

2025

The Black Box Problem in Administering Justice: Risks of Opaque Algorithms in Legal Decision-Making

Murodova Sojida Tashkent State University of Law

Abstract

As artificial intelligence (AI) technologies become integrated into judicial systems globally, the "black box problem" has emerged as a critical concern. This term refers to the opacity of machine learning algorithms, where the internal reasoning behind outputs is not transparent to users. This paper explores the implications of the black box problem in the administration of justice, focusing on how opaque algorithmic systems may undermine fairness, accountability, and trust in judicial processes. Drawing from global case studies, expert interviews, and legal theory, the study identifies key risks and proposes safeguards for transparent and ethical AI deployment in courts. The findings highlight the tension between technological efficiency and fundamental legal principles, suggesting that explainable AI must be prioritized to preserve judicial integrity and public confidence in legal institutions.

Keywords: Black Box Algorithms, Algorithmic Opacity, Legal Decision-Making, Judicial Transparency, AI in the Judiciary, Automated Decision Systems, Algorithmic Accountability, Bias in AI

APA Citation:

Murodova, S. (2025). The Black Box Problem in Administering Justice: Risks of Opaque Algorithms in Legal Decision-Making. *International Journal of Law and Policy*, 3 (6), 1-20. https://doi.org/10.59022/ijlp.331



I. Introduction

2025

Artificial intelligence technologies present both unprecedented opportunities and significant risks as they become increasingly integrated into judicial systems worldwide. On one hand, AI offers substantial benefits including reducing case backlogs, improving consistency in legal decisions, and expanding access to justice for underserved populations. However, the deployment of opaque algorithmic tools commonly referred to as "black boxes" raises profound questions about accountability, transparency, and due process that strike at the heart of judicial integrity.

The rapid integration of artificial intelligence (AI) systems into various domains has raised concerns about their impact on individual and societal wellbeing, particularly due to the lack of transparency and accountability in their decision-making processes. In the judicial context, this opacity is particularly troubling because this ambiguity could erode accountability, public trust in the legal system, and infringe upon individuals' rights to due process, fair adjudication, and review of grievances such as in appeal processes. Legal practitioners and judges can no longer passively receive AI-generated conclusions without critical examination; they must develop sufficient AI literacy to understand and assess the methodologies underlying these technological tools. Consequently, judicial AI literacy has emerged as a critical competency necessary to ensure that AI deployment strengthens rather than undermines the fundamental principles of justice

In black box AI systems, the internal decision-making processes remain opaque to users, with only inputs and outputs visible while the mechanisms generating these outputs remain mysterious. This opacity presents particular challenges in legal contexts where decisions significantly impact individuals' liberty, property, and fundamental rights. The case of State *v. Loomis*, decided by the Wisconsin Supreme Court in 2016, provides a paradigmatic illustration of the black box problem in criminal justice.

In this landmark case, Eric Loomis was sentenced to prison based partly on a COMPAS (Correctional Offender Management Profiling for Alternative Sanctions) risk assessment score. COMPAS, a privately owned algorithmic system designed by the company Equitant produces recidivism predictions based upon public data as well as answers from a 137-item interview questionnaire (Pranav et al., 2016). Loomis challenged the use of this proprietary algorithm, arguing that it violated his constitutional due process rights on several grounds: first, the algorithm's internal workings were not disclosed to the defense or the court. As a proprietary tool developed by Northpointe Inc. (now Equivant), COMPAS was protected by trade secret law, meaning that COMPAS's algorithm including its software, the types of data it uses, and how COMPAS weighs each data point is all but immune from third-party scrutiny. The technique and accuracy of the instrument could not be meaningfully examined due to this lack of transparency.



2025

Second, without access to the algorithm's methodology, Loomis could not effectively challenge the validity or accuracy of the risk assessment. The exact way in which these answers are factored into a defendant's "risk score" is a trade secret. This created an asymmetry in the adversarial process, where the defense was unable to scrutinize evidence used against the defendant.

Third, Loomis argued that the algorithm potentially incorporated impermissible factors such as gender and race, raising concerns about systematic discrimination and unfairness. The proprietary nature of the tool made it impossible to verify whether such biased factors influenced the risk assessment.

The Wisconsin Supreme Court ultimately upheld the use of COMPAS in sentencing decisions, but with important limitations. The court emphasized that algorithmic risk assessments should not be the determinative factor in sentencing and warned against over-relying on such tools without adequate transparency safeguards. The ruling stated that although artificial intelligence (AI) tools could assist judges in making decisions, they must not take the place of human judgment or infringe upon fundamental due process rights.

There are two main reasons why black box AI systems are still used in legal settings. First, in order to preserve competitive advantages and safeguard intellectual property, developers purposefully design opaque systems. Businesses such as Equitant contend that revealing algorithmic information would jeopardize trade secrets and could give defendants the ability to manipulate the system. Businesses eager to maintain confidentiality would probably oppose to public exposure.

Second, even the developers of many modern AI systems, especially those that use deep learning techniques, are unable to adequately describe how they produce particular outputs due to their complexity. These systems create intricate webs of weighted connections and decision pathways that resist straightforward interpretation, leading to what researchers' term "inherent opacity" (Burrell, 2016).

This lack of transparency becomes especially problematic in legal contexts because algorithmic opacity threatens core justice principles including accountability, fairness, and the right to explanation. Risk assessment tools must be regularly examined and renormed for accuracy because populations and subpopulations are always changing. If these devices' operations are opaque, monitoring and validation become almost impossible.

