical judgment, discretion, and knowledge-based practice (Staub-Bernasconi, 2009; Garkisch & Goldkind, 2024). Drawing on the misalignments identified above, we propose guiding questions organized at the intersection of the mandate domains and three accountability requirements: transparency, contestability, and redress (Table 1). Each guiding question asks whether the conditions necessary for meaningful accountability are present within a given mandate domain. These questions do not constitute a comprehensive governance framework but a minimum accountability layer that any AI governance approach must address when applied in social work contexts.

The model also carries a broader structural argument, illustrated in Figure 2. When AI systems override professional expertise (Figure 2a), the professional mandate is weakened, tensions between mandates are distorted, and the triple mandate risks collapsing into a double mandate in which clients face organizational power without professional mediation. When professional expertise governs AI tools (Figure 2b), tensions remain but are manageable, and the triple mandate is preserved.

This structural argument rests on a fundamental premise: the professional mandate is embodied in people, not systems. As AI becomes more influential in social work, the demands on professional judgment increase rather than diminish. Practitioners must now exercise judgment not only about complex human situations but also about when to trust,

question, or reject algorithmic inputs—a more demanding form of professional judgment that requires sustained investment in education, organizational support, and workforce capacity. Efficiency gains from AI must therefore be reinvested in professional capacity rather than used to justify workforce reduction or further automation (Ahn, 2025). AI governance in social work that does not invest in the people who enact it will find the triple mandate hollowed out, the accountability structure collapsed, and the people social work aims to serve will bear the consequences.

5. Discussion and Conclusion

In this paper, we have examined how existing AI governance frameworks operate when applied to social work, a context with distinctive features that differ in important ways from the settings where these frameworks were originally developed. Our analysis suggests that the transfer is not straightforward: the populations social work serves are disproportionately involuntarily involved, structurally marginalized, and less able to exercise the kinds of agency that existing frameworks presume; decisions carry consequences that cannot be reversed by post-hoc correction; and the relational, discretionary nature of practice resists the standardization that algorithmic governance often requires. These are not problems of implementation but questions about the conditions under which accountability itself can be sustained.

By adapting Staub-Bernasconi’s (Staub-Bernasconi, 2007; 2009; 2016) triple mandate from an individual practice prin-

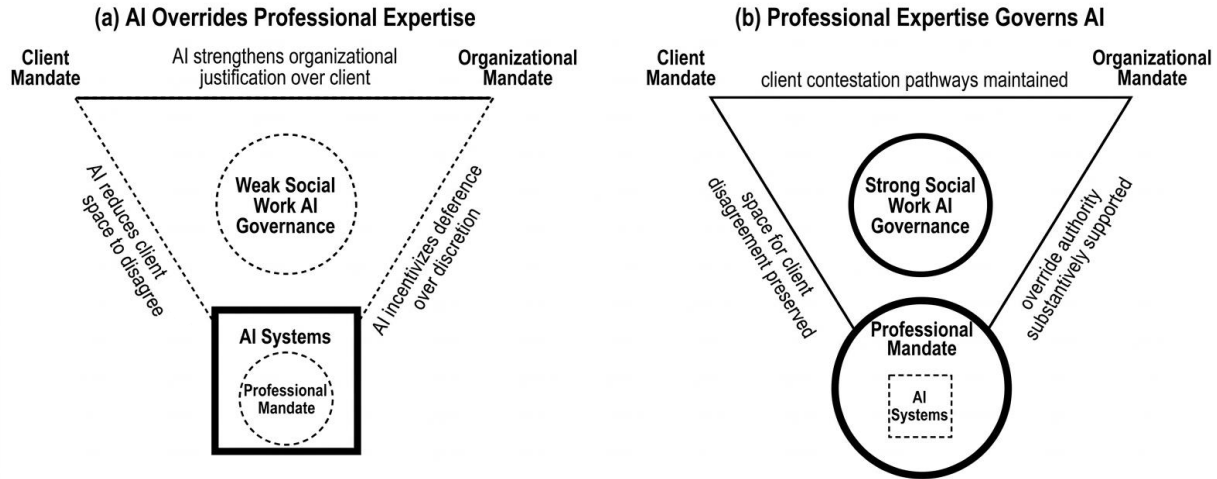


Figure 2. The triple mandate under two AI governance configurations. (a) When AI systems override professional expertise, the professional mandate is weakened (dashed circle), tensions between mandates become unmediated, and the triple mandate collapses into a double mandate in which clients face organizational power without professional mediation. (b) When professional expertise governs AI, the professional mandate remains strong (solid circle), AI tools operate within professional authority, and the triple mandate is preserved with client contestation pathways maintained.

