

# Interpretable AI Techniques for Enhancing Transparency in Algorithmic Decision-Making for Public Welfare Services

KINGSLEY WISDOM AKHIBI

*Abstract- The increasing utilization of artificial intelligence in public welfare services has transformed how eligibility is determined, resources are allocated, and benefits are distributed by governments in ways that promise gains along the dimensions of efficiency, consistency, and scalability. Yet, the deployment of algorithmic decision-making in high-stakes welfare contexts has also generated significant challenges concerning transparency, accountability, fairness, and citizen trust, particularly when opaque or so-called "black-box" models are deployed. This journal critically investigates the role of interpretable AI in meeting these challenges in public welfare systems. Through the use of interdisciplinary literature related to machine learning interpretability and human-centered explanation design, administrative law, and public-sector governance, this research examines how superior operational efficiency can be balanced with democratic accountability through transparent and explainable AI systems. It discusses inherently interpretable model architectures and post-hoc explanation techniques, as well as the design of human-readable explanations that fit different welfare recipients, while assessing a number of methods for evaluating explanations regarding their quality, bias, and even fairness. Situating interpretability within broader governance and legal frames, the paper further underlines its importance for appeal rights, institutional accountability, and procedural justice. Synthesizing technical, social, and regulatory perspectives, the current research demonstrates that interpretability needs to be embedded by design in welfare AI rather than considered as an auxiliary feature. It, therefore, concludes that interpretable, human-centered AI is a necessary component in ensuring that algorithmic welfare decisions will not only be accurate and efficient but also socially legitimate, fair, and trustworthy.*

*Keywords: Artificial Intelligence in Public Welfare; Algorithmic Decision-Making; Interpretable Artificial Intelligence; Explainable AI (XAI); Human-Centered Explanations; Public Sector AI Governance; Welfare Eligibility and Allocation Systems; Algorithmic Fairness and Bias; Citizen Trust and Procedural Justice; Ethical AI in Government; Administrative Law and Automated Decision Systems*

## I. INTRODUCTION

Already, the deployment of artificial intelligence in public welfare services is changing the manner by which governments make eligibility decisions, allocation, or the distribution of certain benefits. This is because manually, eligibility decisions were made by caseworkers or other administrative officers using their discretion or certain established criteria. Today, with the sheer volume of data involved in social welfare systems, there is a drive for the application of algorithms for eligibility systems, for example, or for other social services aimed at improving efficiency, consistency, or speed (OECD, 2021). These systems utilize data from the past or predictions through machine learning algorithms for eligibility for vital services such as social benefits, health care, or housing and unemployment benefits.

While algorithmic systems promise efficiency gains, they pose substantial challenges with respect to transparency, accountability, and trust with citizens. This is because many "AI models are black boxes that generate answers without reason and especially if complex machine learning models are used" (Burrell, 2016). When it comes to public welfare, the implications of such black boxes are critical since "AI systems make decisions about people who need support, but those decisions are incomprehensible to the people concerned and often to the people responsible for running the system too" (O'Neil, 2016). Such challenges with respect to transparency significantly impair accountability, trust, and the appeals process (Zerilli et al., 2019).

This journal explores the usefulness of interpretable AI approaches in improving the transparency of algorithmic decision-making processes for public welfare services. This study is primarily concerned with techniques that enable human understanding of explanations for citizens who are impacted by the

automated eligibility or allocation decisions made by algorithms. This research synthesizes literature on topics such as model interpretability, post-hoc attribution methods, human-centered design principles, or governance frameworks to fill the gap between efficiency and accountability.

The research is driven by three major concerns. Firstly, it is important to ensure that AI systems are technically interpretable to achieve explainable AI systems that are traceable, accountable, or defensible in the context of current administrative or legal systems. Secondly, it is important to provide human-centered explanations that enable people who are impacted by AI systems to understand the rationale behind their implementation. Lastly, current governance regulations stipulate that AI systems must ensure accountability, ethics, or equity in their implementation, making it imperative for these systems to support AI-related transparency.

This journal helps to understand the relationship between the use of explainable AI, citizen-centered design of explanations, and governance, as it presents a broad perspective of how public welfare bodies should make use of AI in a responsible manner. Ultimately, the objective is to build AI systems that are optimized for efficiency but also for social equity, accountability, and trust, ensuring that algorithms are both more robust and more legitimate.

