



U.S. Department of
JUSTICE

Artificial Intelligence and Criminal Justice

Final Report



December 3, 2024

Table of Contents

Foreword.....	3
I. Introduction.....	4
Background	5
The Use of AI in the Criminal Justice System	9
II. Identification & Surveillance.....	12
Introduction	12
Biometrics	13
Focusing on FRT for Identification in Law Enforcement Investigations.....	17
Automated License Plate Recognition	26
III. Forensic Analysis.....	29
Introduction	29
Current Uses of AI in Forensic Analysis.....	31
Future Uses of AI in Forensic Analysis	33
Challenges for AI in Forensic Analysis	36
Recommendations	40
IV. Predictive Policing.....	42
Introduction	42
Uses of Predictive Policing	43
Risks of Predictive Policing	46
Recommendations	49
V. Risk Assessment	54
Uses of Risk Assessment.....	55
Risk Assessment Design	58
Potential Benefits for a More Effective, Equitable, and Efficient Criminal Justice System	60
The Risks of Risk Assessment	62
Recommendations	67
VI. Conclusion & Best Practices.....	70
Foundations for AI Governance	71
Pre-Deployment Measures	73
Post-Deployment Measures.....	75
Legal Disclaimer	77

Foreword

Dear Mr. President,

The Department of Justice (DOJ), in consultation with the Department of Homeland Security (DHS) and the White House Office of Science and Technology Policy (OSTP), submits this report in response to Executive Order 14110, *Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence* (EO 14110).

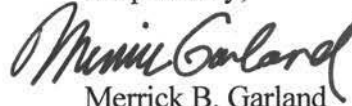
Pursuant to the requirements set forth in EO 14110, this report satisfies the EO's directive in Section 7.1(b) that the Attorney General:

submit to the President a report that addresses the use of AI in the criminal justice system, including any use in:

- (A) sentencing;
- (B) parole, supervised release, and probation;
- (C) bail, pretrial release, and pretrial detention;
- (D) risk assessments, including pretrial, earned time, and early release or transfer to home-confinement determinations;
- (E) police surveillance;
- (F) crime forecasting and predictive policing, including the ingestion of historical crime data into AI systems to predict high-density “hot spots”;
- (G) prison-management tools; and
- (H) forensic analysis.

In each of those areas of the criminal justice system, this report “identif[ies] areas where AI can enhance law enforcement efficiency and accuracy, consistent with protections for privacy, civil rights, and civil liberties” and “recommend[s] best practices for law enforcement agencies, including safeguards and appropriate use limits for AI” and addresses the concerns set forth in both EO 14110 and Executive Order 14074, *Advancing Effective, Accountable Policing and Criminal Justice Practices to Enhance Public Trust and Public Safety* (EO 14074).

In addition, chapter IV of this report—Predictive Policing—is submitted in fulfillment of Section 13(e) of Executive Order 14074. This report satisfies EO 14074's directive that “the Attorney General, the Secretary of Homeland Security, and the Director of OSTP . . . jointly lead an interagency process regarding the use by [Law Enforcement Agencies] of facial recognition technology, other technologies using biometric information, and predictive algorithms.”

Respectfully,

Merrick B. Garland
Attorney General

I. Introduction

Artificial intelligence (AI) use is rapidly transforming the criminal justice system and has the potential to make it more effective, equitable, and efficient. AI use also has the potential to cause harms, amplify disparities, and misdirect resources.

Today, AI in criminal justice predominantly involves conventional statistical analysis, such as regression models.¹ But that is changing. The accelerating pace of AI innovation is leading to increased use of computer vision, natural language processing, and other types of AI across the criminal justice system. Organizations and officials in the criminal justice system are also beginning to deploy generative AI systems.

The policy and technology choices that law enforcement agencies, pretrial and probation services, prison systems, and other criminal justice stakeholders make in the near term will affect millions of Americans. These choices will also set the trajectory for rapid expansion in the scope and scale of AI use throughout the criminal justice system.

On October 30, 2023, President Biden issued Executive Order 14110 on the *Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence* (EO 14110).² EO 14110 advances a coordinated governmentwide approach to responsible adoption of AI. It emphasizes that AI must advance equity and civil rights, respect privacy and civil liberties, and meet government performance objectives.

In order to “promote the equitable treatment of individuals and adhere to the Federal Government’s fundamental obligation to ensure fair and impartial justice for all,” section 7.1(b) of EO 14110 directs the Attorney General, in consultation with the Secretary of Homeland Security and Director of the Office of Science and Technology Policy (OSTP), to submit a report to the President on “the use of AI in the criminal justice system.” The EO enumerates types of AI uses in criminal justice to address, and it further directs that the report “identify areas where AI can enhance law enforcement efficiency and accuracy, consistent with protections for privacy, civil

¹ This report uses a broad definition of the term “artificial intelligence,” consistent with Section 238(g) of the John S. McCain National Defense Authorization Act for Fiscal Year 2019, Pub. L. No. 115-232, and the elaboration provided by OMB Memoranda M-24-10 and M-24-18 in implementing Executive Order 14110. The definition encompasses “[a]ny artificial system that performs tasks under varying and unpredictable circumstances without significant human oversight, or that can learn from experience and improve performance when exposed to data sets,” as well as “[a]n artificial system developed in computer software, physical hardware, or other context that solves tasks requiring human-like perception, cognition, planning, learning, communication, or physical action,” among other types of systems. Furthermore, “no system [is] too simple to qualify as covered AI due to a lack of technical complexity (e.g., the smaller number of parameters in a model, the type of model, or the amount of data used for training purposes).” MEM. FROM SHALANDA D. YOUNG, DIR., OFF. MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, TO HEADS OF EXEC. DEP’T & AGENCIES (Mar. 28, 2024), available at <https://www.whitehouse.gov/wp-content/uploads/2024/03/M-24-10-Advancing-Governance-Innovation-and-Risk-Management-for-Agency-Use-of-Artificial-Intelligence.pdf>; MEM. FROM SHALANDA D. YOUNG, DIR., OFF. MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, TO HEADS OF EXEC. DEP’T & AGENCIES (Sept 24, 2024), available at <https://www.whitehouse.gov/wp-content/uploads/2024/10/M-24-18-AI-Acquisition-Memorandum.pdf>.

² Exec. Order No. 14110, 88 Fed. Reg. 75191 (Oct. 2023), <https://www.federalregister.gov/d/2023-24283>.

rights, and civil liberties” and “recommend best practices for law enforcement agencies, including safeguards and appropriate use limits for AI.”

The Department of Justice submitted an Executive Report in conformity with these requirements on October 29, 2024. This final report addresses these issues in richer detail and provides particularized recommendations.

Background

This report benefits from substantial input from stakeholders inside and outside of the federal government. Throughout 2024, Deputy Attorney General Lisa Monaco convened Justice AI, a series of six roundtable conversations that brought together law enforcement agencies and groups, civil society organizations, companies that develop AI products and services, and academic researchers who study AI use in criminal justice. In addition, the Civil Rights Division hosted four quarterly information exchanges with federal, state, and local agencies addressing civil rights issues associated with AI. The National Institute of Justice received comments for this report through a public request for input, and Department of Justice staff met with stakeholders throughout the process of preparing this report. The Department is grateful to the many experts who shared their valuable perspectives.

This report builds on a decade of U.S. Government initiatives that aim to ensure that AI use is effective, transparent, and respectful of privacy and civil liberties.

- In May 2014, the White House released a report entitled *Big Data: Seizing Opportunities, Preserving Values*.³ The report touched on how algorithmic analysis of large datasets could have “tremendous” benefits for law enforcement activities, while also posing privacy and civil liberties issues.
- In May 2016, the White House issued the follow-up report, *Big Data: Algorithmic Systems, Opportunity, and Civil Rights*.⁴ In a section on criminal justice, the report noted that uses of algorithms could advance public safety and public trust, but they must be “designed and deployed carefully” to prevent “exacerbat[ing] unwarranted disparities.” The report also cautioned that “criminal justice data is notoriously poor” and often “inherently subjective.”
- In October 2016, the National Science and Technology Council released the report *Preparing for the Future of Artificial Intelligence*, which summarized the state of AI in the government and made recommendations about AI governance and safety, noting

³ EXEC. OFF. OF THE PRESIDENT, BIG DATA: SEIZING OPPORTUNITIES, PRESERVING VALUES (May 2014), https://obamawhitehouse.archives.gov/sites/default/files/docs/big_data_privacy_report_may_1_2014.pdf.

⁴ EXEC. OFF. OF THE PRESIDENT, BIG DATA: A REPORT ON ALGORITHMIC SYSTEMS, OPPORTUNITY, AND CIVIL RIGHTS (May 2016), https://obamawhitehouse.archives.gov/sites/default/files/microsites/ostp/2016_0504_data_discrimination.pdf.

that the lack of complete quality data in the criminal justice system risked “exacerb[ing] problems of bias.”⁵

- In December 2020, Executive Order 13960, *Promoting the Use of Trustworthy Artificial Intelligence in the Federal Government*, directed federal agencies to adhere to principles when using AI, including ensuring that AI applications are “consistent with the Constitution and all other applicable laws and policies, including those addressing privacy, civil rights, and civil liberties.”⁶ EO 13960 further directed that AI uses be “accurate, reliable, and effective” as well as “safe, secure, and resilient,” “understandable,” “regularly monitored,” “transparent,” and “accountable” through the implementation of “appropriate safeguards.”
- Also in December 2020, Congress enacted the AI in Government Act, which required the Director of the Office of Management and Budget (OMB) to issue a memorandum that “identif[ies] best practices for identifying, assessing, and mitigating any discriminatory impact or bias on the basis of any classification protected under Federal nondiscrimination laws, or any unintended consequence of the use of artificial intelligence.”⁷
- In January 2021, Congress enacted the National Artificial Intelligence Initiative Act of 2020, which established a National AI Advisory Committee under the Department of Commerce.⁸
- In May 2022, President Biden signed EO 14074, *Advancing Effective, Accountable Policing and Criminal Justice Practices to Enhance Public Trust and Public Safety*.⁹ EO 14074 directed the Attorney General to commission a National Academy of Sciences (NAS) study on facial recognition, other biometric identification, and predictive algorithms used by law enforcement. NAS issued a report on facial recognition in January 2024 and convened a workshop on predictive policing in June 2024.¹⁰

⁵ NAT’L Sci. & TECH. COUNCIL, COMM. ON TECH., EXEC. OFF. OF THE PRESIDENT, PREPARING FOR THE FUTURE OF ARTIFICIAL INTELLIGENCE 30 (Oct. 2016), https://obamawhitehouse.archives.gov/sites/default/files/whitehouse_files/microsites/ostp/NSTC/preparing_for_the_future_of_ai.pdf.

⁶ Exec. Order No. 13960, 85 Fed. Reg. 78939, (Dec. 2020), <https://www.federalregister.gov/documents/2020/12/08/2020-27065/promoting-the-use-of-trustworthy-artificial-intelligence-in-the-federal-government>.

⁷ Pub. L. No. 116-260, div. U, title 1, § 104 (codified at 40 U.S.C. § 11301 note), <https://www.congress.gov/116/plaws/publ260/PLAW-116publ260.pdf>.

⁸ Pub. L. No. 116-617, div. C, title LI § 5104, <https://www.congress.gov/116/crpt/hrpt617/CRPT-116hrpt617.pdf#page=1216>.

⁹ Exec. Order No. 14074, 87 Fed. Reg. 32945 (May. 2022), <https://www.federalregister.gov/d/2022-11810>.

¹⁰ NAT’L ACAD. OF SCIS., ENG’G, & MED., FACIAL RECOGNITION TECHNOLOGY: CURRENT CAPABILITIES, FUTURE PROSPECTS, AND GOVERNANCE (Jan. 2024), <https://nap.nationalacademies.org/catalog/27397/facial-recognition-technology-current-capabilities-future-prospects-and-governance>; NAT’L ACAD. OF SCIS., ENG’G, & MED., LAW ENFORCEMENT USES OF PREDICTIVE POLICING APPROACHES (Nov. 2024), <https://nap.nationalacademies.org/catalog/28037/law-enforcement-use-of-person-based-predictive-policing-approaches-proceedings>.

- In October 2022, the White House Office of Science and Technology Policy (OSTP) published the *Blueprint for an AI Bill of Rights*, which recommended a set of principles for responsible AI use, including ensuring safety and efficacy, protecting against algorithmic discrimination, respecting privacy, providing notice and explanation, and establishing human oversight. The report noted that “[d]esigners, developers, and deployers of automated systems should take proactive and continuous measures to protect individuals and communities from algorithmic discrimination and to use and design systems in an equitable way.”¹¹
- In December 2022, Congress enacted the Advancing American AI Act, which directed the Secretary of Homeland Security to “issue policies and procedures ... to ensure that full consideration is given to ... the privacy, civil rights, and civil liberties impacts of artificial intelligence-enabled systems.”¹² The Act also requires federal agencies to publish their AI use case inventories.
- In May 2023, the National Science and Technology Council released a revised *National Artificial Intelligence Research and Development Strategic Plan* to coordinate and focus federal investments in AI. The plan emphasized the importance of developing AI systems “in a manner that mitigates bias and harm and is done in accordance with the civil rights, civil liberties, and interests of those affected by the system.”¹³
- In October 2023, President Biden signed EO 14110, which set out policy priorities and a whole-of-government approach for responsible AI development and implementation.¹⁴ EO 14110 includes over 100 directives to agencies, including tasking the Attorney General with submitting this report.
- In March 2024, OMB issued Memorandum M-24-10, *Advancing Governance, Innovation, and Risk Management for Agency Use of Artificial Intelligence*. The memo fulfills a statutory requirement of the AI in Government Act of 2020 and a directive of EO 14110, and establishes baseline AI governance requirements for federal agencies, including governance structures, inventories, impact assessments, testing in real-world contexts, independent evaluation, consultation with impacted communities, and ongoing monitoring and risk mitigation.¹⁵

¹¹ EXEC. OFF. OF THE PRESIDENT, BLUEPRINT FOR AN AI BILL OF RIGHTS 5 (Oct. 2022),

<https://www.whitehouse.gov/wp-content/uploads/2022/10/Blueprint-for-an-AI-Bill-of-Rights.pdf>.

¹² Pub. L. No. 117-263, div. G, title LXXII, subtitle B, §§ 7224(a), 7224(d)(1)(B), and 7225 (codified at 40 U.S.C. 11301 note), <https://www.congress.gov/117/plaws/publ263/PLAW-117publ263.pdf>.

¹³ NAT’L SCI. & TECH. COUNCIL, SELECT COMM. ON A.I., EXEC. OFF. OF THE PRESIDENT, NATIONAL ARTIFICIAL INTELLIGENCE RESEARCH AND DEVELOPMENT STRATEGIC PLAN 2023 UPDATE 14 (May 2023), <https://www.whitehouse.gov/wp-content/uploads/2023/05/National-Artificial-Intelligence-Research-and-Development-Strategic-Plan-2023-Update.pdf>.