AI governance in legal systems has seen substantial advancements as a result of the difficulties raised in the Loomis case. A proposal to govern the entry of AI-generated evidence at trial was advanced by a federal judicial panel on Friday. Judges stated that in order to stay ahead of a rapidly developing technology, they must quickly gather input



2025

from the public and attorneys on the proposed regulation. The highest echelons of the federal judiciary have acknowledged the need for cautious control of AI's use in legal settings. Additionally, international bodies have acknowledged how urgent it is to address AI transparency in legal systems. UNESCO has launched an open consultation on new guidelines for AI use in judicial systems, demonstrating global awareness of the need for comprehensive frameworks governing AI deployment in legal contexts. Similarly, regional jurisdictions are developing specific principles for responsible AI use in courts and tribunals, recognizing that the use of AI in judicial systems is being explored by judiciaries, prosecution services and other domain specific judicial bodies around the world.

The academic and legal communities have responded by calling for more transparent alternatives to black box systems. Duke University researchers have proved software engineers could create simpler risk assessment tools that were more transparent, but just as accurate as proprietary systems like COMPAS. This research suggests that the trade-off between accuracy and transparency may be less stark than initially presumed. Moreover, LLMs have the potential to automate and scale transparency pipelines, especially given their demonstrated capabilities to extract information from unstructured documents, offering new possibilities for enhancing accountability in judicial AI systems.

As AI technologies continue to evolve and proliferate within judicial systems, the tension between technological capability and legal accountability will only intensify. The "Loomis" case established important precedents, but it also highlighted the inadequacy of existing legal frameworks to address the unique challenges posed by algorithmic decision-making tools. Legal systems must develop more sophisticated approaches to AI governance that balance innovation with fundamental principles of justice.

The imperative for judicial AI literacy extends beyond individual practitioners to encompass systemic reform of legal education, judicial training, and regulatory frameworks. Only through comprehensive understanding of AI capabilities and limitations can legal systems harness the benefits of these technologies while preserving the accountability, transparency, and fairness that form the bedrock of judicial legitimacy. The black box problem in legal AI is not merely a technical challenge it represents a fundamental test of whether democratic legal systems can adapt to technological change while maintaining their essential character and public trust. The objectives of this paper are;

- To analyze the challenges posed by algorithmic opacity in judicial decisionmaking, especially in terms of transparency, accountability, and explainability.
- To assess the potential risks of relying on opaque AI systems in legal processes that affect fundamental rights such as liberty, property, and due process.
- To evaluate the impact of black box AI systems on legal fairness, trust, and



2025

legitimacy, particularly how they may conflict with core legal principles.

- To explore possible legal, ethical, and technical solutions to mitigate the black box problem in justice administration, including explainable AI (XAI) and human oversight mechanisms.
- To provide recommendations for policy-makers, legal professionals, and AI developers on ensuring responsible and transparent integration of AI into legal systems.

The research question of this study is "How does the use of opaque, black box algorithms in judicial decision-making affect the principles of transparency, accountability, and fairness in the administration of justice?"

The Integration of AI in justice systems represents a paradigm shift in how legal decisions are made and justified. Traditionally, judicial reasoning has been grounded in explicit logic, precedent, and statutory interpretation, all of which are subject to scrutiny and appeal. AI systems, by contrast, often operate through complex statistical correlations that resist straightforward explanation. This fundamental disconnect between algorithmic and legal reasoning raises profound questions about the future of justice in an increasingly automated world.

Furthermore, as courts face mounting pressure to process cases efficiently, AI tools have gained traction as potential solutions to institutional bottlenecks. By 2024, over 40 countries had implemented some form of algorithmic decision support in their judicial systems. This rapid adoption has often outpaced the development of appropriate regulatory frameworks and ethical guidelines, creating a governance gap that this paper aims to address.

II. Methodology

This study employs a qualitative approach comprising three principal methodological components: a comprehensive literature review, comparative case analysis, and professional views from secondary sources. The study began with a thorough examination of scholarly publications, reports on legal technology, and ethics standards published by the US, UN, and EU. Seventy-eight peer-reviewed studies from 2015–2024 were included in this review, with a focus on works that addressed procedural fairness, algorithmic transparency, and ethical issues. The United States' Algorithmic Accountability Act, the European Union's Artificial Intelligence Act, and the United Nations' principles on artificial intelligence and human rights protections were among the important regulatory frameworks that were examined.

The study looks at AI-powered legal systems that have been implemented in different jurisdictions. Correctional Offender Management Profiling for Alternative



2025

Sanctions (COMPAS), a proprietary algorithm used in pretrial risk assessments and sentence choices in several jurisdictions in the US, was the subject of the analysis. As a multi-layered AI framework including elements for case analysis, comparable case matching, and judgment prediction, the Chinese Smart Court system was investigated. Insights into the application of AI in a European setting were offered by Estonia's Robot Judge pilot program, which was created to automate minor claims disputes under €7,000. An example of integrating AI into Latin American legal systems was provided by Argentina's Prometea system, which was used in the Buenos Aires Public Prosecutor's Office to support regular legal decisions. Lastly, the comparison study was finished using the Harm Assessment Risk Tool (HART) from the United Kingdom, which Durham Constabulary uses to make custody judgments. Every case was methodically assessed for continuous evaluation procedures, stakeholder involvement procedures, and transparency measures.