## II. ALGORITHMIC DECISION-MAKING IN PUBLIC WELFARE SERVICES

Public welfare systems are increasingly using decision-making algorithms in order to screen eligibility, prioritize beneficiaries, or allocate scarce public resources. Such systems are increasingly being used in areas such as public benefits management, housing allocation, public health services, unemployment benefits, or tax or subsidy benefits. Various governments opt for such systems with the ultimate aim of creating efficiency in administrative work, expediting processing, dealing with big data, or achieving consistency in decision-making (OECD, 2021).

Usually, the role of algorithmic decision-making in public welfare is played by some form of predictive or classification models that are trained using data from

previous administrative practices. These models make predictions using factors such as income levels, employment, family structure, medical history, or geo-location variables, among other factors, to produce automated or semi-automated decisions. In many places, such decisions either directly decide or significantly influence the decisions of social workers (UK Parliament, 2024). This means that decision-making algorithms are actually driving the decision-making processes in areas that directly influence the welfare of the citizenry.

It should, however, be noted that in more traditional forms of administrative decision-making, the rationale for application of certain rules or criteria is explicitly formalized. This is not necessarily the same with AI systems, where models such as deep learning operate in opaque manners that veil the logical framework that governs the relationship between input data and the final outcome. Indeed, such systems pose challenges in public welfare institutions that require certain standards of fairness, reasonableness, and accountability (Burrell, 2016).

International policy-making bodies have recognized that the application of algorithmic decision-making in the welfare sector is more than just an innovative application of technology, it is a shift in the governance of public services themselves. The OECD (2024) states that “when AI systems influence public benefits, they need to implement principles of administrative law, public governance, and the rights of citizens.” It is no longer a question of applying AI in public welfare services but of ensuring that it is applied in such a manner as to retain legitimacy.

## III. TRANSPARENCY DEFICITS & CHALLENGES OF PUBLIC TRUST

Despite their technical benefits, black-box decision-making systems often pose significant challenges related to transparency in public-welfare institutions. This is observed when users of these services, including public administrators, do not understand either the “how” or “why” of the decision that was made. This creates a public governance issue related to the explainability of public decisions (Zerilli et al., 2019).

Deficits of transparency appear in the following forms. Firstly, many welfare algorithms are either proprietary or developed by third-party contractors, which means that their system documentation is beyond public reach. Secondly, even if it is documented, it is mostly in a technical language that cannot be comprehended by people who are not involved in the processing (AlgorithmWatch, 2020). Lastly, advanced machine learning models can make accurate predictions without necessarily being able to provide a justification for the predictions, commonly referred to as the “black box problem” (Guidotti et al., 2018).

These issues of lack of transparency have substantial effects on trust for the public. Decisions on welfare have inherent sensitivities since they relate to resource availability that is critical for matters of dignity, survival, and socialization. If the public gets wrong decisions like disqualification for welfare or reduced welfare without valid explanations, they end up suspecting the system of being arbitrary or discriminatory (O’Neil, 2016). Research reveals that a lack of explanation influences distrust of the institutions involved, even with systems that work well in statistics, like automated systems (Barredo Arrieta et al., 2020).

Moreover, the lack of explainability makes it difficult for accountability to be established. This is since administrative law systems stipulate that people must be able to contest decisions and should be able to understand the rationale behind the decisions taken (Kaminski, 2019). This is especially true in social welfare decision-making, where the involved group could be less resourced or less knowledgeable about technology. Realizing these challenges, institutions of governance increasingly support that it is imperative for welfare algorithms to be designed with the integration of transparency, as an afterthought aspect. Transparency is clearly related to accountability, fairness, or legitimacy, as mentioned by the OECD (2021). Therefore, requiring techniques for the interpretation of AI is no longer a technical issue but is instead about rebuilding public trust for governance through AI powered welfare algorithms.

Below is a chart to visualize the relative contribution of key transparency deficits in public welfare AI

systems, illustrating that opacity is not purely technical but institutional and communicative.

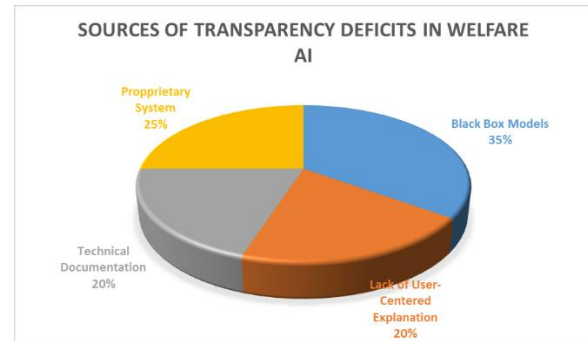


Figure1: Sources of Transparency Deficits in Welfare AI

The pie chart highlights black-box model opacity as the dominant source of transparency failure, followed closely by proprietary system constraints and inaccessible technical documentation. Notably, the significant share attributed to the lack of human-centered explanations supports the argument that transparency deficits are as much a design and governance problem as a computational one.