¹⁴ Exec. Order No. 14110, 88 Fed. Reg. 75191.

¹⁵ MEMORANDUM FROM SHALANDA D. YOUNG, DIR., OFF. MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, TO HEADS OF EXEC. DEP’T & AGENCIES, (Mar. 28, 2024), available at <https://www.whitehouse.gov/wp-content/uploads/2024/03/M-24-10-Advancing-Governance-Innovation-and-Risk-Management-for-Agency-Use-of-Artificial-Intelligence.pdf>.

- In September 2024, OMB issued Memorandum M-24-18, *Advancing the Responsible Acquisition of Artificial Intelligence in Government*.¹⁶ The memo complements M-24-10 with further acquisition guidance, including on how to enable testing of AI and ensure data sources are consistent with privacy and civil liberties protections, in fulfillment of section 7224(d) of the Advancing American AI Act.

This report also builds on work within the Department that has been central to the U.S. government's continuing efforts to ensure the responsible use of AI in criminal justice.

- In November 2023, Deputy Attorney General Monaco announced the establishment of DOJ's Emerging Technology Board (ETB), based on recommendations in the Deputy Attorney General's Comprehensive Cyber Review.¹⁷
- In February 2024, Attorney General Merrick Garland announced the designation of the Department's first Chief AI Officer (CAIO).¹⁸ The CAIO and ETB are charged with developing and overseeing a comprehensive program of AI governance for DOJ, including implementation of EO 14110, the accompanying OMB Memoranda M-24-10 and M-24-18, and the National Security Memorandum on *Advancing the United States' Leadership in Artificial Intelligence; Harnessing Artificial Intelligence to Fulfill National Security Objectives; and Fostering the Safety, Security, and Trustworthiness of Artificial Intelligence*.
- In April 2024, the Chief of the Computer Crime and Intellectual Property Section within the Criminal Division sent a letter to the Copyright Office expressing support for a safe harbor for researchers who conduct independent bias testing on AI systems.¹⁹ The letter noted that the Department benefits from this type of research in its work, including in the context of enforcement actions by the Civil Rights Division that are informed by this research.

¹⁶ MEMORANDUM FROM SHALANDA D. YOUNG, DIR., OFF. MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, TO HEADS OF EXEC. DEPT'S & AGENCIES, (Sept. 24, 2024), available at <https://www.whitehouse.gov/wp-content/uploads/2024/10/M-24-18-AI-Acquisition-Memorandum.pdf>.

¹⁷ Press Release, U.S. Dep't Just., Readout of Deputy Attorney General Lisa Monaco's Trip to New York and Connecticut (Nov. 9, 2023), <https://www.justice.gov/opa/pr/readout-deputy-attorney-general-lisa-monacos-trip-new-york-and-connecticut>.

¹⁸ Press Release, U.S. Dep't Just., Attorney General Merrick B. Garland Designates Jonathan Mayer to Serve as the Justice Department's First Chief Science and Technology Advisor and Chief AI Officer (Feb. 22, 2024), <https://www.justice.gov/opa/pr/attorney-general-merrick-b-garland-designates-jonathan-mayer-serve-justice-departments-first>.

¹⁹ Letter from John T. Lynch, Jr., Chief, Comput. Crime & Intell. Prop. Section, Crim. Div., Dep't Just., to Suzanne V. Wilson, Gen. Couns. & Assoc. Reg. Copyrights, Copyright Off., Libr. Cong., (Apr. 15, 2024), <https://www.copyright.gov/1201/2024/USCO-letters/Letter%20from%20Department%20of%20Justice%20Criminal%20Division.pdf>.

- Also in April 2024, DOJ joined five cabinet-level federal agencies in a pledge to uphold a national commitment to core principles of fairness, equality, and justice as use of AI and other emerging technologies continues to increase.²⁰
- In September 2024, the Criminal Division released an update to its guidance on evaluating corporate compliance programs, setting an expectation that programs will address the role of AI in creating and identifying compliance risks.²¹
- In October 2024, the Criminal Division hosted a symposium on AI. The Principal Deputy Assistant Attorney General for the Division announced that the Department will be convening AI researchers and companies to understand how DOJ can best support independent AI accountability research.²² This initiative follows steps the Department has previously taken to support security research, including a charging policy for good-faith research and guidance to industry.
- In October 2024, DOJ issued its plan for implementing a comprehensive AI governance program, consistent with OMB Memorandum M-24-10.²³ The plan describes DOJ's AI governance program, beginning with a thorough inventory of AI uses across the Department. For AI uses with heightened potential for impact on individuals' rights and safety, it provides for qualitative impact assessments, quantitative testing and ongoing monitoring for performance and biases, risk mitigation, and departmentwide coordinated decision-making.

The Use of AI in the Criminal Justice System

The types of AI uses in criminal justice described in EO 14110 fall into four categories, set forth below.²⁴ This report addresses each in turn, provides recommendations, and addresses the establishment of AI governance programs.

- **Identification and Surveillance.** From recognizing faces, fingerprints, and other biometric identifiers, to tracking license plates and locating gunshots, AI has a wide range of existing and potential applications for identification and surveillance in

²⁰ Press Release, U.S. Dep't Just., Five New Federal Agencies Join Justice Department in Pledge to Enforce Civil Rights Laws in Artificial Intelligence (Apr. 4, 2024), <https://www.justice.gov/opa/pr/five-new-federal-agencies-join-justice-department-pledge-enforce-civil-rights-laws>.

²¹ U.S. Dep't Just., Crim. Div., Evaluation of Corporate Compliance Program (Sept. 2024), <https://www.justice.gov/criminal/criminal-fraud/page/file/937501/dl>.

²² Press Release, U.S. Dep't Just., Readout of the Criminal Division's Symposium on Artificial Intelligence in the Justice Department (Oct. 3, 2024), <https://www.justice.gov/opa/pr/readout-criminal-divisions-symposium-artificial-intelligence-justice-department>.

²³ U.S. Dep't Just., Compliance Plan for OMB Memorandum M-24-10 (Oct. 2024), available at <https://www.justice.gov/media/1373026/dl>.

²⁴ Identification and Surveillance addresses "police surveillance" and "prison-management tools," as directed by EO 14110 sections 7.1(b)(i)(E) and (G). Forensic Analysis covers "forensic analysis," as required by section 7.1(b)(i)(H). Predictive Policing addresses "crime forecasting and predictive policing, including the ingestion of historical crime data into AI systems to predict high-density 'hot spots,'" as required by section 7.1(b)(i)(F). Risk Assessment covers "sentencing," "parole, supervised release, and probation," "bail, pretrial release, and pretrial detention," and "risk assessments, including pretrial, earned time, and early release or transfer to home-confinement determinations" as required by sections 7.1(b)(i)(A) through (D).

criminal justice contexts. These uses of AI can be significantly more accurate²⁵ and efficient than human observations and comparisons, and they can provide entirely new capabilities. But these AI uses also pose concerns, especially related to errors, bias, and privacy. When harms associated with these AI uses occur, they can be serious, including mistaken arrests, with a potential for disproportionate impact on certain communities. There is substantial nationwide variation in policies about whether and how AI may be used for identification and surveillance in criminal justice contexts.

- **Forensic Analysis.** AI can improve the capabilities, speed, and accuracy of forensic analysis. It is already being used to enhance DNA comparison, facilitate tracing of seized drugs, and prioritize electronic evidence, among other applications. Ongoing research suggests that future uses could include analysis of physical and trace evidence, medical evaluations, and assessing crime scenes. Forensic analysis must continue to meet exacting standards of accuracy and transparency to ensure due process and satisfy evidentiary requirements. Uses of AI may pose distinct challenges for meeting these requirements because of the complexity of validating and explaining AI-based forensic analysis, as well as the limitations of the data necessary for enabling these types of analyses.
- **Predictive Policing.** Law enforcement agencies use historical data to forecast the places where crime is likely to cluster and people who are at a higher risk of engaging in or being victims of criminal activity. Fundamental police work includes tracking where and when crimes occur, who is involved, and how crimes and people involved with crimes are connected. Developing accurate predictive models based on these types of data may help more efficiently direct resources—including non-law enforcement resources, such as social services—preventing crimes and decreasing response times. But there are also significant risks associated with predictive policing. The data used for predictive policing may have significant gaps and errors, and it may reflect human biases. Use of models based on that data may entrench existing disparities and result in unintended consequences and unjust outcomes, such as over-policing of certain individuals and communities. Successful place-based predictive policing programs integrate a range of strategies and interventions to promote public safety. At the same time, some law enforcement agencies have shifted away from person-based predictive policing, citing limited value and impact on privacy and civil liberties.
- **Risk Assessment.** Risk assessment tools estimate the likelihood that a certain individual outcome will occur in the criminal justice system, such as recidivating or failing to appear in court. These tools are widely used to inform pretrial release, sentencing, prison classification, probation, parole, and supervision. Used properly, risk assessment tools can be more accurate than human judgment alone and can enable more targeted use of various tools within the criminal justice system. Risk assessment tools can also be more transparent and equitable than human judgments. There are, however, significant risks associated with these tools. Risk assessment tools can be inaccurate, especially when they are not validated on local data or fail to take into

²⁵ This report uses the terms “accurate” and “reliable” as shorthand for the predictive performance of an AI system. The report uses more precise terminology when referring to specific performance metrics and measures, such as the precision or false positive rate of an AI system.

account relevant factors. They may be designed to estimate outcomes that are not directly relevant to the decision being made and thus fail to properly inform decision-makers. Risk assessment tools can perpetuate bias and inequality, since the data used in building models may reflect errors or biases and the development process may not incorporate input from affected communities. Models may also be unnecessarily complex, lack transparency, and apply substantially different categorizations to similar people.

II. Identification & Surveillance

Introduction

AI has the potential to enable agencies across the criminal justice system to more effectively and efficiently identify people based on biometrics, i.e., the measurement and analysis of individual physical characteristics.¹ Advances in computer vision, data mining, and complex pattern comparison tools, combined with decreasing costs of cameras and sensors as well as improvements in computation and data storage, have made AI-based biometric identification less expensive and more widely available.² Some agencies with law enforcement, correctional, and community supervision responsibilities now routinely use AI for biometrics.³

The use of AI for biometric identification also has risks. AI could misidentify individuals, which can misdirect law enforcement efforts and impact the civil rights and civil liberties of affected individuals. The performance of AI for biometric identification may also differ across demographic groups. There have been public reports of seven instances of mistaken arrests associated with the use of facial recognition technology, almost all involving Black individuals.⁴ The collection and use of biometric data also poses privacy risks, especially when it involves personal information that people have shared in unrelated contexts.

The first part of the chapter opens with a description of biometric applications of AI in criminal justice, including automated fingerprint identification systems (AFIS), facial recognition technology (FRT), and iris scanning. The chapter then focuses on FRT use in criminal investigations to demonstrate the nuances of evaluating AI-based identification systems for accuracy and biases, as well as the importance of establishing policy frameworks and addressing impacts on privacy, civil rights, and civil liberties.

The second part of the chapter addresses uses of AI in the conduct of law enforcement surveillance. It focuses on automated license plate recognition (ALPR), an increasingly common practice for identifying vehicles that may be involved in criminal activity, including in ongoing emergencies.

¹ See INTERNATIONAL ORGANIZATION FOR STANDARDIZATION, *Biometrics, in* INFORMATION TECHNOLOGY - VOCABULARY (3rd ed. Int'l Org. for Standardization 2022), <https://www.iso.org/standard/73514.html> (“automated recognition of individuals based on their biological and behavioural characteristics”).

² Consistent with the definition of artificial intelligence in OMB Memoranda M-24-10 and M-24-18, this chapter categorizes identification and surveillance methods that involve algorithms or statistical analysis as AI.

³ The following discussion of different biometric approaches includes references to examples of how these technologies are used by federal, state, and local agencies.

⁴ U.S. COMM’N CIV. RTS., THE CIVIL RIGHTS IMPLICATIONS OF THE FEDERAL USE OF FACIAL RECOGNITION TECHNOLOGY 25 (2024), https://www.usccr.gov/files/2024-09/civil-rights-implications-of-frt_0.pdf; see also NAT’L ACADS. SCIS., ENG’G, & MED., *Facial Recognition Technology: Current Capabilities, Future Prospects, and Governance* 83 (2024), <https://doi.org/10.17226/27397> (“NAS 2024 Report”) (describing six then-known cases of mistaken arrests).

Biometrics

a. Automated Fingerprint Identification Systems

Law enforcement has used forms of fingerprint analysis for over a century.⁵ The two common types of analysis today involve: (1) fingerprints collected from a known source in a controlled environment, commonly referred to as ten-prints, and (2) fingerprints from an unknown source collected from a surface or object, commonly referred to as latent fingerprints.⁶ Analysis generally involves capturing friction ridge patterns on a person's skin, observing features including the locations and types of ridges (minutiae), and then comparing features across fingerprints for similarity.⁷

Automated methods for fingerprint comparison became practical in the 1970s and entered widespread use by the 1990s.⁸ Live scan systems, which capture digital images of fingerprints without applying ink, developed in parallel. While identification based on fingerprints generally involves specialized scanning hardware or high-resolution digital cameras, technology may eventually enable ordinary smartphone cameras to capture latent prints.⁹ Machine learning methods may also enable more advanced forms of fingerprint analysis in the future, such as comparing fingerprints for whether they may be from different fingers of the same person.¹⁰

Comparison of ten-prints is highly automated in practice today, often not involving a human examiner unless fingerprints will be offered as forensic evidence in a prosecution. Analysis

⁵ See Press Release, Fed. Bureau Investigation, Crim. Just. Info. Serv. Div., FBI's Criminal Justice Information Services Division Celebrates 100th Anniversary of National Fingerprint Repository (July 10, 2024), <https://www.fbi.gov/news/press-releases/fbi-s-criminal-justice-information-services-division-celebrates-100th-anniversary-of-national-fingerprint-repository> (“In 1924, the FBI established an Identification Division informally called ‘ident’ for many years. ‘Ident’ gathered prints from police agencies nationwide and manually searched them upon request for matches to criminals and crime evidence.”).

⁶ Fingerprints may include impressions of palms, in addition to individual fingers.

⁷ See John R. Vanderkolk, *Chapter 9: Examination Process*, in NATIONAL INSTITUTE OF JUSTICE, THE FINGERPRINTING SOURCEBOOK 9-13 (Nat’l Inst. of Justice 2011), <https://www.ojp.gov/pdffiles1/nij/225329.pdf> (“The direct or side-by-side comparison of friction ridge details to determine whether the details in two prints are in agreement based upon similarity, sequence, and spatial relationship occurs in the comparison phase.”).