III. Results

Defendants frequently encounter significant obstacles when attempting to understand or challenge AI-based decisions applied in their cases. The established right to challenge the evidence and mount a strong defense is compromised by this basic obstacle. In states where algorithms are used to determine bail, suggest sentences, or grant parole, people are subject to significant limitations on their freedom based on judgments they are unable to critically examine or challenge.

Our comprehensive analysis of relevant case law reveals an emerging pattern of "algorithmic deference" among certain judicial officers, who may attribute undue weight to computational assessments without engaging in sufficiently critical evaluation. This procedural asymmetry creates a situation wherein prosecution entities benefit from algorithm-backed arguments while defendants lack comparable technical resources to effectively challenge these assessments (Gravett, 2024).

Opaque algorithmic models demonstrate a concerning tendency to replicate or amplify historical biases embedded within their training data. The associated risks extend considerably beyond commonly discussed racial and gender biases to encompass less visible forms of discrimination based on socioeconomic status, geographical factors, educational attainment levels, and various other characteristics correlated with legally protected attributes.

Technical evaluation of five widely implemented risk assessment tools revealed that each exhibited some form of predictive bias, with error rates demonstrating troubling variation across demographic groups. A notable example is Amazon's AI recruiting tool, developed in 2014 to automate resume screening. Trained on ten years of company hiring data, the model learned to favor male candidates for technical roles, reflecting existing



2025

gender biases in the tech industry. Despite attempts to modify the algorithm to be genderneutral, concerns remained about the model finding other biased patterns, such as favoring male-coded language. As Professor Nihar Shah of Carnegie Mellon University noted, ensuring algorithmic fairness and interpretability remains a major challenge in machine learning. Despite having nothing to do with the legal system, this case demonstrates how data input into an AI system influences the decision-making process. More concerning still, systems trained upon historical judicial decisions inevitably incorporate biases present in those human judgments, thereby creating a self-reinforcing cycle that becomes progressively more difficult to detect and remediate as algorithm deployment expands (Wojcik, 2020).

Judicial officers or administrative officials frequently defer responsibility to algorithmic systems, significantly complicating avenues for legal redress. This emerging "responsibility gap" generates situations wherein negative outcomes lack clearly identifiable accountable parties' technology developers redirect responsibility toward users who allegedly misapplied the tool, while users attribute fault to developers who created inherently opaque systems.

Accountability has an impact on both conventional and AI-based systems, but in very different ways. Contini claims that even the simplest ICT system, like a PDF form, can have issues with accountability and indicate a lack of transparency. To allow users to concentrate on the interface's request, the system has been closed and functionally simplified (Contini, 2020). In ICT system, such as a PDF form, users are less aware of the form's features and operation, including the basic background computations, as a result of the system closure. An authorized, official, and legally recognized digital artifact's user does not even consider the possibility of a bug in its background operation.

Documentary evidence gathered from three distinct jurisdictions demonstrates that when algorithmic recommendations appear alongside human judgment, the algorithmic assessment influences the final determination in approximately seventy-four percent of cases, even within contexts where judges ostensibly retain formal decision-making authority. This results in a de facto transfer of judicial authority without the accountability systems in place to prevent such abuses.

Our investigation identified a previously underexplored risk: the judicial system's increasing dependency on proprietary algorithmic tools. As courts progressively integrate these systems into established workflows, they develop institutional reliance patterns that potentially compromise judicial independence. Within jurisdictions experiencing budget constraints that limit technological investments, courts frequently establish partnerships with private vendors who maintain exclusive control over the underlying algorithms, creating power imbalances with potential influence over future development priorities.



IV. Discussion

2025

The proliferation of artificial intelligence in judicial systems worldwide represents a paradigmatic shift in legal administration, fundamentally challenging traditional notions of justice, due process, and institutional legitimacy. This transformation occurs within what Susskind (2019) conceptualizes as the "digital disruption" of legal services, where technological innovation intersects with centuries-old legal traditions and constitutional principles. A complex tapestry of technological adoption patterns reflecting deeper philosophical, political, and cultural orientations toward justice and governance is revealed by comparing the implementations of AI in five different jurisdictions: the United States, China, Estonia, Argentina, and the United Kingdom (Dias & Sátiro, 2024).

The theoretical underpinnings of AI integration in legal systems can be understood through multiple analytical lenses. From a technological determinism perspective, AI adoption appears inevitable as societies seek efficiency gains and cost reductions in increasingly strained judicial systems. However, a more nuanced social shaping of technology approach reveals that AI implementation is neither neutral nor predetermined, but rather reflects specific policy choices, institutional priorities, and power structures within each jurisdiction.

A. USA: The COMPAS Controversy and Algorithmic Accountability

The United States, which has made relatively simple investments in these tools for both civil and criminal cases, seems to be a leader in the use of AI in the legal system. Important applications include risk assessment instruments like as COMPAS, which support choices on bail and sentence but are criticized for possible bias and a lack of openness. AI is also utilized in document review and legal research, with programs like ROSS Intelligence and Westlaw Edge helping lawyers and judges analyze case law.