#### IV. CONCEPTUAL FOUNDATIONS OF INTERPRETABILITY IN ARTIFICIAL INTELLIGENCE

Interpretability, in artificial intelligence, designates a human's ability to understand the internal mechanics or decision logics of an AI system. In application domains such as public welfare, interpretability is not a matter of optional design preference but a binding requirement connected to legal accountability, ethical governance, and citizen rights. The concept itself, however, often stands in for or with cognates such as explainability, transparency, and intelligibility and, as such, is in dire need of conceptual clarification.

Interpretability is generally understood as a model's intrinsic comprehensibility—that is, the relationship between inputs and outputs is directly inspectable and interpretable by humans (Lipton, 2016). In contrast, explainability usually describes post-hoc methods of explaining how models work without revealing the model's inner structure. Transparency is a more general concept in governance; it refers to access to information, documentation, and decision-making processes, while intelligibility highlights that

explanations have to be meaningful for those they are given to (Molnar, 2023).

Scholars argue that interpretability needs to be context-sensitive. Doshi-Velez and Kim (2017) stress that what counts as an "interpretable" system brings in who the explanation is for and what decisions it informs. Within public welfare services, interpretability should satisfy several stakeholders at once: citizens wanting understandable reasons, caseworkers needing actionable insights, auditors checking on compliance, and policymakers safeguarding fairness. This multistakeholder requirement makes welfare AI different from commercial applications where explanations may be required for developers or regulators alone.

A key debate in this literature indeed revolves around whether inherently interpretable models should be favored over complex black-box models with post-hoc explanations. Rudin (2019) argues that in high-stakes domains, like welfare, health care, and criminal justice, reliance on post-hoc explanations is inadequate and potentially misleading. Models must be constructed to be interpretable. This reasoning fits very strongly with public welfare contexts within which decisions affect both fundamental rights and important social protections.

Interpretable accountability is also a form of mechanism for institutional accountability from a governance perspective. According to Zerilli et al. (2019), explanations allow decision-makers to justify outcomes, observe errors or bias, and correct systemic failures. Thus, interpretability is not only a matter of transparency at the level of single decisions; rather, it is about enabling continuous oversight and improvement of algorithmic systems embedded in public administration.

#### V. EXPLANATION REQUIREMENTS, HUMAN-CENTRED, FOR WELFARE DECISIONS

While technical interpretability is of course not unimportant, it is not enough. Public welfare systems will need to be human-centered in their explanations—that is, devised around the cognitive, social, and informational needs of those who experience the consequences of algorithmic decisions. This derives

from the fact that beneficiaries of welfare typically represent diverse societies with a large range of levels of education, digital literacy, and vulnerability.

Explanations centered on human needs would focus on clarity, relevance, and actionability. According to Wachter et al., explanations should answer what most citizens care about: "Why did this decision happen to me, and what could change it?" This perspective shifts explanation design away from abstract model descriptions toward individualized counterfactual narratives specifying which factors influenced the decision and how different circumstances might lead to alternative outcomes.

Explanations therefore need to avoid using technical jargon or probabilistic abstractions. According to Barredo Arrieta et al. (2020), explanations need to be in plaintext; their basis should be familiar concepts and be related to the decision context. For instance, an explanation for welfare eligibility has to refer to concrete criteria like income thresholds or household size, rather than statistical weights or feature importance scores. This way, comprehension enhances and arbitrariness decreases.

Other human-centered requirements involve considerations about accessibility and fairness. Explanations should be understandable by people with a disability, available in multiple languages if need be, and sensitive to socio-economic contexts. The European Data Protection Supervisor 2020 goes on to illustrate that meaningful explanation is integral in protecting fundamental rights, especially within those populations that are marginalized and likely to suffer disproportionately at the hands of automated welfare systems.

Lastly, explanations should underpin procedural justice. Research has demonstrated that citizens are more likely to accept negative decisions if the process leading up to the decision is perceived as procedurally fair and respectful, even if the decision itself is unfavorable (Binns, 2018). Transparent, human-readable explanations contribute towards this perception by communicating that decisions are reasoned rather than arbitrary. In welfare contexts, this is essential for procedural legitimacy, wherein trust in public institutions is sustained for long-term algorithmic governance.