ciple to a governance-level framework, we proposed a minimum accountability layer that organizes AI accountability along three mandate domains (client, organizational, and professional) intersecting with three accountability requirements (Bovens, 2007; Busuioc, 2021). The model’s central argument is that the professional mandate functions as the independent ethical fulcrum holding the triple mandate together. When AI weakens this fulcrum, the triple mandate risks collapsing into a double mandate, leaving involuntary, marginalized clients and the state in direct opposition without professional mediation.

We offer this work as a contribution to ongoing dialogue with several implications seem worth noting. First, the diagnostic here may generalize beyond social work. Other high-stakes, relational public-sector domains, such as mental health services, re-entry programs, and public benefits adjudication, share common structural features including power asymmetry, irreversibility, and professional mediation as an ethical fulcrum. Examining whether similar adaptations apply in those settings is an open question for future work. Second, the analysis points to concrete research directions for the machine learning community working on governance and trustworthy AI: designing transparency mechanisms meaningful to non-expert, involuntary users; building contestability pathways that are not gated by institutional power asymmetries; measuring when model deployment reshapes the outcomes being predicted; and developing evaluation frameworks that capture relational and procedural harms alongside accuracy-based metrics. Third, and most broadly, we hope this work contributes to an emerging conversation about how AI governance knowledge, substantial and still

developing, can be thoughtfully adapted to the diverse contexts in which AI systems are now being deployed. This adaptation is unlikely to succeed as a unilateral effort from any single community; it calls for sustained collaboration between AI researchers, domain practitioners, and those most affected by algorithmic systems.

References

Afrouz, R. and Lucas, J. A systematic review of technology-mediated social work practice: Benefits, uncertainties, and future directions. *Journal of Social Work*, 23(5): 953–974, 2023. doi: 10.1177/14680173231165926.

Ahn, E. Who is the human in human-centered AI? *AI & Society*, 2025. doi: 10.1007/s00146-025-02825-6.

Ahn, E., Morstatter, F., Waters-Roman, D., Palmer, L., and McCroskey, J. Qualitative exploration of child welfare workers’ decision-making experiences and perspectives on fairness. *Journal of Public Child Welfare*, 19(1):229–252, 2025. doi: 10.1080/15548732.2024.2312846.

Amnesty International. Xenophobic machines: Discrimination through unregulated use of algorithms in the dutch childcare benefits scandal. Technical Report EUR 35/4686/2021, Amnesty International, 2021.

Attard-Frost, B. and Lyons, K. AI governance systems: A multi-scale analysis framework, empirical findings, and future directions. *AI and Ethics*, 5(3):2557–2604, 2025. doi: 10.1007/s43681-024-00569-5.