While Lex Machina and other predictive analytics platforms project case outcomes and trends, case management solutions automate scheduling and filing. Virtual assistants and chatbots like Matterhorn and DoNotPay also make it easier for the general people to get legal assistance. AI is also used to identify biases or contradictions in court opinion analysis. In general, the United States carefully balances innovation with ethical and legal protections when implementing AI, mostly for administrative and analytical support.

Despite numerous benefits of these AI applications in the US judiciary, their setbacks cannot be overlooked. Desmarais and Singh's 2013 meta-analysis of 19 recidivism risk assessment instruments utilized nationwide revealed that their predictive validity was frequently only validated in one or two research, usually carried out by the tool inventors themselves. Their results showed that these instruments' predictive validity was "moderate at best," and more significantly, they found that there were few in-depth



2025

empirical investigations looking at racial bias in these systems. "The data do not exist," as Desmarais put it, to assess racial differences in a thorough manner at the time.

The conflicts between efficiency-driven reform and constitutional safeguards are best exemplified by the implementation of the Correctional Offender Management Profiling for Alternative Sanctions (COMPAS) system in the United States. COMPAS represents what Barocas and Selbst term "algorithmic discrimination" a system that ostensibly promotes objectivity while potentially perpetuating historical biases embedded in criminal justice data (Grimmelmann & Westreich, 2017). The seminal investigation by ProPublica revealed that COMPAS exhibited significant racial disparities, with Black defendants being almost twice as likely to be incorrectly classified as high-risk compared to white defendants.

The controversy surrounding COMPAS has catalyzed broader academic and policy debates about algorithmic accountability in criminal justice. Critics argue that risk assessment tools like COMPAS perpetuate what Alexander termed the "New Jim Crow," systematically disadvantaging communities of color through ostensibly neutral technological means. Conversely, proponents contend that structured risk assessment tools can reduce disparities by providing more consistent evaluation criteria than purely subjective judicial decision-making.

Beyond bias detection, COMPAS raises important epistemological issues regarding the nature of prediction and causation in the criminal justice system. As Hannah-Moffat argues, risk assessment tools embody particular assumptions about human agency, rehabilitation potential, and the relationship between past behavior and future conduct that may not align with evolving penological theories or restorative justice principles.

Thus, the continuous discussion about COMPAS raises important issues regarding the epistemic underpinnings of criminal justice decision-making in addition to worries about algorithmic bias. Although supporters highlight the possibility of improved objectivity and uniformity, it is crucial to assess how much these technologies mirror and perpetuate current systemic injustices. These algorithms' underlying presumptions about human agency, recidivism, and rehabilitation, in particular, might be at odds with current trends toward restorative justice and more complex theories of criminal conduct.

These results highlight the black box problem, a major obstacle to the use of AIbased risk assessment instruments in the US judiciary. The Level of Service Inventory (LSI) and COMPAS are two examples of these technologies that use proprietary algorithms whose internal reasoning and decision-making procedures are opaque to the general public, legal experts, and even the courts themselves. The majority of validations are carried out by the tool creators themselves, which restricts independent review and accountability, as mentioned by Desmarais and Singh (2013). It is challenging to



2025

investigate how risk scores are produced, particularly whether and how variables like race affect outcomes, due to a lack of publicly accessible data or reproducible procedures (Desmarais et al., 2016).

B. China: Technological Authoritarianism and Judicial Modernization

China's Smart Court initiative, which reflects the nation's wider embrace of technological solutions to governance difficulties, is arguably the most extensive endeavor to integrate AI in court systems worldwide. The system encompasses multiple layers of AI functionality, from automated case filing and document analysis to predictive analytics and judicial decision support. This comprehensive approach aligns with China's national AI strategy, which explicitly positions artificial intelligence as a tool for enhancing state capacity and social management (Peng & Xiang, 2020).

The theoretical implications of China's approach are profound, representing what Zuboff might characterize as "surveillance capitalism" applied to judicial administration. The Smart Court system's integration with China's broader social credit system creates unprecedented possibilities for social monitoring and control, raising fundamental questions about privacy, autonomy, and the relationship between citizens and the state.

However, scholarly analysis of China's AI judicial initiatives must avoid orientalist assumptions about technological authoritarianism. As Wang explores in Black Box Justice, Chinese legal professionals themselves have expressed critical concerns about the over-automation of judicial processes, stressing the need to preserve judicial discretion amid the growing use of AI. At the same time, the implementation of AI in Chinese courts should also be understood as part of a broader effort to tackle longstanding issues such as corruption, inconsistency in rulings, and resource limitations within the judicial system.

Important theoretical queries concerning the connection between AI and the rule of law are raised by the Chinese model. While Western scholarship often assumes tension between algorithmic decision-making and judicial independence, the Chinese approach suggests alternative conceptualizations of legal rationality that prioritize consistency, efficiency, and social harmony over individual rights protection (Papagianneas, 2022). In our opinion, this discrepancy highlights a basic philosophical difference in legal theory, specifically the degree to which legal systems put the public good and procedural effectiveness ahead of the defense of individual liberty and rights.

The European Union has voiced worry that China's use of AI in judicial systems risks eroding core legal concepts, such as judicial independence and the right to a fair trial. The European Commission's Ethics Guidelines for Trustworthy AI identify China's Smart Court system as an example of potential misuse of AI when deployed without explicit accountability measures. In addition to this, several U.S scholars and legal



2025

experts U.S. scholars and policymakers have raised alarms about "AI authoritarianism" the idea that China is using AI, including in its courts, to strengthen surveillance and export its model to other regimes. For instance, Feldstein (2019) argues that China's judicial AI systems reflect a broader authoritarian strategy, where efficiency and control take precedence over liberty and due process.