## VI. INTERPRETABLE MODEL ARCHITECTURES FOR WELFARE DECISION SYSTEMS

Interpretable model architectures are those whose internal logic can be directly understood without the need for external explanation tools. In public welfare decision-making, such models are particularly valuable in allowing administrators, auditors, and citizens to trace how specific inputs lead to specific outcomes. The literature increasingly supports the use of inherently interpretable models in high-stakes domains where decisions affect fundamental rights and access to essential services.

Common interpretable architectures include rule-based systems, decision trees, scoring models, and generalized linear models. Rule-based systems, which encode eligibility criteria as explicit if-then rules, are very similar to traditional administrative decision-making processes. It is thereby easy to align with legal compliance and also easy to explain since the decisions can be justified by referring to clearly defined criteria (Molnar, 2023). Similarly, decision trees provide a transparent decision path; users can trace a sequence of conditional splits that culminate in an outcome.

Another major use for welfare contexts, which also relates to their simplicity and interpretability, is linear and logistic regression models. The coefficients in those models directly show the direction and magnitude of influence that individual variables have on the decision outcome. According to Doshi-Velez and Kim, such models may not always yield predictive performance comparable to more complex alternatives, but they do provide a level of transparency often more appropriate for public-sector uses.

Recent work questions this presumed trade-off of accuracy for interpretability. For example, Rudin shows that for many structured decision-making tasks, such as welfare eligibility and risk assessments, interpretable models yield performance competitive with black-box methods. This negates the rationale for using opaque models in the welfare context, where the social cost of unexplained decisions is very high.

Importantly, interpretable architectures support institutional accountability beyond individual explanations. Since the entire decision logic is transparent, policymakers and oversight bodies can assess whether the model is aligned with policy objectives, legal standards, and equity goals. Such systemic transparency thus allows for proactive detection of biases and policy refinement, reinforcing again how interpretability ought to be recognized as a device of governance rather than as a purely technical feature of the models.

## VII. POST-HOC EXPLANATION TECHNIQUES FOR COMPLEX ALLOCATION AND ELIGIBILITY MODELS

Despite the many advantages of inherently interpretable models, public welfare agencies often use complex machine learning systems in order to deal with high-dimensional data or to capture relationships that are nonlinear. In such scenarios, explanation techniques that are post-hoc need to be used to approximate the rationale associated with model outputs. Post-hoc explanation techniques aimed at providing insights into black-box models do not affect their internal structure (Guidotti et al, 2018).

One such popular post-hoc technique is feature attribution, which explains the specific contribution any variable provides for a particular prediction. Algorithms like LIME and SHAP produce local interpretability by explaining which factors most drove a given decision or not. In welfare contexts, feature attribution helps caseworkers and citizens understand why certain attributes, such as income level or employment history, were decisive. Other post-hoc approaches are surrogate models. These models approximate the behavior of complex systems by simpler, interpretable models like decision trees or linear regressions. While intuitive explanations can be provided using surrogate models, the reliability of such explanations depends on the faithfulness of the surrogate model with regards to the original model's behavior. This is particularly critical in public welfare contexts where possibly wrong explanations can mislead citizens or affect appeal procedures. Counterfactual explanations have recently emerged as a human-centered complement to technical model explanations. Instead of detailing how a model works,

counterfactuals describe what minimal perturbations would have resulted in a different outcome. Thus, a welfare applicant might receive the explanation that eligibility would have been approved if income were slightly lower, or household size larger. This approach is fundamentally closely aligned with citizens' informational needs and promotes actionable understanding. Yet, researchers advise against relying too much on post-hoc explanations when it comes to high-stake public decisions. According to Rudin (2019), post-hoc methods could yield plausible but incorrect explanations, giving the impression of transparency without actual accountability. Post-hoc techniques are best used in moderation and, whenever possible, coupled with interpretable model design and strong governance safeguards.

The chart below compares the relative adoption of explanation techniques used in welfare AI systems, showing how institutions balance interpretability and performance.