⁸ See generally Kenneth R. Moses, *Chapter 6: Automated Fingerprint Identification System (AFIS)*, in NATIONAL INSTITUTE OF JUSTICE, THE FINGERPRINTING SOURCEBOOK 6-1 (Nat’l Inst. Of Justice 2011), <https://www.ojp.gov/pdffiles1/nij/225326.pdf>.

⁹ Robert Pitts et al., *Empirical Comparison of DSLRs and Smartphone Cameras for Latent Prints Photography*, 3 WIREs FORENSIC SCI. (2021), <https://wires.onlinelibrary.wiley.com/doi/epdf/10.1002/wfs2.1391> (“argu[ing] that the cameras equipped in current and future mobile devices are adequate for the purpose of latent print documentation and identification, making it a useful complement, if not a replacement, to DSLRs currently used by crime scene investigators and fingerprint examiners.”); Maryah E. M. Haertel, Eduardo J. Linhares & Andre L. de Melo, *Smartphones for Latent Fingerprint Processing and Photography: A Revolution in Forensic Science*, 3 WIREs FORENSIC SCI. 2 (2021), <https://wires.onlinelibrary.wiley.com/doi/epdf/10.1002/wfs2.1410> (“The use of smartphones in the search and acquisition of latent fingerprints is still new, but various studies show its possibilities.”).

¹⁰ Gabe Guo et al., *Unveiling Intra-person Fingerprint Similarity via Deep Contrastive Learning*, 10 SCI. ADVANCES (2024), <https://www.science.org/doi/10.1126/sciadv.adi0329> (stating that “fingerprints from different fingers of the same person share very strong similarities” and suggesting that intra-person fingerprint similarities “can also help narrow down the candidate list generated by automated fingerprint identification systems.”).

of latent prints, by contrast, usually involves review by a trained examiner because prints may be incomplete or degraded.¹¹ Automated methods may generate possible leads for subsequent analysis and support examiners in comparing fingerprints.

Fingerprint comparison is widely used for other purposes in the criminal justice system. Fingerprint-based checks are the standard for background checks in criminal justice, as well as for other positions of public trust, including teachers, childcare workers, and those in other sensitive occupations. Fingerprints are also the standard for identification based on criminal history record information in the United States. Fingerprint-based verification using 1-2 fingers may be used in certain applications to confirm the identity of an authorized user, track chain-of-custody of certain types of evidence, or limit access to sensitive areas.¹²

The FBI's Next Generation Identification (NGI) system provides fingerprint services to law enforcement agencies nationwide.¹³ NGI contains over 217 million unique fingerprint identity records, over 28 million unique palm print identity records, and over 1.2 million unidentified latent prints. Almost every (if not every) state has its own automated fingerprint identification system, and these systems are also common at local law enforcement agencies.¹⁴

b. Facial Recognition Technology

Facial recognition technology uses methods from computer vision and other areas of AI to isolate and compare faces in photos or video. FRT became available for criminal justice use in the 2000s and became more widely used in the 2010s.¹⁵ Algorithms have rapidly advanced in recent

¹¹ There is ongoing research and debate about statistical models to estimate the likelihood of fingerprint features from population base rates and validation of fingerprint comparison as practiced in particular laboratories. See PRESIDENT'S COUNCIL OF ADVISORS ON SCI. & TECH., FORENSIC SCIENCE IN CRIMINAL COURTS: ENSURING SCIENTIFIC VALIDITY OF FEATURE-COMPARISON METHODS (Sept. 2016), https://obamawhitehouse.archives.gov/sites/default/files/microsites/ostp/PCAST/pcast_forensic_science_report_final.pdf; William Thompson et al., *Latent Fingerprint Examination*, in AAAS, FORENSIC SCIENCE ASSESSMENTS: A QUALITY AND GAP ANALYSIS (AAAS 2017), https://www.aaas.org/sites/default/files/reports/Latent%20Fingerprint%20Report%20FINAL%209_14.pdf; Bradford T. Ulery et al., *Accuracy and Reliability of Forensic Latent Fingerprint Decisions*, 108 PROCS. NAT'L ACAD. SCI. 7733 (2011), <https://www.pnas.org/doi/full/10.1073/pnas.1018707108>.

¹² The fingerprint readers used in these applications often lack the resolution necessary for criminal investigative uses of fingerprint analysis. As discussed below, there is a significant distinction between the 1:N matching that is common in investigations and the 1:1 matching that is common for these applications.

¹³ FBI, FY 2025 President's Budget Request 55 (Mar. 2024), https://www.justice.gov/d9/2024-03/fbi_fy_2025_presidents_budget_narrative_3-5-24_final_1.pdf ("The NGI System services connectivity for 106,981 Federal, State, local, and Tribal law enforcement customers. These customers have existing statutory authorization to conduct background checks using the NGI System; however, only about one third (38,108) of those regularly do.").

¹⁴ Nat'l Inst. of Justice, Latent Fingerprint Interoperability Survey: A National Study of Automated Fingerprint Information Systems (AFIS) Maintained by Law Enforcement Agencies 36 (2014), <https://www.ojp.gov/pdffiles1/nij/247910.pdf> (showing a map of AFIS vendor information for state agencies which excluded Vermont (did not provide an answer), Minnesota and the District of Columbia (did not participate)).

¹⁵ See generally Nat'l Inst. of Justice, History of NIJ Support for Face Recognition Technology (Mar. 5, 2020) at <https://nij.ojp.gov/topics/articles/history-nij-support-face-recognition-technology> (discussing NIJ role in face algorithm research and development); Statement of Jerome M. Pender Before the Senate Judiciary Committee, Subcommittee on Privacy, Technology, and the Law 112 Cong. (2012), <https://archives.fbi.gov/archives/news/testimony/what-facial-recognition-technology-means-for-privacy-and-civil->

years, significantly increasing the ability of FRT systems used in criminal justice to correctly match faces.¹⁶

In the usual design of an FRT system, an algorithm first detects a person's face in a photo or video and extracts the relevant region.¹⁷ Next, it computes a quantitative representation of the face (the "template"). Finally, the algorithm compares templates, producing a "similarity score" for pairs of templates. Some recent FRT systems use deep learning models that integrate these steps.¹⁸

There are, broadly, two types of FRT uses. One-to-one ("1:1") FRT compares a captured ("probe") image to a single other image or template, typically to verify a person's identity. One-to-many (also called "one-to-n" or "1:N") FRT compares a captured image to a database ("gallery") of known images.¹⁹

In the criminal justice system, one-to-one FRT has several applications. The Federal Bureau of Prisons uses one-to-one FRT to confirm employees' identities before entering secure areas of a facility.²⁰ Probation services may use FRT to allow individuals under court-ordered supervision to verify their identity via smartphone rather than requiring physical contact with a probation or pretrial officer.²¹ Similarly, Customs and Border Protection and the Transportation Security Administration use one-to-one FRT to confirm traveler identities.²²

liberties (discussing advances in the NGI program through July 2012); U.S. Gov. Facial Recognition Legal Series (Aug. 31, 2011), https://ucr.fbi.gov/fingerprints_biometrics/biometric-center-of-excellence/files/Forum_1_Minutes.pdf; William Casey et al., *Facial Recognition Technology — Baseline Uses and Legal Challenges: Meeting Minutes*, in U.S. GOVERNMENT FACIAL RECOGNITION LEGAL SERIES (2011), https://ucr.fbi.gov/fingerprints_biometrics/biometric-center-of-excellence/files/Forum_1_Minutes.pdf (noting advances made in facial recognition technology through August 2011). The FBI, for example, began developing the NGI-IPS system in 2008 and began using and providing access to the system in 2011. NGI-IPS became fully operational in 2015, at which point 7 states had access to the system. See GOV'T ACCOUNTABILITY OFF., PUB. NO. GAO-19-579T FACE RECOGNITION TECHNOLOGY: DOJ AND FBI HAVE TAKEN SOME ACTIONS IN RESPONSE TO GAO RECOMMENDATIONS TO ENSURE PRIVACY AND ACCURACY, BUT ADDITIONAL WORK REMAINS (June 2019), <https://www.gao.gov/assets/gao-19-579t.pdf>.

¹⁶ The National Institute of Standards and Technology has performed FRT algorithm evaluations for over 30 years, and currently publishes FRT algorithm evaluations on an ongoing basis through the Facial Recognition Technical Evaluation (FRTE) program. In 2024, FRT algorithms commonly have a below 1% false negative rate with a false positive rate of 3 in 1,000. By comparison, the best performing algorithm in NIST's 2017 challenge had a 22% false negative rate with a false positive rate of 1 in 1000. *Face Technology Evaluation – FTRE/FATE*, Nat'l Inst. Standards & Tech, <https://www.nist.gov/programs-projects/face-technology-evaluations-frtefate>.

¹⁷ See NAS 2024 Report *supra* note 4, at 32-34.

¹⁸ See Mei Wang & Weihong Deng, *Deep Face Recognition: A Survey*, 429 NEUROCOMPUTING 215 (2021), <https://doi.org/10.1016/j.neucom.2020.10.081> (discussing the emergence of deep learning models in FRT in 2012).

¹⁹ One-to-one and one-to-many FRT systems are closely related, because one-to-many FRT systems are often based on one-to-one comparisons. The acceptable levels of performance and demographic differences for these systems may significantly differ by use case.

²⁰ GOV'T ACCOUNTABILITY OFF., PUB. NO. GAO-21-518 FACIAL RECOGNITION TECHNOLOGY: FEDERAL LAW ENFORCEMENT AGENCIES SHOULD BETTER ASSESS PRIVACY AND OTHER RISKS 20, <https://www.gao.gov/assets/gao-21-518.pdf>.

²¹ *Id.* at 19.

²² *Id.* at 19–20; see also DHS Directive 026-11, Use of Face Recognition and Face Capture Technologies, https://www.dhs.gov/sites/default/files/2023-09/23_0913_mgmt_026-11-use-face-recognition-face-capture-technologies.pdf.

One-to-many FRT is used by law enforcement agencies to identify or match people in images and video. These systems usually return a fixed number of candidate matches or candidate matches above a threshold similarity score. The results may be ordered by score and may show a numerical score or a score category.

While there is not comprehensive public data on FRT use by law enforcement nationwide, surveys indicate that FRT is widely used by federal, state, and local agencies.²³ A number of agencies operate their own FRT systems, based on photos from driver's licenses, arrests, and other government interactions. Agencies commonly have access to FRT systems maintained by other agencies, such as FBI's Next Generation Identification-Interstate Photo System (NGI-IPS), which can in turn incorporate results from other agencies. NGI-IPS, for example, incorporates results from 17 state agencies and two federal agencies and encompasses over 67 million arrest photos.²⁴ Commercial vendors also offer FRT services to law enforcement agencies which, as discussed further below, can heighten privacy impacts.

At the federal level, law enforcement agencies use FRT in support of their missions and pursuant to applicable policies. The FBI, for instance, uses one-to-many FRT to help identify perpetrators, victims, and witnesses as part of authorized investigations of criminal offenses.²⁵ The FBI also uses FRT to help identify and locate missing persons or other at-risk individuals, such as abducted children, or victims of child sexual abuse or human trafficking, and to identify deceased or incapacitated individuals. These uses of FRT can be faster and more efficient than other investigative methods, and they can provide unique leads that may not have been available or may have been impractical to obtain through other avenues. The FBI also uses FRT to structure and organize large volumes of lawfully obtained photo or video data, allowing investigators to more efficiently interpret the collected data.²⁶ Under the Department of Justice interim FRT policy, uses of FRT must be lawful and consistent with other DOJ policies. Among other requirements, FRT results alone may not be relied upon as the sole proof of a person's identity; activity protected by the First Amendment may not be the sole basis for using FRT; and personnel who use or approve FRT systems must receive relevant training; among other requirements.²⁷

²³ See GAO-21-518 *supra* note 20 (finding that, in a GAO survey of 42 federal agencies with law enforcement responsibilities, 20 used FRT between 2015 and 2020); CLARE GARVIE, ALVARO BEDOYA & JONATHAN FRANKLE, THE PERPETUAL LINE-UP (2016) 93, 97 <https://www.perpetuallineup.org/sites/default/files/2016-12/The%20Perpetual%20Line-Up%20-%20Center%20on%20Privacy%20and%20Technology%20at%20Georgetown%20Law%20-%2020121616.pdf> (reporting that, based on records requests to 106 state and local law enforcement agencies, at least 53 used, previously used, or planned to use FRT). There is limited data available about FRT use by Tribal law enforcement agencies.

²⁴ U.S. Dep't Just., Written Testimony in Connection with the United States Commission on Civil Rights' Examination of Civil Rights Implications of the Federal Use of Facial Recognition Technology (Mar. 21, 2024).

²⁵ The FBI's use of FRT is governed by the Department of Justice's interim FRT policy, which, among other requirements, mandates that FRT results alone may not be relied upon as sole proof of identity. Rather, an individual's identity must be confirmed through other analysis and/or investigation.

²⁶ This use case is a variation of 1:N FRT, where the gallery of images is drawn from the collected evidence in an investigation rather than from an established repository. The use case supports organizing and triaging media that has been collected in an investigation, which may be voluminous and unorganized.

²⁷ U.S. Dep't Just., Written Testimony in Connection with the United States Commission on Civil Rights' Examination of Civil Rights Implications of the Federal Use of Facial Recognition Technology (Mar. 21, 2024).

State and local law enforcement agencies similarly use FRT to support investigations. There is wide variation in applicable state laws, local ordinances, and law enforcement agency policies.²⁸ Differences include whether and when FRT may be used, protections for expressive activities, quality reviews of results from FRT, how the results from FRT may be used, which probe images may be used, which databases may be searched, what training is necessary, what information is recorded and audited, what information must be disclosed in discovery, and what public transparency must be provided. At one end of the spectrum, some agencies may use FRT under generally applicable laws and policies, but without a law or policy specific to FRT. At the other end, some jurisdictions have entirely prohibited law enforcement agencies from using FRT.

c. Iris Scanning

Iris scanning examines the unique tissue patterns in the donut-shaped part of an eye surrounding the pupil. Iris patterns do not appear to meaningfully change over time, and are protected by the cornea, limiting the potential for damage or mutilation.²⁹ Iris as a biometric modality is relatively new compared with FRT and other biometric modalities, with national-level matching capabilities coming online at the FBI in just the last 5 years.³⁰ The FBI's NGI Iris Service has over 3 million sets of iris images from over 2 million people.³¹

Iris scans are well-suited for identity confirmation in custodial settings, because they are highly accurate when collected properly and can be taken either from a short stand-off distance or without removing handcuffs. Iris scans can also be an effective supplement to other identification modalities for immigration and border screening.