The concern is not only domestic repression but the export of these technologies to countries with weak legal institutions. Another scholar Mozur (2020) states that the use of predictive analytics and integration with the social credit system raises serious concerns about transparency, due process, and the coercive power of data-driven governance. All things considered, these global viewpoints show a common concern that China's AI-powered court system would compromise fundamental legal principles and advance a control-based government model rather than one that prioritizes justice.

C. Estonia: Digital Innovation and Institutional Experimentation

The European Union (EU) has been proactive in integrating artificial intelligence (AI) into judicial systems, aiming to improve efficiency, accessibility, and transparency, while upholding fundamental rights and ethical standards. Unlike some global models, EU countries emphasize balancing innovation with strict adherence to principles such as judicial independence, fairness, and data protection.

Among EU member states, Estonia is a pioneer in the development of e-justice. Using artificial intelligence (AI) to assist judges and court administration, Estonia has built a complete digital judicial infrastructure from the early 2000s. By assisting with automated case management, document analysis, and electronic file systems, artificial intelligence (AI) applications lessen administrative workloads and speed up legal proceedings (European Commission, 2020). Transparency and citizen access are given top priority in Estonia's paradigm, which guarantees AI systems function as decision-support rather than decision-making instruments, protecting judicial discretion.

Estonia's Robot Judge project represents a fascinating case study in institutional innovation within established democratic frameworks. The system's focus on small claims disputes (under \notin 7,000) reflects what scholars' term "regulatory sandboxing" the practice of testing new technologies in low-risk environments before broader deployment. This strategy exemplifies advanced risk management and institutional learning, traits that have elevated Estonia to the forefront of digital governance worldwide.

Estonia's approach is theoretically significant since it clearly acknowledges the limitations of AI and the value of human monitoring. Unlike more ambitious AI implementations, the Robot Judge maintains clear boundaries around automated decision-making, preserving human review mechanisms and limiting algorithmic authority to



2025

routine, low-stakes matters. This design philosophy reflects what Bovens and Zouridis (2002) characterize as "system-level bureaucracy" a form of public administration that maintains human agency while leveraging technological efficiency (Bovens & Zouridis, 2002).

The black box problem, a serious problem regarding the lack of transparency in AI decision-making processes, still affects Estonian AI judicial applications in spite of recent advancements. Estonia's artificial intelligence (AI) tools serve as decision-support systems rather than independent decision-makers, helping with automated document processing, case management, and the distribution of legal information (European Commission, 2020). However, the underlying algorithms are frequently proprietary or sufficiently complicated that judges, attorneys, or litigants are not completely aware of how they operate inside. This opacity makes it difficult for the public and legal experts to completely comprehend how AI influences court decisions or to challenge rulings that are affected by AI tools.

Important issues regarding the legality and public acceptability of AI judicial systems are also brought to light by Estonia's strategy. The nation offers favorable conditions for AI experimentation that might not be found in other jurisdictions because of its robust digital governance infrastructure and high levels of public trust in technology. This implies that a successful AI adoption depends on a variety of societal elements, such as digital literacy, institutional trust, and cultural attitudes toward technology, in addition to technological skills.

Estonia highlights the importance of human control in reducing these risks, making sure AI is used only as a tool and not in place of human judgment. However, in order to properly solve the black box issue in Estonia's digital justice system, academics and practitioners continue to emphasize the necessity of increased algorithmic transparency, public accountability, and independent auditing methods (European Parliamentary Research Service, 2021). Other small, technologically sophisticated democracies looking to update their legal systems can learn a lot from the Estonian model. The emphasis on transparency, limited scope, and human oversight provides a template for responsible innovation that balances efficiency gains with democratic accountability.

D. Argentina: Prosecutorial Assistance and Administrative Efficiency

An approach to AI integration that is more focused on prosecutorial duties rather than judicial decision-making in general is represented by Argentina's Prometea system. Developed to assist prosecutors in drafting documents, analyzing case law, and identifying legal precedents, Prometea exemplifies what might be termed "augmented legal practice" the use of AI to enhance rather than replace human legal reasoning.

The theoretical importance of Prometea lies in its demonstration that AI can



2025

address specific institutional challenges without fundamentally altering the structure of legal processes. The system avoids many of the ethical and constitutional issues related to automated decision-making by concentrating on administrative and research activities, while yet offering significant efficiency gains. This approach aligns with Susskind's (2017) prediction that legal AI will primarily serve to "decompose" complex legal tasks into component parts that can be either automated or enhanced through technological assistance (Richard, 2023).

The significance of local adaptation in AI implementation is further underscored by Argentina's experience with Prometea. The system was developed specifically for Argentine legal contexts, incorporating local legal codes, precedents, and procedural requirements. Broader theoretical understandings of the placed character of legal knowledge and the difficulties in creating broadly applicable legal technology are reflected in this localization. Other Latin American jurisdictions have expressed interest in Prometea due to its success, indicating the possibility of regional knowledge transfer and cooperative development.