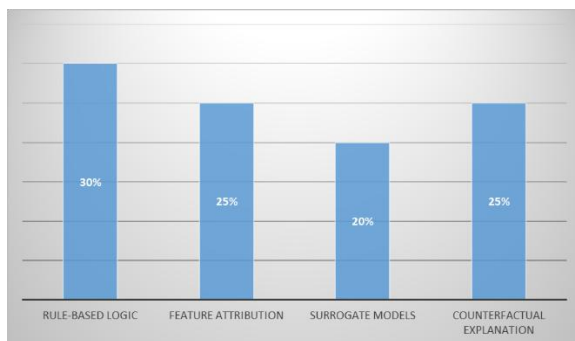


Figure 2: Relative Use of Explanation Techniques in Welfare AI

The bar chart illustrates that while rule-based logic remains prominent, post-hoc methods such as feature attribution and counterfactual explanations are widely used to compensate for complex models. The relatively high use of counterfactual explanations aligns the emphasis on human-centered explanation design, as these techniques directly answer citizens' core questions about eligibility and changeability. However, the chart also implicitly supports the caution that reliance on post-hoc explanations must be moderated to avoid creating an illusion of accountability.

## VIII. DESIGNING HUMAN-READABLE EXPLANATIONS FOR AFFECTED

Designing human-readable explanations for algorithmic decisions in public welfare services is a complex task that involves repackaging technical reasoning into a form that is comprehensible to ordinary citizens. Human-readable explanations are distinct from technical explanations in that their key application is not for technical inspection but for communication with laypeople. This is important in welfare, where the welfare of vulnerable sections of the community hinges on the outcome of such decisions.

It is important to provide effective explanations for humans using the principles of narrative clarity, context relevance, and procedural clarity. Molnar (2023) states that the explanation of a prediction should be limited to the factors that influence that prediction, as it would not be helpful to specify everything that influences that prediction. For example, instead of pointing out the various input features, the explanation should point out the important features that influence the prediction, such as exceeding income or lack of contribution history.

A more human-readable design applies best in the context of counterfactual explanations. Using the counterfactual explanation approach is helpful to people since it reveals what would happen if an individual or group of individuals had achieved a certain outcome differently (Wachter et al. 2017). A welfare context could be used by identifying what would happen if income levels were low or if there were an update of reports. It is important to note that such explanations not only clarify but also empower.

Visual and structured explanation aids are also beneficial for readability. According to Barredo Arrieta et al., it is essential to provide explanations through basic tables, decision trees, or bullet-point explanations that establish a direct relationship between the criteria and the outcomes (2020). Such methods are important in digital welfare systems, which should offer readable explanations. Moreover, it is crucial that the explanation does not employ probabilistic or statistic language to provide clarity on decision-making processes.

From a governance point of view, explanations that are human-readable add to procedural justice and improve the legitimacy of institutions. People believe that the decision taken is just, even if it is not in their favor, if they are treated with respect and given good explanations. Organizations involved in public welfare should ensure that their explanations are human-readable so that citizens continue to trust the institutions (Binns, 2018).

#### IX. ASSESSMENT OF EXPLANATION QUALITY IN PUBLIC WELFARE SITUATIONS

It is important to assess the quality of explanations of public welfare AI systems for ensuring that the objectives related to transparency are not superficially dealt with. It is not possible to evaluate the quality of explanations for AI systems using technical criteria only, but it should be carried out by a socio-technical framework that takes into account factors such as accuracy, utility, fairness, as well as public perception (Doshi-Velez & Kim, 2017).

One important aspect of explanation assessment is fidelity: the closer the explanation is to mirroring the behavior of the system, the better the explanation is. This is important since high-fidelity explanations ensure that the explanations given to the public are true indicators of the decision-making process, as it helps to avoid misled or incomplete explanations (Guidotti et al., 2018). In welfare institutions, low-fidelity explanations could hamper the appeals system with potential lawsuits.

Understandability is another important metric here. An explanation could be technically accurate but still not understandable to the target group of users. Molnar (2023) argues that the quality of explanation should be judged with respect to users' capability to understand and apply the information given correctly. This is particularly important in public welfare institutions, where testing of explanations among various users is often a consideration.

Perceptions of fairness and trustworthiness are also significant evaluative criteria. Research clearly shows that explanations affect not only understanding but also perceptions of decision legitimacy (Zerilli et al., 2019). A clear articulation of consistent criteria and

expressions of uncertainty can add to perceptions of trustworthiness, while vagueness or terminological language can contribute to declining trust.

Lastly, evaluation should look into the issue of institutional usability. Explanations should facilitate work for caseworkers, auditors, and authorities responsible for decision correction. AlgorithmWatch (2020) points out that it is essential for explanation evaluation to be incorporated into the processes of evaluation of impact. This way, explanation quality would become a permanent governance issue, no longer just a design solution.