Focusing on FRT for Identification in Law Enforcement Investigations

While FRT use by law enforcement agencies has significant benefits in developing leads, it also poses challenges for responsible use and governance of technology in criminal justice.

FRT poses significant privacy concerns because of the quantity, and likely long retention period, of data required for the system to be effective. FRT enables identifying people without

²⁸ See generally Mailyn Fidler & Justin (Gus) Hurwitz, *An Overview of Facial Recognition Technology Regulation in the United States*, in CAMBRIDGE HANDBOOK OF FACIAL RECOGNITION IN THE MODERN STATE (Rita Matulionyte & Monika Zalierute, eds., 2024), <https://doi.org/10.1017/9781009321211.018>; GARVIE ET AL., *supra* note 23 at 121–50; Jameson Spivack & Clare Garvie, *A Taxonomy of Legislative Approaches to Face Recognition in the United States*, in REGULATING BIOMETRICS: GLOBAL APPROACHES AND OPEN QUESTIONS (Amba Kak ed., 2023), <https://ainowinstitute.org/wp-content/uploads/2023/09/regulatingbiometrics-spivack-garvie.pdf>; JAKE LAPERRUQUE, CEN. DEM. & TECH., LIMITING FACE RECOGNITION SURVEILLANCE: PROGRESS AND PATHS FORWARD, (Aug. 23, 2022), <https://cdt.org/insights/limiting-face-recognition-surveillance-progress-and-paths-forward/>.

²⁹ *Iris Recognition*, NEC (Sept. 22, 2021), <https://www.nec.com/en/global/solutions/biometrics/iris/index.html> (“A person’s iris pattern is unique and remains unchanged throughout life. Also, covered by the cornea, the iris is well protected from damage, making it a suitable body part for biometric authentication.”).

³⁰ *The Eyes Have It: Iris Biometric Added to Next Generation Identification System* (Dec. 11, 2020), <https://www.fbi.gov/news/stories/fbi-adds-iris-biometric-to-next-generation-identification-system-121120>.

³¹ FBI, FY 2025 Budget Request *supra* note 13 at 57 (“As of November 30, 2023, the NGI Iris Service consists of over 3.3 million sets of iris images representing more than 2.6 million unique identities.”).

interacting with them or an object on which they left their fingerprints or DNA. The scale of FRT databases can be large, and the effort and cost of running FRT searches can be low.

Civil liberties are another area of concern. For instance, FRT could be misused to enable identification of people engaged solely in protected expressive activity. Going back to the Founding era, the United States has a rich tradition of anonymous civic discourse and protest, where privacy facilitates the expression of ideas and the assembly of groups. As the Supreme Court has recognized, there is a “vital relationship between freedom to associate and privacy in one’s associations.”³²

Civil rights are another significant issue, in part due to possible biases in FRT systems and how they are used. For example, as noted above, public reporting indicates that there have been seven documented instances of mistaken arrests associated with the use of facial recognition technology, almost all involving Black individuals.³³ Many FRT systems deployed in the United States have higher false match (i.e., false positive) rates when applied to racial minorities, including people who are Black, Native American, Asian American, and Pacific Islanders.³⁴ Research has also demonstrated that FRT systems tend to perform worse on women, children, and the elderly, and some FRT algorithms used in the United States have biases also associated with eyewear, hairstyle, and other attributes.³⁵ Testing by the National Institute of Standards and Technology (NIST), discussed further below, indicates that in the last five years, some FRT developers have made significant progress in addressing differences in performance associated with demographics. Low absolute false match rates and lower relative rates—at least in the controlled settings of NIST’s testing of recent algorithms—now exist across demographics, which include gender, age, and race.

The data used to train FRT systems is an important contributing factor. An FRT system generally performs best on faces that are similar to the faces used when training the system. Race, gender, age, and other attributes are often readily apparent from faces, unlike fingerprints and irises, which compounds the risk that these types of biometric systems will reflect demographic biases from training data. Research has demonstrated an “other-race” effect in FRT systems where, for example, systems built in the United States can perform better on white faces and systems

³² *NAACP v. Alabama ex rel. Patterson*, 357 U.S. 449, 462 (1958).

³³ NAS 2024 Report *supra* note 4.

³⁴ NAS 2024 Report *supra* note 4; *Face Technology Evaluation – FTRE/FATE*, NAT’L INST. STANDARDS & TECH, <https://www.nist.gov/programs-projects/face-technology-evaluations-frtefate>.

³⁵ *E.g.*, NAS 2024 Report *supra* note 4; NIST *supra* note 34; Cynthia M. Cook et al., *Demographic Effects in Facial Recognition and Their Dependence on Image Acquisition: An Evaluation of Eleven Commercial Systems*, 1 IEEE TRANSACTIONS ON BIOMETRICS, BEHAV., & IDENTITY SCI. 32, 32-41 (2019), <https://ieeexplore.ieee.org/document/8636231>; Pawel Drozdowski et al., *Demographic Bias in Biometrics: A Survey on an Emerging Challenge*, 1 IEEE TRANSACTIONS ON TECH. & SOC’Y 89 (2020), <https://ieeexplore.ieee.org/document/9086771>; Philipp Terhörst et al., *A Comprehensive Study on Face Recognition Biases Beyond Demographics*, 3 IEEE TRANSACTIONS ON TECH. & SOC’Y. 16, 16-30 (2022), <https://ieeexplore.ieee.org/stamp/stamp.jsp?arnumber=9534882>. Related research has demonstrated similar disparities in other computer vision applications, such as gender classification. *E.g.*, Joy Buolamwini & Timnit Gebru, *Gender Shades: Intersectional Accuracy Disparities in Commercial Gender Classification*, 81 PROCEEDINGS MACH. LEARNING RSCH. 1, 1 (2018), <https://proceedings.mlr.press/v81/buolamwini18a/buolamwini18a.pdf>.

developed in Asia can perform better on Asian faces.³⁶ Balancing the demographics of training datasets can be a valuable step in improving performance disparities, though additional possible sources of bias remain in both FRT systems and how they are used.³⁷

Agencies considering use of FRT must grapple with difficult policy decisions. Agencies should develop and enforce policies regarding the use of FRT and set clear rules for how and when the technology may be used, including guardrails that protect civil rights and liberties. Policies should, among other things, be transparent, be reflected in public documentation to the extent possible, include requirements for evaluating FRT uses, and provide for ongoing monitoring and mitigation of risks. In particular, policies should address the topics and uses described in greater detail below.³⁸

a. Algorithm Evaluation

Before a law enforcement agency begins using an FRT system, it is essential to understand whether the benefits and risks of the system are appropriate for the intended use. Evaluating the algorithms in an FRT system is an important step and includes quantifying how well the system correctly matches faces when a match exists in the database, how well it rejects incorrect matches, and how these types of performance differ across demographic groups.³⁹ An FRT system with better performance can provide greater value in investigations and reduce the likelihood of harmful consequences.

NIST's ongoing Face Recognition Technology Evaluation (FRTE) program provides valuable algorithm performance benchmarks. The program includes one-to-many performance testing on over 350 FRT algorithms and one-to-one demographic testing (race, gender, and age) on over 500 algorithms.⁴⁰ The results of NIST's testing can help law enforcement agencies understand and compare algorithms.

As valuable as it is, the NIST testing program has important limitations. The datasets used in testing predominantly consist of images from a controlled or semi-controlled environment, where the subject is close to the camera, the subject is looking at or near the camera, and lighting is adequate. These datasets may be sufficiently representative for some criminal justice

³⁶ See PATRICK GROTH ET AL., NAT'L INST. STANDARDS & TECH, FACE RECOGNITION VENDOR TEST (FRVT) PART 3: DEMOGRAPHIC EFFECTS (2019) <https://doi.org/10.6028/NIST.IR.8280>; P. Jonathan Phillips et al., *An Other-race Effect for Face Recognition Algorithms*, 8 ACM TRANSACTIONS ON APPLIED PERCEPTION 1 (2011), <https://dl.acm.org/doi/10.1145/1870076.1870082>; NIST *supra* note 34.

³⁷ See Valeriia Cherepanova et al., *A Deep Dive into Dataset Imbalance and Bias in Face Identification*, AIES (2023) <https://dl.acm.org/doi/fullHtml/10.1145/3600211.3604691>.

³⁸ Law enforcement agencies can use existing templates to guide policy development, including the FRT Policy Development Template. *Face Recognition Policy Development Template*, BUREAU JUST. ASSISTANCE (Dec. 2017), <https://bja.ojp.gov/sites/g/files/xyckuh186/files/Publications/Face-Recognition-Policy-Development-Template-508-compliant.pdf>.

³⁹ In NIST's 1:N FRTE program, the core metrics are the false negative identification rate and the false positive identification rate. There are many other possible performance and bias metrics to compute for FRT systems, beyond those included in NIST's valuable program, like for other AI systems. There are also other types of bias to consider, such as people with disabilities. Testing with other metrics and for other types of biases may be appropriate depending on the intended use.

⁴⁰ See NIST *supra* note 34. As NIST explains in its demographic testing methods, it is possible for many FRT algorithms to extrapolate one-to-many demographic results from one-to-one results.

applications where test results may be indicative of the system's performance, such as FRT use to match booking photos to driver's license photos.

The datasets used in NIST's current testing may, however, not be representative of law enforcement investigative uses of FRT. These uses often involve images that are not from controlled environments, such as still images from surveillance camera footage. The subject may be distant, looking away, dimly or unevenly lit, or located below the camera. The camera may be of low quality or introduce distortions, and the image or the subject's face may be occluded. These significant differences in context make it difficult to generalize the results of NIST's current testing to FRT performance and biases when used in law enforcement investigations.

NIST is in the process of reestablishing a line of performance testing for FRT use on images from videos, including from surveillance cameras. This type of testing, the Face In Video Evaluation (FIVE), may be more representative of law enforcement investigation settings.⁴¹

Law enforcement agencies considering use of FRT should establish testing requirements for performance and biases. This testing should, to the extent possible, be representative of real-world deployment contexts and follow standardized methodologies. NIST's testing program is an important starting point, and the ISO/IEC 19795 standards on biometric performance testing also provide valuable guidance. Agency policies should specify the nature of and benchmarks for testing and should ensure retention and disclosure of testing results to the extent feasible. Policies should also require that vendors provide evaluation results for the system version that is procured by agencies, not results for a previous or adjusted versions of the system, and that the documentation and results from vendors be sufficient to allow for independent evaluation and/or auditing.

Continuous monitoring, discussed further below, can provide additional information about the real-world performance and biases of FRT systems. When pre-deployment testing is not fully representative of how an FRT system will be used, post-deployment monitoring is especially important.

At the federal level, OMB Memorandum M-24-10 identifies law investigative uses of FRT as AI use cases that presumptively require heightened risk management practices, including performance and bias testing in real-world conditions. OMB Memorandum M-24-18 further directs federal agencies that procure FRT capabilities to ensure that they have been tested by NIST, where practicable. State, local, Tribal, and territorial law enforcement agencies should also implement these practices, and federal grantmaking agencies should require these practices when providing financial support for the procurement or use of FRT, accounting for the differing missions and resources of grant recipients.⁴²

⁴¹ FIVE will report performance measurements for FRT algorithms but, importantly, may not include measurements of demographic differences.

⁴² For example, in some instances, it may be appropriate for a state, local, Tribal, or territorial law enforcement agency to evaluate a prospective FRT use on the basis of real-world FRT testing conducted by, on behalf of, or in coordination with other law enforcement agencies, provided that the testing is representative of the agency's FRT uses.

b. Database Selection

Agencies may have access to multiple FRT systems with different databases, and some FRT systems have the capability to run searches against multiple databases. Depending on the agency and use case, searches can be run against criminal records, driver's license photographs, other agencies' databases, or commercial databases.

While the likelihood of developing a useful investigative lead generally increases with an expanded search, so too does the likelihood of an FRT system returning candidate matches with high similarity scores—and possibly high visual similarity for human reviewers—that are not the person to be identified. A larger set of databases or data can also increase the risk that a user may unintentionally exceed their authority by searching a dataset or accessing a photo for an unauthorized purpose.

Agency policies should clearly articulate to users what datasets are available for each type of search, as well as the purposes for which a search of a given dataset is authorized. Agency agreements to access external data sources should also contain appropriate restrictions on the use of the data being accessed. Agencies should establish processes to log FRT uses, enabling auditing to ensure that searches have been conducted for authorized purposes.

In selecting databases to use, agencies should give careful consideration to how the underlying data was collected. Some commercial FRT services make use of databases that contain millions, or even billions of images scraped from social media and other online services and websites. This repurposing of personal photographs, in a context different from the ones in which they were originally created and shared—potentially without consent and contrary to expectations—potentially raises questions of law, policy, and ethics. These types of commercial services also generally have less reliable, less complete, and less current information associated with photos than do FRT systems based on government identification records, which can misdirect investigative efforts. Some law enforcement agencies prohibit the use of these types of systems, and others permit their use only in certain types of investigations. OMB Memoranda M-24-10 and M-24-18 specifically direct federal agencies to carefully consider whether and when use of these types of FRT systems is appropriate.

Law enforcement agencies should not use FRT systems trained on photos or built with other information that was collected in violation of laws, federal government guidance, or agency policy. Agencies should also specifically articulate the authority that permits the collection of FRT biometric data or associated personally identifiable information, which should be reflected in public documentation whenever possible.

c. Use of Facial Recognition

Law enforcement agencies substantially differ in their policies regarding when an investigation may make use of facial recognition. At some agencies, facial recognition is available for all criminal investigations. At others, only certain types of investigations may make use of FRT, such as for violent crimes and child safety. As noted above, at some law enforcement agencies, FRT use is prohibited.

There is further divergence in the predication standards that law enforcement agencies implement for using facial recognition. At some agencies, investigators have discretion about turning to FRT. At others, investigators must meet a reasonable suspicion standard. And at some agencies, investigators must have probable cause to conduct an FRT search.

Real-time use of FRT is another area where law enforcement agencies have differing policies. Some prohibit it, while others allow it but restrict use directed at protected speech activities.

Law enforcement agencies that use FRT should establish policies that clearly specify when FRT may be used. These policies should be public and easily accessible to the greatest extent possible.

Agency policies should describe the types of investigations in which FRT use is appropriate, taking into account factors such as the type of the criminal offense, the likelihood of generating a true match, the evaluated performance of the FRT system, and the quality of relevant data. Agency policies should also delineate the circumstances in which it is appropriate to conduct FRT searches, such as to provide a lead for identifying a witness, perpetrator, victim, or a person who is missing or otherwise believed to be at risk of harm. Policies should additionally specify what predication is necessary to take the step of an FRT search.