- 440 Báez, J. C., Ahn, E., Tamietti, A., Victor, B. G., and Gold-
 441 kind, L. Clinical social workers' perceptions of large
 442 language models in practice. *Journal of Evidence-Based*
 443 *Social Work*, 23(1):42–63, 2026.
- 444
 445 Batool, A., Zowghi, D., and Bano, M. AI governance: A
 446 systematic literature review. *AI and Ethics*, 5(3):3265–
 447 3279, 2025.
- 448
 449 Berrick, J. D. *The Impossible Imperative: Navigating the*
 450 *Competing Principles of Child Protection*. Oxford Uni-
 451 versity Press, 2018.
- 452
 453 Birkstedt, T., Minkinen, M., Tandon, A., and Mäntymäki,
 454 M. AI governance: Themes, knowledge gaps and future
 455 agendas. *Internet Research*, 33(7):133–167, 2023.
- 456
 457 Bovens, M. Analysing and assessing accountability: A
 458 conceptual framework. *European Law Journal*, 13(4):
 459 447–468, 2007.
- 460
 461 Brown, A., Chouldechova, A., Putnam-Hornstein, E., Tobin,
 462 A., and Vaithianathan, R. Toward algorithmic account-
 463 ability in public services: A qualitative study of affected
 464 community perspectives on algorithmic decision-making
 465 in child welfare services. In *Proceedings of the 2019 CHI*
 466 *Conference on Human Factors in Computing Systems*, pp.
 467 1–12, 2019.
- 468
 469 Bureau of Labor Statistics. Social workers. U.S.
 470 Department of Labor, Occupational Outlook
 471 Handbook, 2025. [https://www.bls.gov/
 472 ooh/community-and-social-service/
 473 social-workers.htm](https://www.bls.gov/ooh/community-and-social-service/social-workers.htm).
- 474
 475 Busuioc, M. Accountable artificial intelligence: Holding
 476 algorithms to account. *Public Administration Review*,
 477 2020.
- 478
 479 Busuioc, M. Accountable artificial intelligence: Holding
 480 algorithms to account. *Public Administration Review*, 81
 481 (5):825–836, 2021.
- 482
 483 Butcher, J. and Beridze, I. What is the state of artificial
 484 intelligence governance globally? *RUSI Journal*, 164
 485 (5–6):88–96, 2019.
- 486
 487 Camilleri, M. A. Artificial intelligence governance: Ethical
 488 considerations and implications for social responsibility.
 489 *Expert Systems*, 41(7):e13406, 2024.
- 490
 491 Carter, S. O. and Dale, J. G. Social bias in AI: Re-coding
 492 innovation through algorithmic political capitalism. *AI &*
 493 *Society*, 2025.
- 494
 495 Charette, R. N. Michigan's MiDAS unemployment sys-
 496 tem: Algorithm alchemy created lead, not gold. *IEEE*
 497 *Spectrum*, January 2018.
- 498
 499 Chenot, D. The vicious cycle: Recurrent interactions among
 500 the media, politicians, the public, and child welfare ser-
 501 vices organizations. *Journal of Public Child Welfare*, 5
 502 (2–3):167–184, 2011.
- 503
 504 Corrêa, N. K., Galvão, C., Santos, J. W., Del Pino, C., Pinto,
 505 E. P., Barbosa, C., Massmann, D., Mambrini, R., Galvão,
 506 L., Terem, E., and De Oliveira, N. Worldwide AI ethics:
 507 A review of 200 guidelines and recommendations for AI
 508 governance. *Patterns*, 4(10):100857, 2023.
- 509
 510 David, G. Artificial intelligence: Opportunities and chal-
 511 lenges for public administration. *Canadian Public Ad-
 512 ministration*, 67(3):388–406, 2024.
- 513
 514 Devlieghere, J., Gillingham, P., and Roose, R. Dataism
 515 versus relationshipism: A social work perspective. *Nordic*
 516 *Social Work Research*, 12(3):328–338, 2022.
- 517
 518 Etehad, M. and Winton, R. Gabriel Fernandez case prosecu-
 519 tion. Los Angeles Times, 2017.
- 520
 521 Eubanks, V. *Automating Inequality: How High-Tech Tools*
 522 *Profile, Police, and Punish the Poor*. St. Martin's Press,
 523 2019.
- 524
 525 Farrell, H. AI as governance. *Annual Review of Political*
 526 *Science*, 28(1):375–392, 2025.
- 527
 528 Fjeld, J., Achten, N., Hilligoss, H., Nagy, A., and Srikumar,
 529 M. Principled artificial intelligence: Mapping consensus
 530 in ethical and rights-based approaches to principles for
 531 AI. *SSRN Electronic Journal*, 2020.
- 532
 533 Fong, K. Child welfare involvement and contexts of poverty.
 534 *Children and Youth Services Review*, 72:5–13, 2017.
- 535
 536 Forenza, B. and Eckert, C. Social worker identity: A pro-
 537 fession in context. *Social Work*, 63(1):17–26, 2018.
- 538
 539 Garkisch, M. and Goldkind, L. Considering a unified model
 540 of artificial intelligence enhanced social work: A system-
 541 atic review. *Journal of Human Rights and Social Work*,
 542 10(1):23–42, 2024. doi: 10.1007/s41134-024-00326-y.
- 543
 544 Garrett, P. M. 'magic moments': AI and the 'disappearance'
 545 of social work ethics? *The British Journal of Social Work*,
 546 2025. doi: 10.1093/bjsw/bcaf230.
- 547
 548 Gasser, U. and Almeida, V. A. F. A layered model for
 549 AI governance. *IEEE Internet Computing*, 21(6):58–62,
 550 2017.
- 551
 552 Gillingham, P. Can predictive algorithms assist decision-
 553 making in social work with children and families? *Child*
 554 *Abuse Review*, 28(2):114–126, 2019.