#### X. BIAS, FAIRNESS, AND SOCIAL EQUITY IN INTERPRETABLE WELFARE ALGORITHMS

Issues of bias and fairness are key considerations in applying algorithms for decision-making in public welfare services. This is because welfare algorithms are often developed using historical administrative data, such as data that is influenced by existing social inequalities or by unequal distribution of public services. If such data is applied without critical scrutiny, then the issue is that the algorithms developed could end up entrenching structural inequalities despite their objectivity (O'Neil, 2016).

Interpretable AI is crucial in pointing out fairness concerns and dealing with them. Interpretability helps in understanding decision-making processes, thereby allowing policymakers to understand the influence of sensitive variables like income proxies, geographic factors, or living arrangements on policy outcomes. According to Selbst et al. (2019), abstraction in AI systems tends to remove social context, causing fairness solutions to be technically sound but ineffective socially. This is where the role of interpretation is crucial in re-establishing a link between technology and social implications.

Feature-level explanations are especially valuable in social welfare applications. Consider, for example, a machine learning system that gives considerable emphasis to data at the neighborhood level. This could indirectly introduce racial or socio-economic biases even if protected features are not explicitly modeled. Explainable models or techniques can help reveal hidden features, allowing institutions to re-evaluate

their application in relation to equity principles (Guidotti et al., 2018).

On more normative ground, fair welfare systems require more than just statistical parity; they also need procedural and distributive justice. Writing on fairness in welfare systems, Binns (2018) highlights that fairness is value-laden and grounded in political theory. Interpretability helps with fairness by allowing societal debate about what should count in welfare systems.

Finally, interpretability is a sort of safety net that prevents invisible harm. This is because interpreter AI is responsible for allowing for the discussion of inequitable trends that can then be fixed. This is important since without this type of interpretation, the fix for inequalities may end up being superficial.

#### XI. GOVERNANCE, ACCOUNTABILITY, AND LEGAL CONSIDERATIONS OF EXPLAINABLE WELFARE AI

AI applications for public welfare services are limited by certain legal frameworks that prioritize accountability, reason-giving, and the right to contest administrative decisions of the government. Explainable or interpretable AI models are important for ensuring that public law principles are adhered to.

Today, governance bodies understand that it is important for algorithmic systems to be auditable and explainable from start to finish. This is related to accountability by the OECD (2021), who asserted that public institutions must explain their decisions not only within their own institutions but among the people. Explainable AI helps reconcile this need by allowing decisions to be traced from data to results.

Scholars of law discuss the relevance of explanation in ensuring the promotion of human rights. Kaminski (2019) explains that despite the differences of the “right to explanation” among various nations, administrative justice still demands that people must be able to comprehend the rationale of administrative decisions that influence their lives. In the welfare state context, opacity in AI systems hinders appeal rights as citizens are not able to discern any disparity or mistake.

However, accountability is also related to institutional responsibility. When decisions become automated, there is a possibility of diffused accountability, where accountability or responsibility is transferred from human institutions to automated systems or their providers or vendors. It is possible to resist or mitigate this challenge of accountability through interpretability, ensuring that human decision-makers become accountable for the automated decisions (Zerilli et al, 2019),.

Guidelines on governance underscore the importance of proactive measures for safeguarding. This is cited by the European Commission’s Ethics Guidelines for Trustworthy AI (2020) or European Data Protection Supervisor (2020), which state that it is important for explainability to be “integrated into system design, not bolted on as an afterthought.” Within public welfare services, it is important for accountability systems to be functioning from the start of their implementation to minimize any legal concerns, while public trust is maintained.

#### XII. CHALLENGES OF INSTITUTIONAL AND TECHNICAL INTEGRATION

The integration of interpretable AI systems into public welfare services is faced with various institutional as well as technical challenges that go beyond the development of the models. This is because public welfare institutions are often characterized by complex organizational structures that have information systems, data infrastructure, and limited administrative capacities. These factors can hamper the adoption of the interpretable AI systems.