A policy should also describe the level of supervisory review, if any, necessary before conducting a search. Policies should also distinguish between the different types of FRT use and account for the varying levels of risk of harm. For example, using FRT to identify an unknown perpetrator need not be treated the same as using FRT to organize and triage collected media.

In addition to specifying when FRT may be used, policies should specify when it may not be used. At minimum, policies should prohibit the use of FRT solely based on constitutionally protected activities (e.g., a First Amendment-protected protest), to facilitate unlawful discrimination, or in any other way that would be inconsistent with legal requirements or other policies.

d. Image Quality

The quality of the probe image submitted to an FRT system is a critical factor in whether the system can return a match. Many factors can affect the performance and biases of an FRT system, including photo resolution and clarity; the subject's pose and attire; lighting; occlusions; and the camera's position, sensor, and lens. Additionally, modifications to the probe image—including to size, aspect ratio, or coloration—could potentially have an adverse impact on FRT results. In general, searches with lower-quality images are less likely to return matches. In a law enforcement investigation, however, only lower-quality images may be available.

Law enforcement agencies can mitigate risks associated with probe photo selection by establishing minimum quality criteria.⁴³ Where possible, these criteria should be set by an

⁴³ The Facial Identification Scientific Working Group (FISWG)—a consortium of state, local, federal, and international law enforcement agencies as well as FRT vendors and academics—provides voluntary image quality

independent entity with expertise in FRT, such as a testing or standards-setting organization. These criteria may differ depending on the type of use, the feasibility of alternative investigative steps, and agency resources. An agency might conclude, for example, that the minimum quality criteria for locating a missing person or identifying a victim of child sex trafficking should be different from the criteria for identifying a witness to a nonviolent crime.

The policies that law enforcement agencies establish for FRT should also address use of probe images that are not photographs of the person to be identified. These images may include sketch drawings, generated images, or images of people who are “lookalikes” for the subject. Searches with these types of images can be more prone to incorrect matches,⁴⁴ so a policy should establish when (if ever) they are permitted and should require heightened safeguards when they are used.

Law enforcement agencies should also implement policies that prohibit the use of FRT with probe images that were collected in violation of law or another applicable policy.

e. Quality Control for Results

The FRT systems used in law enforcement investigations are inexact, and a system may return candidate matches that are not the subject.⁴⁵ Law enforcement agencies should apply a minimum similarity threshold for candidate results, which may vary depending on the nature of the investigation and should only be overridden in exigent circumstances.

Human review is also essential. Law enforcement agencies should require that an examiner who is trained to compare faces and mitigate bias—and who ideally is independent of the case team—manually reviews results before they are used in an investigation. When there are multiple candidate results from an FRT system, an examiner should review all top results, and similarly a case team should consider all candidates returned by an examiner before focusing on one candidate result for further investigation.

standards. *Image Factors to Consider in Facial Image Comparison*, FACIAL IDENTIFICATION SCI. WORKING GRP. (May 28, 2021),

https://fiswg.org/fiswg_image_factors_to_consider_in_facial_img_comparison_v1.0_2021.05.28.pdf

The International Standardization Organization has worked with other organizations to develop the ISO/IEC 30137 series, which outlines effective video system performance for FRT and other uses. *See generally ISO/IEC JTC 1/SC 37: Biometrics*, INT’L STANDARDIZATION ORG. & INT’L ELECTROTECHNICAL COMM’N, <https://www.iso.org/committee/313770.html>. *ISO/IEC 30137-1:2024: Information technology — Use of biometrics in video surveillance systems*, INT’L STANDARDIZATION ORG. & INT’L ELECTROTECHNICAL COMM’N (2024), <https://www.iso.org/standard/87734.html>.

⁴⁴ *See, e.g.*, NAT’L INST. STANDARDS & TECH, NIST INTERAGENCY REPORT 8009, FACE RECOGNITION VENDOR TEST 4 (2014), <https://nvlpubs.nist.gov/nistpubs/ir/2014/NIST.IR.8009.pdf> (documenting high error rates in sketch recognition); Christian Galea & Reuben A. Farrugia, *Forensic Face Photo-Sketch Recognition Using a Deep Learning-Based Architecture*, 24 IEEE SIGNAL PROCESSING LETTERS 1586 (2017), <https://ieeexplore.ieee.org/document/8025793> (discussing performance deep learning for face photo-sketch recognition); CLARE GARVIE, GARBAGE IN, GARBAGE OUT (2019), <https://www.flawedfacedata.com/> (“Even the most detailed sketches make poor face recognition probe images. . . . The most likely outcome of using a forensic sketch as a probe photo is that the system fails to find a match—even when the suspect is in the photo database available to law enforcement.”).

⁴⁵ *See* NAS 2024 Report *supra* note 4, at 1,6.

Law enforcement agencies should take steps to minimize the risks of automation bias (i.e., examiner deference to system output) and confirmation bias (i.e., reinforcing an examiner’s beliefs about a subject or the system). These steps could include removing similarity scores or ranking information from the results shown to examiners, reminding examiners of the limitations of FRT systems, and providing appropriate training on facial comparison and mitigating biases (discussed further below).

f. Uses of FRT Results

In addition to specifying when and how investigators can use FRT, law enforcement agency policies should establish permissible uses of results. Because of the limitations of FRT, policies should specify that FRT search results should be considered a lead and not sufficient to establish probable cause or a positive identification without corroboration. Policies should describe when and how FRT can support probable cause. This can be a complex issue where FRT plays a role in witness identification, which may also involve comparing faces. At least one law enforcement agency prohibits conducting a lineup based solely on an FRT investigative lead without independent and reliable evidence linking the suspect to the crime.⁴⁶

g. End-to-End Evaluation and Continuous Monitoring

When law enforcement agencies use FRT systems in investigations, it is essential to understand the technical performance and biases of these systems, as noted above. It is also important to understand the broader context for and impacts of FRT use, including why investigators use it, how it affects investigations, and how it affects the people who appear in results. This type of “end-to-end” operational evaluation is a best practice for AI governance and, at the federal level, encouraged by OMB Memorandum M-24-10.

Law enforcement agencies should consider implementing end-to-end evaluation for uses of FRT. This evaluation could address questions like: What are the types of cases where investigators turn to FRT, and why do they use FRT instead of other investigative methods? How valuable are the leads from FRT in advancing investigations? How often does FRT generate a lead for investigators that could not have been developed otherwise, or that would have taken considerably more time or resources otherwise? How often does FRT lead investigators to focus on a person who is later determined to not be relevant to an investigation? Answering basic questions like these can be important for evaluating the benefits and risks of FRT use, and can help reinforce community trust by demonstrating the practical impacts of FRT.

Continuous monitoring is another form of evaluation that law enforcement agencies should consider. An FRT system’s behavior is affected by both probe and gallery photos. If there are changes in either type of photo—for example, if a law enforcement agency starts focusing on a particular type of photo as evidence or if a state changes its driver’s license photo format—that can impact the performance and biases of the system. FRT vendors also update their algorithms, which can also affect performance and biases. Continuously keeping track of real-world

⁴⁶ U.S. COMM’N CIV. RTS., THE CIVIL RIGHTS IMPLICATIONS OF THE FEDERAL USE OF FACIAL RECOGNITION TECHNOLOGY 112-13 (2024), https://www.usccr.gov/files/2024-09/civil-rights-implications-of-frt_0.pdf (Stmt. of Vice Chair Nourse) (discussing Detroit Police Department practices and policies).

performance and bias statistics can substantiate a system's ongoing value and can alert an agency if there are changes that need attention. OMB Memorandum M-24-10 requires federal law enforcement agencies to implement continuous monitoring for FRT uses.

h. Restrictions on FRT Use and Logging

Policies should prohibit the use of FRT systems unless the agency has approved the system, the user, and the use. Policies and procedures should prohibit and establish consequences for unauthorized or improper use of a FRT biometric system (including examples discussed above, such as use based solely on constitutionally protected activity). Agencies should retain detailed internal logs of FRT system use for auditing and ensuring compliance with requirements.

i. Transparency Regarding FRT Use

As feasible, law enforcement agencies should adopt policies that require them to publicly disclose their use of FRT use, including details of the system in use, and the nature and purpose of the use. Agencies should also engage with community stakeholders about FRT and, to the extent possible, provide transparent responses about how they use FRT.

At the federal level, OMB Memorandum M-24-10 requires these practices for law enforcement uses of FRT. Federal law enforcement agencies must provide public transparency about FRT uses in an annual AI inventory. Agencies are also required to engage with stakeholders to obtain their input.

j. Data Management

Law enforcement agency policies on FRT should describe the collection, management, storage, and retention requirements for requests, probe images, and results. Policies should also describe security, privacy, recordkeeping, and audit requirements. Agencies should ensure compliance with any connected system policies, including FBI's Criminal Justice Information Services (CJIS) Security Policy.⁴⁷

k. Training

Documented and effective training is critical for the successful implementation of FRT. A policy should establish training requirements for personnel who will interact with FRT or its results. Training should address all aspects of an agency's FRT policy, including when investigators may run an FRT search and how the results may be used. Training should also provide background on the technology and its limitations, including how incorrect matches can occur and may disproportionately affect certain demographic groups.

⁴⁷ See generally BUREAU JUST. ASSISTANCE, FACE RECOGNITION POLICY DEVELOPMENT TEMPLATE (Dec. 2017), <https://bja.ojp.gov/sites/g/files/xyckuh186/files/Publications/Face-Recognition-Policy-Development-Template-508-compliant.pdf> (discussing best practices and applicable standards for security).

Training related to FRT should also cover the risks of human biases, such as automation and confirmation bias, when using FRT.⁴⁸ These biases can cause personnel to place undue weight on certain results, and strategies are available to mitigate the risk.⁴⁹ Demographic bias (explicit or implicit) can also affect human judgment and should be addressed in training. Research has shown that the innate ability to recognize faces varies widely and that people less reliably identify others from a different race.⁵⁰ This bias may compound with biases in FRT algorithms.⁵¹

Training for FRT use in law enforcement should be appropriate to a person's role and should convey the information necessary for responsibly submitting a probe image, analyzing results from an FRT system, and using the results in law enforcement activities. Roles for training that may be common across law enforcement agencies include:⁵²

- Facial Examiner: Compares a probe photo to candidate matches from an FRT system to develop possible investigative leads.
- Collector (Includes Investigators): Obtains probe images for use with an FRT system.
- Facial Reviewer (Includes Investigators): Reviews results of an FRT search adjudicated by a facial examiner.
- Supervisor: Oversees personnel involved with FRT and ensures compliance with law and policy.

Failure to properly train individuals who interact with FRT systems can increase the risk of potential errors at each step of the facial recognition process, which could ultimately impact individuals' privacy, civil rights, and civil liberties.

Automated License Plate Recognition

License Plate Readers (LPRs) are cameras with computer vision capabilities designed to detect and capture information from license plates within their field of view. Computer vision is an area of AI that can identify patterns and objects in images. In the criminal justice context, LPR cameras can be mounted on patrol vehicles to identify vehicles in connection with criminal investigations, including stolen vehicles, vehicles owned by wanted persons, vehicles involved in

⁴⁸ See Reva Schwartz et al., NAT. INST. STANDARDS & TECH. SPECIAL PUBL'N 1270, TOWARDS A STANDARD FOR IDENTIFYING AND MANAGING BIAS IN ARTIFICIAL INTELLIGENCE 26 (2022), <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1270.pdf>, (discussing systemic biases, statistical and computational biases, and human biases).

⁴⁹ See Samuel Peterson et al., RAND, *Finding a Broadly Practical Approach for Regulating the Use of Facial Recognition by Law Enforcement* 32 (Feb. 15, 2023), https://www.rand.org/pubs/research_reports/RRA2249-1.html (discussing sources of bias, including human decisionmaking).

⁵⁰ Jacqueline G. Cavazos et al., *Learning Context and the Other-Race Effect: Strategies for Improving Face Recognition*, 157 VISION RSCH., Apr. 2019, 169, 169–83, <https://nij.ojp.gov/library/publications/learning-context-and-other-race-effect-strategies-improving-face-recognition>.

⁵¹ See Peterson et al., *supra*, note 49.

⁵² See FACIAL IDENTIFICATION SCI. WORKING GRP., GUIDE FOR ROLE-BASED TRAINING IN FACIAL COMPARISON (2020), https://fiswg.org/fiswg_guide_for_role-based_training_in_facial_comparison_v1.0_20200717.pdf (listing different roles in a facial comparison environment including assessor, reviewer, examiner, manager, supervisor, collector, technical reviewer, and trainer).

an ongoing emergency like kidnapping, or vehicles being operated in an unlawful manner, enabling timely law enforcement action. LPR data also enables retrospectively tracking the movements of a vehicle of interest.

LPR cameras can be mounted in fixed locations to identify vehicles entering or exiting sensitive locations, including prisons or other government facilities, as well as known locations of criminal activity to identify potential criminals or enable facility security activities.

The collected license plate numbers can be cross-referenced against law enforcement databases to identify vehicles. That information can be used to identify vehicles that may require law enforcement action, such as stolen vehicles or suspect vehicles.⁵³ LPR cameras, including both fixed and vehicle mounted cameras, can assist in locating vehicles during Amber alerts, Ashanti alerts, silver alerts, or similar emergency situations.⁵⁴

Some LPR systems are operated by law enforcement agencies, using their own cameras typically on their own patrol vehicles. Other LPR systems are commercial services, with networks of participating cameras (e.g., at parking lots) that agencies can subscribe to.

LPR systems are in widespread use. Almost all large local law enforcement agencies have an LPR program, as do many smaller agencies.⁵⁵ Law enforcement agencies collect billions of LPR records per year.

Like biometrics used for identification, LPR systems can serve as an effective tool for monitoring vehicles in connection with criminal investigations. LPR successes include apprehending violent offenders and rescuing abducted children.⁵⁶ The use of LPR systems, however, comes with some risks. For example, LPR systems can misread plates or misidentify stolen vehicles.⁵⁷ In addition, as with all large stores of data, agencies may create privacy risks if they do not take the necessary steps to properly secure LPR data and dispose of it after it is no longer needed.

LPR systems, like FRT systems, must be accompanied by strong policy and procedural guardrails to ensure appropriate use.

⁵³ LPR systems in the United States are operated by both law enforcement agencies as well as commercial providers who operate the system as a service. As a result, the database referenced and system data retention policies are specific to each LPR system.

⁵⁴ See U.S. DEP'T JUST, BUREAU JUST. ASSISTANCE, FACT SHEET: NATIONAL ASHANTI ALERT NETWORK (2021), <https://bja.ojp.gov/sites/g/files/xyckuh186/files/media/document/National-Ashanti-Alert-Network-Fact-Sheet.pdf> ("Ashanti Alerts, once implemented, can provide rapid dissemination of information to law enforcement agencies, media, and the public about adults who have been reported missing, along with suspect information in cases of abduction.").