- 495 Gillingham, P. and Graham, T. Big data in social welfare:
496 The development of a critical perspective on social work's
497 latest 'electronic turn'. *Australian Social Work*, 70(2):
498 135–147, 2017.
- 499
500 Gray, S. and Farrington, C. Opinion: Undoing the harm of
501 MiDAS' fraud designations. *The Detroit News*, 2018.
- 502
503 Hailu, E. M., Maddali, S. R., Snowden, J. M., Carmichael,
504 S. L., and Mujahid, M. S. Structural racism and adverse
505 maternal health outcomes: A systematic review. *Health
506 & Place*, 78:102923, 2022.
- 507
508 Hu, L. What is new, and what is old, in fairness and machine
509 learning. *ACM Journal on Responsible Computing*, 2025.
- 510
511 Jobin, A., Ienca, M., and Vayena, E. The global landscape
512 of AI ethics guidelines. *Nature Machine Intelligence*, 1
513 (9):389–399, 2019.
- 514
515 Keddell, E. Algorithmic justice in child protection: Sta-
516 tistical fairness, social justice and the implications for
517 practice. *Social Sciences*, 8(10):281, 2019.
- 518
519 Kolt, N., Shur-Ofry, M., and Cohen, R. Lessons from com-
520 plex systems science for AI governance. *Patterns*, 6(8):
521 101341, 2025.
- 522
523 Kuziemski, M. and Misuraca, G. AI governance in the pub-
524 lic sector: Three tales from the frontiers of automated
525 decision-making in democratic settings. *Telecommunica-
526 tions Policy*, 44(6):101976, 2020.
- 527
528 Lee, J. Y., Ahn, E., Xu, A., Yang, Y., Chang, Y., Cha,
529 H., and Ammari, T. Artificial intelligence in applied
530 family research involving families with young children:
531 A scoping review. *Family Relations*, 2024.
- 532
533 Lipsky, M. *Street-Level Bureaucracy: Dilemmas of the
534 Individual in Public Services*. Russell Sage Foundation,
535 30th anniversary expanded edition, 2010.
- 536
537 Mäntymäki, M., Minkkinen, M., Birkstedt, T., and Viljanen,
538 M. Defining organizational AI governance. *AI and Ethics*,
539 2(4):603–609, 2022.
- 540
541 Mäntymäki, M., Minkkinen, M., Birkstedt, T., and Vilja-
542 nen, M. Putting AI ethics into practice: The hourglass
543 model of organizational AI governance. *arXiv preprint
544 arXiv:2206.00335*, 2023.
- 545
546 Mao, F. Robodebt: Illegal Australian welfare scheme inves-
547 tigation. *BBC News*, 2023.
- 548
549 McCradden, M. D., Mazwi, M. L., and Oakden-Rayner, L.
Can an accurate model be bad? *Patterns*, 6(4):101205,
2025.
- McNellan, C. R., Gibbs, D. J., Knobel, A. S., and Putnam-
Hornstein, E. The evidence base for risk assessment tools
used in U.S. child protection investigations: A systematic
scoping review. *Child Abuse & Neglect*, 134:105887,
2022.
- Metcalf, J., Moss, E., Watkins, E. A., Singh, R., and Elish,
M. C. Algorithmic impact assessments and accountabil-
ity: The co-construction of impacts. In *Proceedings of
the 2021 ACM Conference on Fairness, Accountability,
and Transparency*, pp. 735–746, 2021.
- Mishna, F., Milne, E., Bogo, M., and Pereira, L. F. Respond-
ing to COVID-19: New trends in social workers' use
of information and communication technology. *Clinical
Social Work Journal*, 49(4):484–494, 2021.
- Mišić, J., Van Est, R., and Kool, L. Good governance of
public sector AI: A combined value framework for good
order and a good society. *AI and Ethics*, 5(5):4875–4889,
2025.
- Mittelstadt, B. D., Allo, P., Taddeo, M., Wachter, S., and
Floridi, L. The ethics of algorithms: Mapping the debate.
Big Data & Society, 3(2), 2016.
- Morozov, E. *To Save Everything, Click Here: The Folly of
Technological Solutionism*. PublicAffairs, 2013.
- National Association of Social Workers. Code of ethics
of the National Association of Social Workers. NASW
Press, 2021.
- Neil, R. and Zanger-Tishler, M. Algorithmic bias in criminal
risk assessment: The consequences of racial differences
in arrest as a measure of crime. *Annual Review of Crimi-
nology*, 8(1):97–119, 2025.
- Novelli, C., Taddeo, M., and Floridi, L. Accountability in
artificial intelligence: What it is and how it works. *AI &
Society*, 39(4):1871–1882, 2024.
- Office of Management and Budget. Advancing governance,
innovation, and risk management for agency use of arti-
ficial intelligence. Technical Report OMB Memorandum
M-24-10, Executive Office of the President, 2024.
- Park, S., Ahn, E., Ahn, T.-H., Ahn, S., Park, S., Kwon, E.,
Ahn, S., and Yang, Y. Artificial intelligence and aging
in place: A scoping review of current applications and
future directions. *The Gerontologist*, 65(6), 2025.
- Raji, I. D., Smart, A., White, R. N., Mitchell, M., Gebru, T.,
Hutchinson, B., Smith-Loud, J., Theron, D., and Barnes,
P. Closing the AI accountability gap: Defining an end-to-
end framework for internal algorithmic auditing. Techni-
cal report, arXiv, 2020.