One of the major issues that need to be addressed is data quality and availability. Usually, welfare data is incomplete, inconsistent, or gathered for purposes that are not analytical. Molnar (2023) agrees that the input variables need to be of significant interpretation for the explanation to make sense. If the data is poorly defined or noisy, Molnar asserts that the explanation could be fallacious or too trivial for understanding. Instances where the data is public or gathered in welfare settings, the presence of past policy shifts adds to the issue of data interpretation. This is attributed to data. (Molnar, 2023)



Institutional capability is also a hindrance. Many public institutions do not have staff with the capability to develop, test, or ensure the explainability of AI systems. This means that they end up depending on external companies for such systems, which could hamper the issue of explainability if the models or systems are assumed to be intellectual property (AlgorithmWatch, 2020).

Additionally, the culture of an organization influences the degree of integration of caseworkers and administrators. These individuals can lack confidence in algorithmic systems if there is a lack of alignment between explanations given by the systems and their own professional judgment or practices. According to OECD (2024), for there to be effective integration of AI in an organization, training must be carried out to enable public servants to critically engage with AI systems. This can be achieved by using interpretable systems that enable public servants to understand the decision-making process of AI systems.

Finally, the issue of scalability or performance needs to be taken into consideration. Though more interpretable models can be easier to audit or understand, they could pose some challenges in dealing with big data. This is because there is a need to strike a good balance between interpretability and other performance considerations without undermining the objectives of ensuring interpretation (Doshi-Velez & Kim, 2017).

### XIII. DESIGN PRINCIPLES FOR TRANSPARENT & TRUSTWORTHY WELFARE AI SYSTEMS

To build trustworthy and transparent AI systems for public welfare purposes, it is important to make interpretability and explainability an integral part of the AI system life cycle, right from development through to deployment and evaluation. Instead, the key is to provide transparency by design, as is emphasized by best practices (OECD, 2021).

An important design tenet is the preference for necessarily interpretable models, whenever possible. As suggested by Rudin (2019), high-risk social welfare issues require models whose reasoning is verifiable for examination and justification by direct scrutiny. However, if necessarily complex models

must be used, then they should be supplemented by high-quality explanation systems that are themselves independently validatable.

Another is stakeholder-centered design. An explanation should be designed with the needs of the stakeholders in mind, such as citizens, social work staff, auditors, and policy-makers. Barredo Arrieta et al. (2020) advocate for multi-layered explanation structures with citizen-friendly summary levels for the former group, accompanied by more technical levels for the latter group.

Documentation and auditability are also critical. This is because documentation of the assumptions made by the models, the data used, and the criteria adopted for decision-making can ensure monitoring and accountability. This is emphasized by the European Data Protection Supervisor (2020), who considers it important to make use of explanation logs and audit trails to ensure that there is compliance with the law after deployment.

Lastly, open assessment of the AI systems increases trustworthiness. Open welfare AI systems should be regularly evaluated for their outcomes, audits for any bias, and user feed-backs. AlgorithmWatch (2020) emphasizes that openness is an iterative process that demands institutional support over time. Once the principles are followed, public welfare bodies can implement welfare AI systems that provide efficiency through effective AI while maintaining democracy.

### XIV. FUTURE DIRECTIONS FOR INTERPRETABLE AI IN PUBLIC WELFARE DECISION-MAKING

Future of Explainable AI: It is being shaped by new technological developments, regulatory concerns, as well as demands for citizen engagement in governance. A number of solutions are being developed to improve the explainability of AI systems without significantly diminishing their efficiency. One of the areas of research is the development of hybrid models that leverage the use of purely interpretable models with the careful application of complex techniques. Such models are aimed at identifying complex patterns in big data while still maintaining interpretability for high-importance decision-making areas (Rudin, 2019). This is achieved through the

application of explainable models that work in tandem with predictive models for optimized efficiency without losing accountability. Second, participatory design of explanation for explanation effectiveness is important. Currently, research suggests that citizens, along with other caseworkers, should be encouraged to participate in the creation of explanation forms and measures of evaluation (Barredo Arrieta et al., 2020). This ensures that not only is the explanation technically sound, but it is also significant for citizens who are influenced by the outputs of an algorithm. Such participatory design techniques include workshops, surveys, and online dashboards. Improved methodologies for evaluation are also important. Currently, the literature suggests that there is a need for standardized evaluation criteria that evaluate both technical soundness and human understanding of explanations (Doshi-Velez & Kim, 2017). Future work is likely to explore the inclusion of cognitive as well as social aspects into evaluation methodologies, such as fairness perception, trust, or behavior effect of explanations. This interdisciplinary methodology is expected to ensure that welfare AI systems adhere to both legal specifications and social demands. Lastly, the emergence of regulations and governance models is likely to influence the course of interpretable AI applications for welfare services. This is with respect to regulations such as the European AI Act, advice from OECD, or other regulatory authorities that stipulate more accountability, fairness, and clarity in critical applications of AI (OECD, 2024; European Data Protection Supervisor, 2020). Implementation of such regulations would need proactive governance models such as independent audit trails.