⁵⁵ U.S. DEP'T JUST., OFF. JUST. PROGS., BUREAU JUST. STATS., LOCAL POLICE DEPARTMENTS, 2013: EQUIPMENT AND TECHNOLOGY (2015), <https://bjs.ojp.gov/content/pub/pdf/lpd13et.pdf> ("An estimated 17% (about 2,000) of departments used automated license plate readers in 2013. This total included more than three-quarters of the departments serving 100,000 or more residents.").

⁵⁶ ANGEL DIAZ & RACHEL LEVINSON-WALDMAN, AUTOMATIC LICENSE PLATE READERS: LEGAL STATUS AND POLICY RECOMMENDATIONS FOR LAW ENFORCEMENT USE, BRENNAN CTR. (2020), <https://www.brennancenter.org/our-work/research-reports/automatic-license-plate-readers-legal-status-and-policy-recommendations>.

⁵⁷ *Id.*

LPR Risks and Mitigation

a. Procedural Risks

Agencies should establish clear policies that outline where, when, and for what purpose an LPR camera can be placed to enhance criminal justice operations. Policies and procedures should also clearly address retention periods for LPR records, as well as the circumstances, if any, in which they can be searched.

b. Accuracy and Reliability

Accuracy in LPR technology is critical to avoid false positives, misidentifications, and misdirected law enforcement actions. Agencies should prioritize data quality, with standards and procedures for identifying and correcting errors.

c. Safeguarding Privacy and Data Security

Protecting the security of LPR data is a priority. This includes implementing robust data encryption, access restrictions, and audit mechanisms to prevent unauthorized access and misuse. Agencies should adopt strong data security protocols to reinforce public confidence in LPR technology while safeguarding sensitive information.

d. LPR Data Sharing

Agencies should establish data sharing protocols if LPR data is shared with other agencies or third-party organizations. Data sharing agreements should include safeguards to protect privacy, and ensure all parties uphold the same standards of data security and ethical use.

e. Training

Proper training for using LPR systems is essential to ensure law enforcement personnel are aware of legal, ethical, and operational standards. This training should cover privacy protections, the importance of data accuracy, and the need to prevent bias in the use of LPR technology.

III. Forensic Analysis

Introduction

Forensic analysis of physical, digital, and multimedia evidence is central to criminal investigation and litigation. In recent years, uses of AI—including statistical techniques and machine learning¹—have accelerated the trend in forensic disciplines from subjective judgment in analyzing and interpreting forensic evidence toward more objective approaches.² This paradigm shift has potential to improve the reproducibility and accuracy of forensic methods, mitigate the possible human biases and examiner variation that may affect forensic analysis, and reinforce public trust in the criminal justice system.³ AI may also provide new capabilities and reduce the time and cost of analysis processes.

Professionals in forensic science have a responsibility to rigorously validate methods of analysis, and that responsibility remains with uses of AI.⁴ Properly designed validation studies can empirically demonstrate that a forensic method is reproducible and accurate, both in principle and

¹ As noted in the Introduction, definitions of AI vary substantially, with some broadly encompassing statistics and others emphasizing recent advances in machine learning. This chapter considers statistical and machine learning methods together because the opportunities and challenges for forensic analysis are broadly similar, and because the approach is consistent with the definition of artificial intelligence in OMB Memoranda M-24-10 and M-24-18. *See* R. SHUTE ET AL., WHAT FSSP LEADERS SHOULD KNOW ABOUT ARTIFICIAL INTELLIGENCE AND ITS APPLICATION TO FORENSIC SCIENCE 4 (Nat'l Inst. Just. 2023), <https://forensiccoe.org/private/65cfa81c601c4> (noting that automated and semi-automated systems used in forensic analysis require similar assessment, regardless of whether they are categorized as AI).

² *See* EXECUTIVE OFFICE OF THE PRESIDENT, PRESIDENT'S COUNCIL OF ADVISORS ON SCIENCE AND TECHNOLOGY, FORENSIC SCIENCE IN CRIMINAL COURTS: ENSURING SCIENTIFIC VALIDITY OF FEATURE-COMPARISON METHODS 46-54 (2016) ("PCAST Report"), https://obamawhitehouse.archives.gov/sites/default/files/microsites/ostp/PCAST/pcast_forensic_science_report_final.pdf (describing differences between subjective and objective forensic methods and explaining that objective methods are generally preferable). *See also* NATIONAL RESEARCH COUNCIL, STRENGTHENING FORENSIC SCIENCE IN THE UNITED STATES: A PATH FORWARD (2009), <https://nap.nationalacademies.org/catalog/12589/strengthening-forensic-science-in-the-united-states-a-path-forward>; JOSÉ ALMIRALL ET AL., AM. ASS'N ADVANCEMENT SCI., FORENSIC SCIENCE ASSESSMENTS: A QUALITY AND GAP ANALYSIS: FIRE INVESTIGATION (2017), https://www.aaas.org/sites/default/files/s3fs-public/reports/Fire%2520Investigation_0.pdf; WILLIAM THOMPSON ET AL., AM. ASS'N ADVANCEMENT SCI., FORENSIC SCIENCE ASSESSMENTS: A QUALITY AND GAP ANALYSIS: LATENT FINGERPRINT EXAMINATION (2017), https://www.aaas.org/sites/default/files/s3fs-public/reports/Latent%2520Fingerprint%2520Report%2520FINAL%2520_14.pdf; Jonathan J. Koehler et al., *The Scientific Reinvention of Forensic Science*, 120 PROCS. NAT'L ACAD. SCI. e2301840120 (2023), <https://doi.org/10.1073/pnas.2301840120>.

³ *See, e.g.*, Mark Barash et al., *Machine Learning Applications in Forensic DNA Profiling: A Critical Review*, 69 FORENSIC SCI. INT.: GENETICS 13 (2024).

⁴ *E.g.*, U.S. DEP'T JUST., CODE OF PROFESSIONAL RESPONSIBILITY FOR THE PRACTICE OF FORENSIC SCIENCE (2016) ("DOJ Code"), https://www.justice.gov/sites/default/files/code_of_professional_responsibility_for-the_practice_of_forensic_science_08242016.pdf ("5. Conduct research and forensic casework using the scientific method or agency best practices. Where validation tools are not known to exist or cannot be obtained, conduct internal or inter-laboratory validation tests in accordance with the quality management system in place."; "8. Conduct examinations that are fair, unbiased, and fit-for-purpose."; "10. Ensure interpretations, opinions, and conclusions are supported by sufficient data and minimize influences and biases for or against any party.").

as applied in a particular case.⁵ Validation can also enable forensic practitioners to convey valuable context for the results of forensic analysis, such as the likelihood of results occurring by chance or reflecting errors.

The use of AI in forensic science can add complexity to validation. Models built from data can have nuanced performance characteristics and incorporate demographic biases. Models can also be sensitive to subtle differences between the data used during development and the data encountered in real-world use, such as differences in how data is collected and prepared for forensic analysis. The implementation of AI systems can be complicated and proprietary, which can make it difficult to validate that a system implements a forensic method as intended. These possible sources of additional complexity reinforce the importance in forensic science of responsible practices for developing, validating, and revalidating methods of analysis.

Explainability is also important in forensic science.⁶ Practitioners have a responsibility to explain in a straightforward manner the data that they analyzed, the methods that they applied, and the interpretations, observations, and conclusions that resulted from applying the methods to the data.⁷ AI models may not be readily understandable by humans and may learn from correlations in data that are difficult to discern and not necessarily causal.⁸ Differences like these may affect how stakeholders in the criminal justice system are able to explain forensic processes that involve AI.

Expert oversight is another critical dimension of forensic science.⁹ AI should be a complement to the expertise of forensic practitioners, such as by recommending next steps for human consideration, checking human analysis, or providing a basis on which an expert might

⁵ See PCAST Report, *supra* note 2, at 42-43 (explaining the distinction between “foundational validity” and “validity as applied”).

⁶ The term “explainability” has a particular technical meaning in evaluating artificial intelligence systems. Here, this report uses it in the colloquial sense.

⁷ Professional codes in forensic science often address explanation of data, methods, and conclusions. *E.g.*, DOJ Code, *supra* note 4 (“12. Prepare reports and testify using clear and straightforward terminology, clearly distinguishing data from interpretations, opinions, and conclusions. Reports should disclose known limitations that are necessary to understand the significance of the findings.”; “15. Honestly communicate with all parties (the investigator, prosecutor, defense, and other expert witnesses) about all information relating to their analyses, when communications are permitted by law and agency practice.”). See EXECUTIVE OFFICE OF THE PRESIDENT, NAT’L SCI. & TECH. COUNCIL, STRENGTHENING THE FORENSIC SCIENCES 24 (2014), https://obamawhitehouse.archives.gov/sites/default/files/microsites/ostp/NSTC/forensic_science_may_2014.pdf (noting that, in a review of “more than 45 codes of ethics in use by various forensic science organizations,” a commonality was “the need to . . . provide clear and objective testimony”). Courts may also interpret legal standards for the admissibility of expert evidence, discussed further below, to require forms of explanation. *E.g.*, *Zenith Elecs. Corp. v. WH-TV Broad. Corp.*, 395 F.3d 416, 419 (7th Cir. 2005) (explaining that “[a]n expert must offer good reason to think that his approach produces an accurate estimate using professional methods, and this estimate must be testable,” such that “[s]omeone else using the same data and methods must be able to replicate the result,” and that if an expert “could or would not explain how his conclusions met [these] requirements, he was not entitled to give expert testimony”).

⁸ Leo Breiman, *Statistical Modeling: The Two Cultures (with Comments and a Rejoinder by the Author)*, 16 STAT. SCI. 199-231 (2001), <https://doi.org/10.1214/ss/1009213726>.

⁹ *E.g.*, DOJ Code, *supra* note 4 (“9. Make and retain contemporaneous, clear, complete, and accurate records of all examinations, tests, measurements, and conclusions, in sufficient detail to allow meaningful review and assessment by an independent professional proficient in the discipline.”).

arrive at a conclusion.¹⁰ Practitioners have an essential role in interpreting the output from AI systems, forming and explaining conclusions, and (when necessary and appropriate) offering expert testimony in litigation. Training for practitioners can also ensure that forensic methods involving AI are used properly and that results from analysis are accurately characterized.

The first part of this chapter describes how AI is currently being used in forensic analysis, as well as research suggesting areas of forensic science where AI may be used in the future. The breadth of use cases speaks to the immense potential of AI for forensic science. The second part of the chapter discusses challenges for AI use in forensic science and means of mitigating the risks associated with AI use. The chapter closes with a set of recommendations for forensic research and practice.

Current Uses of AI in Forensic Analysis

Several factors have limited the integration of AI into forensic analysis. These include the complex nature of forensic science, the standards for admissibility of forensic evidence in litigation, limited availability of high-quality and real (or realistic) data from forensic analysis settings, and resource constraints on forensic researchers and practitioners. As a result, real-world use of AI in forensics is presently limited outside of a few contexts.

As discussed in the Identification and Surveillance chapter, AI is widely used in biometric analysis, including fingerprint, palm print, iris, and face comparison. These methods complement other means of identification, and they are often used to retrieve candidate matches from a database for further forensic analysis or law enforcement investigation.

In DNA analysis, probabilistic genotyping can enable analysts to interpret complex samples that contain small amounts of DNA, mixed DNA, or damaged DNA that would be challenging for traditional analysis methods.¹¹ The statistical methods used for this type of forensic analysis typically estimate a likelihood ratio, comparing the probability of DNA observations given one proposition (e.g., that a defendant contributed genetic material to a sample) to the probability of observations given a stated alternative proposition (e.g., that an unknown person in a relevant population contributed to the sample).¹² The particular statistical analysis, population data, and

¹⁰ H. Swofford & C. Champod, *Implementation of Algorithms in Pattern & Impression Evidence: A Responsible and Practical Roadmap*, 3 FORENSIC SCI. INT'L: SYNERGY (2021), <https://doi.org/10.1016/j.fsisy.2021.100142>.

¹¹ See Michael D. Coble & Jo-Anne Bright, *Probabilistic Genotyping Software: An Overview*, 38 FORENSIC SCI. INT'L: GENETICS 219 (2019). As noted at the outset of the chapter, definitions of artificial intelligence differ. Some practitioners, for example, may not consider probabilistic genotyping to be a type of AI. See R. SHUTE ET AL., *supra* note 1, at 3 (describing disagreement about whether probabilistic genotyping is AI and noting it is at minimum part of an “automated system” that can exhibit bias).

¹² See, e.g., Peter Gill et al., *DNA Commission of the International Society of Forensic Genetics: Recommendations on the Evaluation of STR Typing Results That May Include Drop-out and/or Drop-in Using Probabilistic Methods*, 6 FORENSIC SCI. INT'L: GENETICS 679 (2012), <http://dx.doi.org/10.1016/j.fsigen.2012.06.002>; Hannah Kelly et al., *A Description of the Likelihood Ratios in the Probabilistic Genotyping Software STRmix*, 2 WIREs FORENSIC SCIENCE (2020), <https://doi.org/10.1002/wfs2.1377>.

assumptions underlying probabilistic genotyping methods can vary, leading to different results.¹³ The extent of validation research also differs by probabilistic genotyping method.¹⁴

Forensic genetic genealogy is another type of DNA analysis that can involve statistical models.¹⁵ These methods of analysis can enable the generation of leads by comparing a DNA sample to a large database of samples and estimating possible genealogical relationships.

Narcotics tracing is another area of forensics where AI is in use today. The Drug Enforcement Administration (DEA), for example, uses machine learning models to classify the geographic region of origin for samples of heroin and cocaine.¹⁶ The DEA's system was developed with authentic drug samples and can detect anomalies in analysis and low-confidence results. This program is valuable for understanding trends in drug trafficking, though the results are not presently used as evidence in court.

AI features are increasingly common in tools for forensic analysis of digital and multimedia evidence.¹⁷ Computer vision techniques, for example, can assist forensic examiners in searching large and unorganized collections of photos and videos for specific content, such as weapons, nudity, or violence, that may be helpful to an investigation. Natural language processing methods can similarly help identify files and communications that relate to particular topics. Machine translation can support examiners who are analyzing evidence that involves multiple languages. AI can be valuable, in some circumstances, for identifying and analyzing evidence that may have been created or modified by AI (e.g., “deepfake” images, videos, and audio).¹⁸ These uses of AI for analysis of digital and multimedia evidence are predominantly, for now, in support of investigative steps rather than expert conclusions offered in court.