- 550 Redden, J., Dencik, L., and Warne, H. Datafied child wel-
 551 fare services: Unpacking politics, economics and power.
 552 *Policy Studies*, 41(5):507–526, 2020.
- 553 Roberts, D. E. Child protection as surveillance of African
 554 American families. *Journal of Social Welfare and Family*
 555 *Law*, 36(4):426–437, 2014.
- 557 Robles, P. and Mallinson, D. J. Advancing AI governance
 558 with a unified theoretical framework: A systematic review.
 559 *Perspectives on Public Management and Governance*,
 560 2025.
- 562 Rooney, R. H. and Mirick, R. (eds.). *Strategies for Work*
 563 *with Involuntary Clients*. Columbia University Press,
 564 third edition, 2018.
- 565 Staub-Bernasconi, S. Soziale Arbeit: Dienstleistung oder
 566 Menschenrechtsprofession? In Lob-Hüdepohl, A. and
 567 Lesch, W. (eds.), *Ethik Sozialer Arbeit – Ein Handbuch*,
 568 pp. 20–54. UTB/Schöningh, 2007.
- 570 Staub-Bernasconi, S. Social work as a discipline and pro-
 571 fession. In Leskošek, V. (ed.), *Theories and Methods of*
 572 *Social Work: Exploring Different Perspectives*, pp. 9–30.
 573 Faculty of Social Work, University of Ljubljana, 2009.
- 574 Staub-Bernasconi, S. Social work and human rights—
 575 linking two traditions of human rights in social work.
 576 *Journal of Human Rights and Social Work*, 1(1):40–49,
 577 2016.
- 579 Taylor, B. J. Risk-managing decision-making: A psycho-
 580 social rationality model. *The British Journal of Social*
 581 *Work*, 51(7):2819–2838, 2021.
- 583 Van Amsterdam, W. A. C., Van Geloven, N., Krijthe, J. H.,
 584 Ranganath, R., and Cinà, G. When accurate prediction
 585 models yield harmful self-fulfilling prophecies. *Patterns*,
 586 6(4):101229, 2025.
- 587 Varon, J. and Peña, P. Artificial intelligence and consent: A
 588 feminist anti-colonial critique. *Internet Policy Review*, 10
 589 (4), 2021.
- 591 Wang, X., Oussalah, M., Niemilä, M., Ristikari, T., and
 592 Virtanen, P. Towards AI-governance in psychosocial care:
 593 A systematic literature review analysis. *Journal of Open*
 594 *Innovation: Technology, Market, and Complexity*, 9(4):
 595 100157, 2023.
- 597 Whittaker, A. and Havard, T. Defensive practice as ‘fear-
 598 based’ practice: Social work’s open secret? *British*
 599 *Journal of Social Work*, 46(5):1158–1174, 2016.
- 600 Wirtz, B. W., Weyerer, J. C., and Sturm, B. J. The dark
 601 sides of artificial intelligence: An integrated AI gover-
 602 nance framework for public administration. *International*
 603 *Journal of Public Administration*, 43(9):818–829, 2020.
- 604 Wirtz, B. W., Weyerer, J. C., and Kehl, I. Governance of
 artificial intelligence: A risk and guideline-based integra-
 tive framework. *Government Information Quarterly*, 39
 (4):101685, 2022.