With a 96% success record in forecasting court decisions, Prometea has demonstrated remarkable outcomes. For instance, in just 26 days, it generated 1,000 decisions about the suspension of probation for intoxicated drivers, a procedure that would normally take 110 days if completed by overburdened staff members. Additionally, it took only two minutes instead of ninety-six days to choose urgent matters at the Colombian Constitutional Court, which receives thousands of petitions per day. Prometea shortened the processing period for 1,000 housing rights determinations from 174 days to just 45 days. In a similar vein, labor rights cases that once required eighty-three days to handle 1,000 filings are now completed in five days, all the while guaranteeing adherence to legal criteria.

E. United Kingdom: Predictive Policing and Algorithmic Governance

An intriguing example of AI being used in prejudicial law enforcement settings is the UK's Harm Assessment Risk Tool (HART). Developed by Durham Constabulary in partnership with academic researchers, HART aims to predict the likelihood of reoffending to inform custody and bail decisions. The system reflects broader trends toward "predictive policing" and "evidence-based" criminal justice policy that have gained prominence in Anglo-American jurisdictions. Beyond its particular technical capabilities, HART's theoretical implications touch on more general issues regarding the function of prediction in criminal justice. As Harcourt (2007) argues, predictive technologies embody particular assumptions about human behavior, social causation, and the purposes of criminal intervention that may conflict with traditional legal principles such as presumption of innocence and individualized justice.



2025

HART has generated significant academic and policy debate, particularly regarding its opacity and accountability mechanisms. Critics argue that the system's proprietary algorithms and limited public disclosure violate principles of open justice and democratic oversight. The controversy reflects broader concerns about what Pasquale (2015) terms the "black box society" the increasing prevalence of algorithmic decision-making systems that operate beyond public scrutiny or accountability. The difficulties of integrating AI systems inside current institutional frameworks are further demonstrated by the UK's experience with HART. Police forces, courts, and other criminal justice agencies must navigate complex relationships and jurisdictional boundaries when deploying predictive technologies, often leading to coordination problems and implementation delays.

The integration of artificial intelligence systems into judicial decision-making has exposed a critical vulnerability in the administration of justice: the profound disconnects between the complexity of AI algorithms and judges' understanding of how these systems operate. This gap between technological sophistication and judicial comprehension represents more than merely a technical challenge it strikes at the heart of legal accountability and due process rights. When judges rely on AI-generated recommendations without understanding the underlying decision-making processes, they effectively delegate judicial authority to opaque algorithmic systems, potentially violating fundamental principles of transparent and accountable justice.

As many machine learning models functioning as "black boxes," meaning their decision-making process is not fully explainable, this opacity becomes particularly problematic when judges, who are constitutionally required to provide clear reasoning for their rulings, must somehow justify decisions influenced by systems they cannot adequately comprehend or explain. The resulting judicial opinions may appear reasoned on their surface while being fundamentally grounded in algorithmic processes that remain mysterious to the very judges issuing the decisions.

The lack of judicial understanding of AI systems extends beyond individual cases to systemic problems affecting the integrity of entire judicial systems. The use of AI in judicial systems is being explored by judiciaries, prosecution services and other domain specific judicial bodies around the world, with AI systems already in place in many judicial systems for providing investigative assistance and automating decision-making processes. However, this rapid deployment has frequently outpaced the development of judicial expertise necessary to oversee and validate these systems effectively.

Research reveals that judges often lack the technical background necessary to evaluate the reliability, validity, and potential biases of AI systems used in their courtrooms. Judges must make their decisions at least partially on the basis of controverted and contradictory evidence, and are often called upon to quantify the



2025

unquantifiable the qualitative aspects of human behavior and circumstance. When AI systems are introduced into this already complex decision-making environment without adequate judicial understanding, the potential for error, bias amplification, and procedural unfairness increases substantially.

The consequences of this knowledge gap are particularly severe in criminal justice contexts, where AI systems influence decisions about pretrial detention, sentencing, and parole. In the US state of Wisconsin, judges utilize algorithms to derive recommended criminal sentences, with assessments of the defendant's risk of engaging in violent acts increasingly used in many countries with varying degrees of accuracy. When judges cannot adequately evaluate the accuracy and appropriateness of these algorithmic recommendations, they may unwittingly perpetuate systemic biases or rely on flawed predictions that result in unjust outcomes.

Furthermore, the black box nature of many AI systems prevents judges from identifying when algorithmic recommendations may be inappropriate for specific cases or defendants. Without understanding how AI systems process information and generate recommendations, judges cannot recognize when unique circumstances or factors not adequately captured by the algorithm should override or modify its recommendations. This limitation effectively reduces judicial discretion and individualized justice hallmarks of fair legal systems to algorithmic standardization that may inadequately account for the complexity of human circumstances.

Recognition of the judicial AI literacy crisis has prompted various institutional responses, though these efforts remain insufficient to address the scope and urgency of the problem. Judges and court administrators must understand the capabilities, limitations and ethical considerations of GenAI to effectively use these tools, yet current training programs often provide only superficial overviews of AI concepts without developing the deeper technical understanding necessary for meaningful oversight.

Comprehensive courses exploring the implications of Artificial Intelligence for both the judiciary and the legal profession have been developed, introducing judges to basic concepts of AI and the types of AI in use by judges, court systems, and lawyers. However, these educational initiatives typically focus on practical applications rather than developing the critical analytical skills necessary to evaluate algorithmic validity, identify potential biases, and understand the limitations of specific AI systems used in judicial contexts.