¹³ See John Buckleton et al., *A Diagnosis of the Primary Difference Between EuroForMix and STRmix™*, 69 J. FORENSIC SCI. 40 (2024), <https://doi.org/10.1111/1556-4029.15387>; Peter Gill et al., *A Review of Probabilistic Genotyping Systems: EuroForMix, DNASTatX and STRmix™*, 12 GENES 1559 (2021), <https://www.mdpi.com/2073-4425/12/10/1559>; Susan A. Greenspoon et al., *A Tale of Two PG Systems: A Comparison of the Two Most Widely Used Probabilistic Genotyping Systems in the United States*, 69 J. FORENSIC SCI. 1840 (2024), <https://doi.org/10.1111/1556-4029.15571>.

¹⁴ See the discussion below for further detail about validation of probabilistic genotyping software.

¹⁵ See FORENSIC TECH. CTR. OF EXCELLENCE., AN INTRODUCTION TO FORENSIC GENETIC GENEALOGY TECHNOLOGY FOR FORENSIC SCIENCE SERVICE PROVIDERS (Nat'l Inst. of Just. 2022), <https://forensiccoe.org/private/66291221e66ec>; U.S. DEP'T JUST., INTERIM POLICY ON FORENSIC GENETIC GENEALOGICAL DNA ANALYSIS AND SEARCHING (2019), <https://www.justice.gov/olp/page/file/1204386/dl>.

¹⁶ U.S. DEP'T JUST., 2023 AI USE CASE INVENTORY, <https://www.justice.gov/open/file/1305831/dl>.

¹⁷ E.g., Inseyets, CELLEBRITE, <https://cellebrite.com/en/cellebrite-inseyets/>; Magnet Axiom, MAGNET FORENSICS, <https://www.magnetforensics.com/products/magnet-axiom/>. There is a broad range of possible additional applications of AI for digital and multimedia forensics. See Johannes Fährndrich et al., *Digital Forensics and Strong AI: A Structured Literature Review*, 46 FORENSIC SCI. INT'L: DIGITAL INVESTIGATIONS (2023), <https://doi.org/10.1016/j.fsidi.2023.301617>.

¹⁸ E.g., Video Authentication Software, MEDEX FORENSICS, <https://medexforensics.com/medex-platform/>. The performance of AI-based tools for analyzing possible AI-generated or AI-modified content varies significantly, and the value of these tools depends on the nature of the evidence and the investigative context.

Future Uses of AI in Forensic Analysis

Recent and ongoing research has demonstrated a range of additional possible future applications of AI in forensic analysis.

a. Pattern and Trace Evidence

AI may be well suited for assisting experts in the comparison and categorization of some types of pattern and trace evidence.¹⁹ Recent research has shown promise in applications of AI for assisting experts in analyzing toolmarks on bullets²⁰ and cartridges,²¹ impressions of footwear outsoles,²² fragments of glass,²³ traces of automotive paint,²⁴ and ignitable liquids,²⁵ among other types of evidence.

¹⁹ Trace evidence involves material that has been transferred between objects, people, or the environment.

²⁰ E.g., Eric Hare et al., *Algorithmic Approaches to Match Degraded Land Impressions*, 16 LAW, PROBABILITY & RISK 203-21 (2017), <https://doi.org/10.1093/lpr/mgx018>; Susan Vanderplas et al., *Comparison of Three Similarity Scores for Bullet LEA Matching*, 308 FORENSIC SCI. INT'L (2020), <https://www.sciencedirect.com/science/article/pii/S0379073820300293>; Pattranit Pisantanaroj et al., *Automated Firearm Classification From Bullet Markings Using Deep Learning*, 8 IEEE ACCESS 78236 (2020), <https://ieeexplore.ieee.org/document/9076037>.

²¹ E.g., Xiao Hui Tai & William F. Eddy, *A Fully Automatic Method for Comparing Cartridge Case Images*, 63 J. FORENSIC SCI. 440 (2018), <https://onlinelibrary.wiley.com/doi/full/10.1111/1556-4029.13577>; Joseph Roth et al., *Learning-based Ballistic Breech Face Impression Image Matching*, 2015 IEEE 7TH INTERNATIONAL CONFERENCE ON BIOMETRICS THEORY, APPLICATIONS AND SYSTEMS (2015), <https://ieeexplore.ieee.org/document/7358774>.

²² E.g., Soyoung Park & Alicia Carriquiry, *An Algorithm to Compare Two-dimensional Footwear Outsole Images Using Maximum Cliques and Speeded-up Robust Feature*, 13 STAT. ANALYSIS & DATA MINING: THE ASA DATA SCI. J. 188 (2020), <https://onlinelibrary.wiley.com/doi/full/10.1002/sam.11449>; Hana Lee et al., *An Automated Alignment Algorithm for Identification of the Source of Footwear Impressions with Common Class Characteristics*, 17 STAT. ANALYSIS & DATA MINING: THE ASA DATA SCI. J. (2024), <https://dl.acm.org/doi/10.1002/sam.11659>; Gautham Venkatasubramanian et al., *Quantitative Evaluation of Footwear Evidence: Initial Workflow for an End-to-end System*, 66 J. OF FORENSIC SCI. 2232 (2021), <https://onlinelibrary.wiley.com/doi/full/10.1111/1556-4029.14802>; Gautham Venkatasubramanian et al., *Comparing footwear impressions that are close non-matches using correlation-based approaches*, 66 J. FORENSIC SCI. 2232 (2021), <https://onlinelibrary.wiley.com/doi/10.1111/1556-4029.14658>; Moonsoo Jang & Soyoung Park, *A Finely Tuned Deep Transfer Learning Algorithm to Compare Outsole Images*, 16 STAT. ANALYSIS & DATA MINING: THE ASA DATA SCI. J. 511 (2023), <https://onlinelibrary.wiley.com/doi/abs/10.1002/sam.11636>; Zhijian Wen et al., *Shoeprint Image Retrieval and Crime Scene Shoeprint Image Linking by Using Convolutional Neural Network and Normalized Cross Correlation*, 63 SCI. & JUST. 439 (2023), <https://doi.org/10.1016/j.scijus.2023.04.014>; Samia Shafique et al., *CriSp: Leveraging Tread Depth Maps for Enhanced Crime-Scene Shoeprint Matching*, ARXIV (2024), <https://arxiv.org/abs/2404.16972>.

²³ E.g., Omer Kaspi et al., *Toward Developing Techniques—Agnostic Machine Learning Classification Models for Forensically Relevant Glass Fragments*, 63 J. CHEMICAL INFO. & MODELING 87 (2022), <https://pubs.acs.org/doi/full/10.1021/acs.jcim.2c01362>; Grzegorz Zadora, *Glass Analysis for Forensic Purposes—A Comparison of Classification Methods*, 21 J. CHEMOMETRICS 174 (2007), <https://analyticalsciencejournals.onlinelibrary.wiley.com/doi/10.1002/cem.1030>.

²⁴ E.g., Francis Kwofie et al., *Transmission Infrared Microscopy and Machine Learning Applied to the Forensic Examination of Original Automotive Paint*, 76 APPLIED SPECTROSCOPY 118 (2021), <https://journals.sagepub.com/doi/10.1177/00037028211057574>; George P. Affadu-Danful et al., *Raman Spectroscopy to Enhance Investigative Lead Information in Automotive Clearcoats*, 77 APPLIED SPECTROSCOPY 1064 (2023), <https://journals.sagepub.com/doi/full/10.1177/00037028231186838>.

²⁵ E.g., Christian Bogdal et al., *Recognition of Gasoline in Fire Debris Using Machine Learning*, 331 FORENSIC SCI. INT'L (2022), <https://doi.org/10.1016/j.forsciint.2021.111146>; Michael E. Sigman et al., *Validation of Ground Truth Fire Debris Classification by Supervised Machine Learning*, 26 FORENSIC CHEMISTRY (2021), <https://doi.org/10.1016/j.fore.2021.100358>.

b. Drug Evidence

Analysis of seized drug samples is another promising area. Research has demonstrated that AI can assist with classifying fentanyl analogs and related compounds,²⁶ marijuana varieties,²⁷ and novel psychoactive substances.²⁸ These approaches generally combine established chemistry methods for identifying components of compounds, such as mass spectrometry, with machine learning methods to analyze the components and make categorizations.

c. Forensic Medicine, Pathology, and Anthropology

AI may assist with assessing injuries and injury mechanics. Analyzing a photograph of a bruise, for example, may enable estimating the date of the injury.²⁹

AI may also be able to supplement expert analysis of human remains.³⁰ Recent publications show that it is feasible to estimate sex,³¹ age,³² and population affinity³³ from images and 2D and 3D computed tomography scans of skeletal and dental remains. AI may assist with identifying decedents, including through post-mortem iris recognition³⁴ and association of remains with

²⁶ E.g., Travon Cooman et al., *Evaluation and Classification of Fentanyl-related Compounds Using EC-SERS and Machine Learning*, 68 J. FORENSIC SCI. 1520 (2023), <https://onlinelibrary.wiley.com/doi/10.1111/1556-4029.15285>; Phillip Koshute et al., *Machine Learning Model for Detecting Fentanyl Analogs from Mass Spectra*, 27 FORENSIC CHEMISTRY 100379 (2022), <https://www.sciencedirect.com/science/article/abs/pii/S2468170921000758>.

²⁷ E.g., Austin McDaniel et al., *Toward the Identification of Marijuana Varieties by Headspace Chemical Forensics*, 11 FORENSIC CHEMISTRY 23-31 (2018), <https://doi.org/10.1016/j.forc.2018.08.004>.

²⁸ Swee Liang Wong et al., *Screening Unknown Novel Psychoactive Substances Using GC-MS based Machine Learning*, 34 FORENSIC CHEMISTRY 100499 (2023), <https://www.sciencedirect.com/science/article/abs/pii/S2468170923000358?via%3Dihub>.

²⁹ E.g., Jhonatan Tirado & David Mauricio, *Bruise Dating Using Deep Learning*, 66 J. FORENSIC SCI. 336 (2020), <https://pmc.ncbi.nlm.nih.gov/articles/PMC7821214/>.

³⁰ See Laurent Tournois et al., *Artificial Intelligence in the Practice of Forensic Medicine: A Scoping Review*, 138 INT'L J. LEGAL MED. 1023 (2023), <https://pmc.ncbi.nlm.nih.gov/articles/PMC11003914/>; Nicola Galante et al., *Applications of Artificial Intelligence in Forensic Sciences: Current Potential Benefits, Limitations and Perspectives*, 137 INT'L J. LEGAL MED. 445 (2022), <https://doi.org/10.1007/s00414-022-02928-5>; Andrej Thurzo et al., *Use of Advanced Artificial Intelligence in Forensic Medicine, Forensic Anthropology and Clinical Anatomy*, 9 HEALTHCARE 1545 (2021), <https://www.mdpi.com/2227-9032/9/11/1545>; Micayla C. Spiros & Sherry Nakhaeizadeh, *We Think There's Been a Glitch: Artificial Intelligence and Machine Learning in Forensic Anthropology*, 7 FORENSIC ANTHROPOLOGY (2024), <https://journals.upress.ufl.edu/fa/article/view/2827>.

³¹ Javier Venema et al., *Employing Deep Learning for Sex Estimation of Adult Individuals using 2D Images of the Humerus*, 35 NEURAL COMPUTING & APPLICATIONS 5987 (2023), <https://link.springer.com/article/10.1007/s00521-022-07981-0>; Tomoyuki Seo et al., *Sex Estimation Using Skull Silhouette Images from Postmortem Computed Tomography by Deep Learning*, 14 SCI. REPS. 22689 (2024), <https://www.nature.com/articles/s41598-024-74703-y>; Yongjie Cao et al., *Use of Deep Learning in Forensic Sex Estimation of Virtual Pelvic Models from the Han Population*, 7 FORENSIC SCIS. RES. 540 (2022), <https://academic.oup.com/fsr/article/7/3/540/6987953>.

³² Juan Carlos Gámez-Granados et al., *Automating the Decision Making Process of Todd's Age Estimation Method from the Pubic Symphysis with Explainable Machine Learning*, 612 INFORMATION SCI. 514 (2022), <https://www.sciencedirect.com/science/article/pii/S0020025522010301?via%3Dihub>; NICHOLAS P. HERRMANN ET AL., INVESTIGATION OF SUBADULT DENTAL AGE-AT-DEATH ESTIMATION USING TRANSITION ANALYSIS AND MACHINE LEARNING METHODS (Off. Just. Progs. 2023), <https://www.ojp.gov/pdffiles1/nij/grants/306558.pdf>.

³³ David Navega et al., *AncesTrees: Ancestry Estimation with Randomized Decision Trees*, 129 INT'L J. LEGAL MED. 1145, 1145-1153 (2015); G. Richard Scott et al., *rASUDAS: A New Web-Based Application for Estimating Ancestry from Tooth Morphology*, 1 FORENSIC ANTHROPOLOGY 18, 18-31 (2018).

³⁴ Aidan Boyd et al., *Post-Mortem Iris Recognition—A Survey and Assessment of the State of the Art*, 8 IEEE ACCESS 136570, 136570 (2020).

radiographs, dental records, photos, and 3D scans.³⁵ Where remains are incomplete, AI may help experts impute missing measurements.³⁶ Machine learning methods may also assist experts in categorizing and annotating images of remains,³⁷ and may be able to suggest certain causes of death, such as head trauma³⁸ or drowning.³⁹

The ability of AI tools to effectively support forensic analysis in these ways will be contingent on the quality of photographs or imagery, among other factors.

d. Forensic Biology

AI has the promise of providing experts with new capabilities in probabilistic genetic analysis. These methods may improve the accuracy of genetic sequencing,⁴⁰ overcome gaps in degraded samples,⁴¹ support inference of the number of contributors to a DNA mixture,⁴² and enable predictions about physical appearance based on genetic information.⁴³

³⁵ David C. Cornett et al., *Effects of Postmortem Decomposition on Face Recognition*, IEEE 10TH INT'L CONF. BIOMETRICS THEORY, APPLICATIONS, AND SYS. (2019), <https://ieeexplore.ieee.org/document/9185971>; A. Valsecchi et al., *Skeleton-ID: AI-driven Human Identification*, IEEE CONF. ARTIFICIAL INTEL. (2023), <https://ieeexplore.ieee.org/abstract/document/10195123>.

³⁶ Jinyong Pang & Xiaoming Liu, *Evaluation of Missing Data Imputation Methods for Human Osteometric Measurements*, 181 AM. J. BIOLOGICAL ANTHROPOLOGY 666 (2023), <https://onlinelibrary.wiley.com/doi/10.1002/ajpa.24787>.