## XV. CONCLUSION

The welfare AI of the future must be found where technological development meets humanist solutions, as well as effective governance. This can be achieved through the advancement of understandable models, participatory explanation approaches, and evaluative frameworks that ensure public welfare institutions make a positive contribution in terms of resource distribution in an increasingly algorithmic world.

## REFERENCES

- [1] AlgorithmWatch. (2020). Automating Society Report 2020. <https://algorithmwatch.org/en/automating-society/>
- [2] Barredo Arrieta, A., et al. (2020). Explainable Artificial Intelligence (XAI): Concepts, taxonomies, opportunities and challenges toward responsible AI. *Information Fusion*, 58, 82–115. <https://arxiv.org/abs/1910.10045>
- [3] Binns, R. (2018). Fairness in machine learning: Lessons from political philosophy. *Proceedings of FAT*. <https://arxiv.org/abs/1712.03586>
- [4] Burrell, J. (2016). How the machine ‘thinks’: Understanding opacity in machine learning algorithms. *Big Data & Society*. <https://journals.sagepub.com/>
- [5] Doshi-Velez, F., & Kim, B. (2017). Towards a rigorous science of interpretable machine learning. *arXiv*. <https://arxiv.org/abs/1702.08608>
- [6] European Commission. (2020). Ethics Guidelines for Trustworthy AI. <https://digital-strategy.ec.europa.eu/>
- [7] European Data Protection Supervisor. (2020). *Assessing the impact of artificial intelligence on fundamental rights*. <https://edps.europa.eu/>
- [8] Guidotti, R., et al. (2018). A survey of methods for explaining black box models. *ACM Computing Surveys*, 51(5). <https://arxiv.org/abs/1802.01933>
- [9] Holzinger, A., et al. (2019). What do we need to build explainable AI systems for the medical domain? *arXiv*. <https://arxiv.org/abs/1712.09923>
- [10] Kaminski, M. E. (2019). The right to explanation, explained. *Berkeley Technology Law Journal*. <https://papers.ssrn.com/>
- [11] Kim, B., et al. (2018). Interpretability beyond feature attribution: Quantitative testing. *ICML*. <https://arxiv.org/abs/1811.06166>
- [12] Lipton, Z. C. (2016). The mythos of model interpretability. *ICML Workshop*. <https://arxiv.org/abs/1606.03490>
- [13] London, A. J. (2019). Artificial intelligence and black-box medical decisions. *Hastings Center Report*. <https://onlinelibrary.wiley.com/>

- [14] Molnar, C. (2023). *Interpretable Machine Learning: A Guide for Making Black Box Models Explainable*.  
<https://christophm.github.io/interpretable-ml-book/>
- [15] OECD. (2021). *Algorithmic transparency and accountability in the public sector*.  
<https://www.oecd.org/>
- [16] OECD. (2024). *Governing with Artificial Intelligence: Are Governments Ready?*  
<https://www.oecd.org/>
- [17] O'Neil, C. (2016). *Weapons of math destruction*. Crown Publishing.
- [18] Ribeiro, M. T., Singh, S., & Guestrin, C. (2016). "Why should I trust you?" Explaining predictions of any classifier. *KDD*.  
<https://arxiv.org/abs/1602.04938>
- [19] Rudin, C. (2019). Stop explaining black box machine learning models. *Nature Machine Intelligence*. <https://www.nature.com/>
- [20] Selbst, A. D., et al. (2019). Fairness and abstraction in sociotechnical systems. *FAT\**.  
<https://arxiv.org/abs/1810.08810>
- [21] UK Parliament. (2024). *Automated decision-making in the public sector*.  
<https://commonslibrary.parliament.uk/>
- [22] Wachter, S., Mittelstadt, B., & Russell, C. (2017). Counterfactual explanations without opening the black box. *Harvard Journal of Law & Technology*. <https://papers.ssrn.com/>
- [23] World Economic Forum. (2022). *Global AI Governance Toolkit*. <https://www.weforum.org/>
- [24] Zerilli, J., et al. (2019). Transparency in algorithmic and human decision-making. *Philosophy & Technology*.  
<https://link.springer.com/>