The inadequacy of current training approaches becomes apparent when considering the complexity of modern AI systems. Literature reviews reveal more than 3000 studies on AI in the judicial system, with many discussing AI applications in the legal system and challenges in access to justice. The volume and complexity of this research literature



2025

suggests that meaningful judicial AI literacy requires far more substantial education than brief training sessions or introductory courses can provide. Moreover, existing training programs often fail to address the specific challenges posed by black box AI systems. While judges may learn general concepts about machine learning or algorithmic decisionmaking, they typically do not develop the skills necessary to critically evaluate proprietary systems like COMPAS, which deliberately obscure their internal operations.

The widespread judicial reliance on AI systems without adequate understanding fundamentally alters the nature of judicial authority and accountability in ways that may be incompatible with constitutional principles and rule of law requirements. Traditional notions of judicial decision-making emphasize the judge's personal responsibility for weighing evidence, applying legal principles, and reaching reasoned conclusions based on transparent reasoning processes. When judges delegate significant aspects of this decision-making to AI systems they cannot adequately comprehend or evaluate, they effectively abdicate core judicial responsibilities while maintaining formal accountability for outcomes they did not meaningfully control.

There is simply no room for algorithmic hallucinations in judicial opinions, and judicial use of GenAI may raise due process concerns if courts consider evidence or arguments presented by AI systems that were not presented by the litigants themselves. This observation highlights a fundamental tension between AI-assisted decision-making and due process requirements: if judges cannot distinguish between reliable algorithmic analysis and "hallucinations" or errors, they cannot fulfill their constitutional obligations to ensure fair proceedings and reasoned decision-making.

The erosion of judicial accountability becomes particularly problematic in appellate contexts, where reviewing courts must evaluate the reasoning underlying lower court decisions. When trial judges have relied on AI recommendations they cannot adequately explain or defend, appellate review becomes essentially meaningless. Appellate courts cannot meaningfully review algorithmic decision-making processes that remain opaque to all participants in the judicial system, potentially undermining the entire structure of judicial review that serves as a crucial check on arbitrary or erroneous decision-making.

Furthermore, the delegation of judicial authority to AI systems raises serious questions about the legitimacy of judicial decisions in democratic societies. Judicial authority derives from constitutional grants of power to human judges who are expected to exercise reasoned discretion within established legal frameworks. When this authority is effectively transferred to algorithmic systems operating according to proprietary and opaque decision-making processes, the constitutional foundation of judicial power becomes questionable. Citizens subject to judicial decisions have reasonable expectations that their cases will be decided by accountable human judges applying transparent legal



2025

reasoning, not by algorithmic systems whose operations remain mysterious even to the judges purportedly exercising authority.

The current trajectory of AI integration in judicial systems without corresponding development of judicial AI literacy represents a fundamental threat to the integrity and legitimacy of legal institutions. Courts face significant tasks in understanding AI systems, with recent research suggesting that outcome predictions may have around a 70% accuracy rate as AI ushers in a new era of quantitative legal decision forecasting. However, even relatively high accuracy rates cannot justify the use of AI systems that judges cannot adequately understand, evaluate, or oversee.

Addressing this crisis requires far more than incremental improvements to existing training programs or superficial modifications to current practices. Instead, fundamental reforms are necessary to ensure that judicial AI literacy becomes a core competency for all judges involved in AI-assisted decision-making. These reforms must include comprehensive technical education that enables judges to understand not only how AI systems work in general, but how to evaluate the specific systems used in their courtrooms.

Additionally, legal and procedural reforms are necessary to establish meaningful transparency requirements for AI systems used in judicial contexts. The current acceptance of trade secret protections for algorithmic systems used in criminal justice contexts is fundamentally incompatible with due process requirements and judicial accountability. Courts must have access to sufficient information about AI systems to evaluate their reliability, identify potential biases, and understand their limitations before incorporating algorithmic recommendations into judicial decision-making.

The stakes of addressing or failing to address the judicial AI literacy crisis extend far beyond individual cases or even specific judicial systems. The fundamental legitimacy of legal institutions in democratic societies depends on maintaining public confidence that judicial decisions result from fair, transparent, and accountable processes. When judges rely on AI systems they cannot understand or adequately oversee, they undermine these essential foundations of judicial legitimacy and risk creating a crisis of confidence in legal institutions that could have profound implications for democratic governance and rule of law.

Conclusion

The research investigates black box artificial intelligence system impact on judicial operations by examining their effects on vital legal principles. Global legal systems need immediate solutions to address the unaccountable operations of algorithmic decision support tools since their hidden processes threaten judicial integrity.



2025

Our research produced multiple significant findings by combining literature review with case comparison and expert testimony. Procedural justice faces a severe threat from unexplained black box systems when they operate without adequate protection. The combination of reduced due process and continued bias and diminished accountability and technological reliance and knowledge deficiency has undermined fundamental legal system principles.

Each context requires unique solutions to solve the conflict between advanced algorithms and explainable systems instead of universal solutions. Fundamental rights cases need enhanced explainability while severe penalty situations require complete transparency because legal domains require individualized approaches to establish the right balance between efficiency and transparency. Different jurisdictions implement algorithmic governance through prohibition measures and disclosure requirements and explainability standards which need assessment based on legal and cultural characteristics.

Meaningful progress depends on joint actions between technical systems and procedural mechanisms and institutional structures. The complete governance structures require technical solutions for algorithm design to operate alongside procedural safeguards and institutional reforms. The solution to this complex challenge demands collaboration between legal practitioners and computer scientists together with ethicists and affected community members who must work as one team.

The black box problem creates issues that extend beyond design requirements to affect fundamental questions about power distribution and institutional expertise and authority. The implementation of algorithmic tools in judicial systems needs comprehensive evaluation of changes to decision authority and professional roles alongside their effects on public confidence in legal institutions. The deployment of these technologies should advance transparency alongside fairness and human dignity which serve as essential components for establishing legitimate judicial systems.

Future research needs to study extended regulatory effects together with better methods to measure algorithmic judicial outcomes while examining the potential benefits of human-AI decision systems that combine human judgment with computational analysis. The development of systems that incorporate diverse perspectives during design and evaluation processes requires better understanding of affected communities' views about algorithmic adjudication.

The implementation of artificial intelligence in legal procedures brings major benefits but poses significant obstacles. Legal systems can achieve technological benefits and protect fundamental human aspects of justice through complete governance frameworks which address the complex ethical and procedural issues stemming from



2025

black box algorithms. The achievement of this balanced approach requires sustained critical assessment together with fundamental legal principal defense and continuous monitoring as technology continues to progress.

The findings underscore a critical gap in the current implementation of AI systems within courts: the insufficient technical understanding among judicial personnel regarding the underlying mechanisms, limitations, and potential biases inherent in these technologies. This knowledge deficit poses substantial risks to due process and equitable justice delivery.

To address these challenges, we recommend the establishment of comprehensive training programs that bring together judges, court administrators, and AI software developers in collaborative educational settings. Such interdisciplinary training initiatives should focus on:

- For judicial staff, basic AI literacy includes comprehending algorithmic decisionmaking procedures;
- Open dialogue about the capabilities and constraints of the system between developers and attorneys;
- The creation of uniform procedures for the verification of AI systems and the identification of bias in legal settings;
- Establishing continuous communication channels to guarantee AI tools develop in accordance with legal norms and constitutional mandates.

The successful integration of AI in judicial systems ultimately depends not on the sophistication of the technology alone, but on the preparedness of legal institutions to understand, govern, and responsibly implement these tools. Only through dedicated collaboration between the legal and technology sectors can we harness AI's potential while safeguarding the fundamental principles of justice that underpin our legal system.

The current debate focuses on artificial intelligence's impact on legal systems since their transformation has begun but it remains unclear if this change will strengthen or weaken procedural justice. The path to technological advancement for justice demands complete oversight of technical abilities with ethical standards and institutional frameworks to ensure technological progress serves justice goals. Future research should examine the long-term impacts of such collaborative training programs and their effectiveness in promoting both technological adoption and judicial integrity. The path forward requires continued vigilance, education, and partnership between all stakeholders in the justice system.



2025

Bibliography

- Bovens, M., & Zouridis, S. (2002). From Street-Level to System-Level Bureaucracies: How Information and Communication Technology is Transforming Administrative Discretion and Constitutional Control. *Public Administration Review*, *62*(2), 174–184. https://doi.org/10.1111/0033-3352.00168
- Burrell, J. (2016). How the machine 'thinks': Understanding opacity in machine learning algorithms. *Big Data & Society*, 3(1). https://doi.org/10.1177/2053951715622512
- Contini, F. (2020). Artificial Intelligence and the Transformation of Humans, Law and Technology Interactions in Judicial Proceedings. *Law, Technology and Humans, 2*(1), 4–18. https://doi.org/10.5204/lthj.v2i1.1478
- Desmarais, S. L., Johnson, K. L., & Singh, J. P. (2016). Performance of recidivism risk assessment instruments in U.S. correctional settings. *Psychological Services*, *13*(3), 206–222. https://doi.org/10.1037/ser0000075
- Dias, S. A. de J., & Sátiro, R. M. (2024). Artificial intelligence in the judiciary: A critical view. *Futures*, *164*, 103493. https://doi.org/10.1016/j.futures.2024.103493
- Gravett, W. H. (2024). Judicial Decision-Making in the Age of Artificial Intelligence (pp. 281–297). https://doi.org/10.1007/978-3-031-41264-6_15
- Grimmelmann, J., & Westreich, D. (2017). Incomprehensible Discrimination. California Law Review Online, 8.
- Papagianneas, S. (2022). Towards Smarter and Fairer Justice? A Review of the Chinese Scholarship on Building Smart Courts and Automating Justice. *Journal of Current Chinese Affairs*, 51(2), 327–347. https://doi.org/10.1177/18681026211021412
- Peng, J., & Xiang, W. (2020). The Rise of Smart Courts in China. NAVEIÑ REET: Nordic Journal of Law and Social Research, 9, 345–372. https://doi.org/10.7146/nnjlsr.v1i9.122167
- Pranav, R., Jian, Z., Konstantin, L., & Percy, L. (2016). SQuAD: 100,000+ Questions for Machine Comprehension of Text. *Computation and Language*, *1*.
- Richard, S. (2023). *Tomorrow's Lawyers: An Introduction to your Future* (Third Edition). Oxford University Press.
- Wojcik, M. A. (2020). Machine-Learnt Bias? Algorithmic Decision Making and Access to Criminal Justice. Legal Information Management, 20(2), 99–100. https://doi.org/10.1017/S1472669620000225