³⁷ AUDRIS MOCKUS & DAWNIE WOLFE STEADMAN, ICPUTRD: IMAGE CLOUD PLATFORM FOR USE IN TAGGING AND RESEARCH ON DECOMPOSITION (Off. Just. Progs. 2020), <https://www.ojp.gov/pdffiles1/nij/grants/255312.pdf>; Sara Mousavi et al., *Machine-Assisted Annotation of Forensic Imagery*, 2019 IEEE INT'L CONF. ON IMAGE PROCESSING 1595 (2019), <https://ieeexplore.ieee.org/document/8803068>; Jack Garland et al., *Identifying Gross Post-mortem Organ Images Using a Pre-trained Convolutional Neural Network*, 66 J. FORENSIC SCI. 630 (2021), <https://onlinelibrary.wiley.com/doi/10.1111/1556-4029.14608>.

³⁸ Jack Garland et al., *Identifying Fatal Head Injuries on Postmortem Computed Tomography Using Convolutional Neural Network/Deep Learning: A Feasibility Study*, 65 J. FORENSIC SCI., 2019 (2020), <https://onlinelibrary.wiley.com/doi/full/10.1111/1556-4029.14502>.

³⁹ Noriyasu Homma et al., *A Deep Learning Aided Drowning Diagnosis for Forensic Investigations Using Post-Mortem Lung CT Images*, 42ND ANNUAL INT'L CONF. OF THE IEEE ENG'G MED. & BIOLOGY SOC'Y (2020), <https://ieeexplore.ieee.org/document/9175731>.

⁴⁰ August E. Woerner et al., *Reducing Noise and Stutter in Short Tandem Repeat Loci with Unique Molecular Identifiers*, 51 FORENSIC SCI. INT'L: GENETICS 102459 (2021), [https://www.fsigenetics.com/article/S1872-4973\(20\)30231-3/abstract](https://www.fsigenetics.com/article/S1872-4973(20)30231-3/abstract); Michael A. Marciano et al., *A Hybrid Approach to Increase the Informedness of CE-based Data Using Locus-specific Thresholding and Machine Learning*, 35 FORENSIC SCI. INT'L: GENETICS 26 (2018), [https://www.fsigenetics.com/article/S1872-4973\(17\)30313-7/abstract](https://www.fsigenetics.com/article/S1872-4973(17)30313-7/abstract).

⁴¹ Meng Huang et al., *A Machine Learning Approach for Missing Persons Cases with High Genotyping Errors*, 13 FRONTIERS GENETICS 1 (2022) <https://www.frontiersin.org/journals/genetics/articles/10.3389/fgene.2022.971242>.

⁴² Michael A. Marciano & Jonathan D. Adelman, *PACE: Probabilistic Assessment for Contributor Estimation—A Machine Learning-based Assessment of the Number of Contributors in DNA Mixtures*, 27 FORENSIC SCI. INT'L: GENETICS 82 (2017), <https://pubmed.ncbi.nlm.nih.gov/28040630/>; Hamdah Alotaibi et al., *DNA Profiling: An Investigation of Six Machine Learning Algorithms for Estimating the Number of Contributors in DNA Mixtures*, 12 INT'L J. ADVANCED COMPUT. SCI. & APPLICATIONS 1 (2021), <http://dx.doi.org/10.14569/IJACSA.2021.0121115>.

⁴³ Maria-Alexandra Katsara et al., *Evaluation of Supervised Machine-learning Methods for Predicting Appearance Traits from DNA*, 53 FORENSIC SCI. INT'L: GENETICS 102507 (2021), <https://doi.org/10.1016/j.fsigen.2021.102507>.

Analysis of other types of biological evidence may also benefit from AI. Applying AI to microscopy images, for example, may assist experts with locating sperm cells in sexual assault evidence.⁴⁴

Advances in forensic biology could be particularly valuable for investigations of violent crimes, where a perpetrator or victim may be more likely to leave biological evidence.

e. Forensic Toxicology

In toxicological analysis, AI may enable experts to identify unknown compounds.⁴⁵ Automated analysis may also be able to assess the toxicity of compounds, supporting treatment decisions made by medical professionals.⁴⁶

f. Crime Scene Analysis

Complex crime scenes may involve thousands of photographs with nuanced attributes for forensic examiners to review and document. AI can support experts' work by automatically categorizing crime scene photographs based on visible items of interest, such as drugs and weapons.⁴⁷ AI may also be able to improve the quality of crime scene imagery for expert analysis, such as for scenes that are underwater.⁴⁸ In some instances, automated tools may be able to aid analysts in interpreting evidence within crime scene photographs, such as categorizing the potential mechanism that caused a blood spatter pattern.⁴⁹

Challenges for AI in Forensic Analysis

The preceding chapter on Identification and Surveillance and the later chapter on Risk Assessment describe several challenges that are equally applicable to uses of AI in forensic analysis. In short, high-quality and representative data, rigorous and independent testing for performance and biases, ongoing monitoring, and established policies and oversight are all critical

⁴⁴ Raffael Golomingi et al., *Sperm Hunting on Optical Microscope Slides for Forensic Analysis with Deep Convolutional Networks – A Feasibility Study*, 56 FORENSIC SCI. INT'L: GENETICS 102602 (2022), <https://www.sciencedirect.com/science/article/pii/S1872497321001393>.

⁴⁵ Toshali D. Wankhade et al., *Artificial Intelligence in Forensic Medicine and Toxicology: The Future of Forensic Medicine*, 14 CUREUS (2022), <https://doi.org/10.7759/cureus.28376>; Zhoumeng Lin & Wei-Chun Chou, *Machine Learning and Artificial Intelligence in Toxicological Sciences*, 189 TOXICOL. SCI. 7 (2022), <https://doi.org/10.1093/toxsci/kfac075>.

⁴⁶ See, e.g., Thi Tuyet Van Tran et al., *Artificial Intelligence in Drug Toxicity Prediction: Recent Advances, Challenges, and Future Perspectives*, 63 J. CHEM. INFO. & MODELING 2628 (2023), <https://doi.org/10.1021/acs.jcim.3c00200>.

⁴⁷ Joshua Abraham et al., *Automatically Classifying Crime Scene Images Using Machine Learning Methodologies*, 39 FORENSIC SCI. INT'L: DIGITAL INVESTIGATION 301273 (2021), <https://doi.org/10.1016/j.fsidi.2021.301273>; Amaljith Sreekumar et al., *Weapons and Related Object Classification in Digital Forensic Using Machine Learning*, 14TH INT'L CONF. ON COMPUTING COMM'NS & NETWORKING TECHS. 1 (2023), <https://ieeexplore.ieee.org/abstract/document/10307988>.

⁴⁸ Rosella Paba et al., *Optimizing Underwater Visual Records for Crime Scene Investigations in Water with Clear to Reduced Visibility*, 6 FORENSIC SCI. INT'L: SYNERGY 100329 (2023), <https://doi.org/10.1016/j.fsisyn.2023.100329>.

⁴⁹ Yu Liu et al., *Automatic Classification of Bloodstain Patterns Caused by Gunshot and Blunt Impact at Various Distances*, 65 J. FORENSIC SCI. 729 (2020), <https://doi.org/10.1111/1556-4029.14262>.

for uses of AI in criminal justice, including in support of forensic science. This section notes several challenges that could be particularly acute for the use of AI in forensic analysis.

a. Datasets

A large volume of high-quality data is essential for developing and evaluating AI uses in forensic analysis. If there is insufficient data, or if data has errors or gaps, that can undermine the potential value of AI for forensic science and create risks of harm from mistaken conclusions. Data that is properly distributed across a range of demographics and scenarios is also critical. This helps prevent bias in AI uses, leading to more consistent and fair forensic analysis.

Forensic analysis often depends on specialized data, such as data collected with dedicated equipment or from samples that are not readily available.⁵⁰ Preparing forensic datasets that are adequate for AI can also be expensive and labor intensive. Coordination by forensic experts worldwide to provide highly accurate and representative data may be essential for some uses of AI in forensic science.⁵¹

Respecting privacy is another important consideration for forensic science datasets. The data can be personal and sensitive, necessitating appropriate safeguards.

b. Validation

As noted earlier, analysis methods used in forensic science should be carefully studied for validity in principle, validity as applied in particular cases, and ongoing validity. Validation is essential for understanding the reliability and limitations of evidence analysis. Experts should be able to effectively characterize the accuracy and errors of AI used in forensic science, including the possibility of demographic disparities. When a forensic method involves AI, performance testing, bias testing, and continuous monitoring are all important steps for validation. Review of the implementation of an AI system, including the source code, can also be appropriate in some circumstances. At the federal level, OMB Memorandum M-24-10 requires testing and monitoring for certain use cases of AI used in support of forensic analysis. OMB Memorandum M-24-18 further requires use of independent testing data to the extent practicable, prohibits contractual restrictions on agency disclosure of testing methods and results, and encourages consideration of open-source AI implementations.

As an example, the path to widespread adoption of probabilistic genotyping highlights the importance and value of rigorous validation for new uses of AI in forensic science. When

⁵⁰ Because forensic analysis can involve specialized data, practitioners should be especially attentive to the risks of overfitting, where a model may learn correlations in training data that do not generalize to data in real-world uses. Overfitting can degrade performance and lead to unexpected model behaviors. At the federal level, OMB Memorandum M-24-10 specifically cautions agencies to consider the risks of overfitting.

⁵¹ See Toshali D. Wankhade et al., *Artificial Intelligence in Forensic Medicine and Toxicology: The Future of Forensic Medicine*, 14 CUREUS 1, 4 (2022), <https://doi.org/10.7759/cureus.28376>. This type of data may be less necessary for certain types of analysis for digital and multimedia evidence, where the possible range of formats is determined by readily available software.

probabilistic genotyping began regularly appearing in forensic DNA analysis in the mid-2010s,⁵² the President’s Council of Advisors on Science and Technology (PCAST) expressed optimism that these methods could improve on prior approaches to analyzing samples with complex DNA mixtures, while also emphasizing the need for further validation.⁵³ PCAST concluded that studies had demonstrated validity in principle, but only for some mixtures of genetic material, and that additional steps were necessary to demonstrate validity as applied by laboratories. PCAST also expressed concern that validation studies had predominantly been carried out in collaboration with the vendors of probabilistic genotyping software. Scholars criticized probabilistic genotyping software as opaque and at risk of implementation errors, and they called attention to restrictive licensing agreements and trade secret protections that could inhibit independent validation.⁵⁴

The forensic science community responded with a concerted effort to further substantiate probabilistic genotyping methods through independent and collaborative research. Laboratories worldwide coordinated in studies using the data available to each, reinforcing that particular methods could, in particular circumstances, have validity as they are deployed.⁵⁵ Vendors of probabilistic genotyping software improved access to their tools and source code for review by opposing legal teams.⁵⁶ While there is ongoing debate about the extent of validation for probabilistic genotyping—a 2021 draft comprehensive report by NIST concluded that gaps remain,⁵⁷ and a 2024 workshop by the National Academies reflected ongoing stakeholder concerns⁵⁸—these steps have been positive and important.

⁵² See SCIENTIFIC WORKING GROUP ON DNA ANALYSIS METHODS, GUIDELINES FOR THE VALIDATION OF PROBABILISTIC GENOTYPING SYSTEMS (2015), https://www.swgdam.org/files/ugd/4344b0_22776006b67c4a32a5ffc04fe3b56515.pdf.

⁵³ PCAST Report, *supra* note 2.

⁵⁴ E.g., Rebecca Wexler, *Life, Liberty, and Trade Secrets: Intellectual Property in the Criminal Justice System*, 70 STAN. L. REV. 1343, 1343 (2018), <https://review.law.stanford.edu/wp-content/uploads/sites/3/2018/06/70-Stan.-L.-Rev.-1343.pdf>; Andrea Roth, *Machine Testimony*, 126 YALE L.J. 1972, 1972 (2017), https://www.yalelawjournal.org/pdf/RothFinal_c4o97on1.pdf. See also Rediet Abebe et al., *Adversarial Scrutiny of Evidentiary Statistical Software*, 2022 ACM CONF. ON FAIRNESS, ACCOUNTABILITY, & TRANSPARENCY 1733 (2022), <https://doi.org/10.1145/3531146.3533228>.

⁵⁵ Safia Boodoosingh et al., *An Inter-laboratory Comparison of Probabilistic Genotyping Parameters and Evaluation of Performance on DNA Mixtures from Different Laboratories*, 71 FORENSIC SCI. INT’L: GENETICS 103046 (2024), <https://doi.org/10.1016/j.fsigen.2024.103046>; John M. Butler et al., *NIST Interlaboratory Studies Involving DNA Mixtures (MIX05 and MIX13): Variation Observed and Lessons Learned*, 37 FORENSIC SCI. INT’L: GENETICS 81 (2018), <https://doi.org/10.1016/j.fsigen.2018.07.024>; John M. Butler et al., *DNA Mixture Interpretation: A NIST Scientific Foundation Review*, NAT’L INST. STANDARDS & TECH. (2021), <https://nvlpubs.nist.gov/nistpubs/ir/2021/NIST.IR.8351-draft.pdf> (surveying validation studies of probabilistic genotyping software). See Jo-Anne Bright et al., *Internal Validation of STRmix™ – A Multi Laboratory Response to PCAST*, 34 FORENSIC SCI. INT’L: GENETICS 11 (2024), <https://doi.org/10.1016/j.fsigen.2018.01.003>.

⁵⁶ E.g., TrueAllele® Source Code Access, CYBERGENETICS, <https://www.cybgcn.com/support/code-access/>; Access to STRmix™ Software By Defence Legal Teams, STRMIX, <https://www.strmix.com/assets/STRmix/STRmix-PDFs/Access-to-STRmix-Software-by-Defence-Legal-teams-March-2022-v2.pdf>.

⁵⁷ John M. Butler et al., *DNA Mixture Interpretation: A NIST Scientific Foundation Review*, NAT’L INST. STANDARDS & TECH. (2021), <https://nvlpubs.nist.gov/nistpubs/ir/2021/NIST.IR.8351-draft.pdf>.

⁵⁸ NAT’L ACADS. SCIS., ENG’G & MED., *Law Enforcement Use of Probabilistic Genotyping, Forensic DNA Phenotyping, and Forensic Investigative Genetic Genealogy Technologies: Proceedings of a Workshop* (2024), <https://doi.org/10.17226/27887>.

These developments with probabilistic genotyping highlight the value of proactively carrying out large-scale interlaboratory validation for new uses of AI in forensic analysis. They also suggest that forms of access to AI implementation, such as models or source code, may be appropriate for validation.

c. Explainability

As noted earlier, it is important for forensics experts to be able to explain the analytical methods they apply and how they obtain results. The types of AI models that are used in forensic science today are generally interpretable, such that an expert could describe how inputs combine to arrive at an output.⁵⁹ It is foreseeable, though, that forensic science may begin to involve AI models that can be more difficult to understand, such as deep neural networks. These models could possibly be less persuasive in court and could undermine public confidence in forensic analysis. When considering models that are not interpretable, forensic practitioners should explore the feasibility of using explainability methods that can provide some understanding of model behavior.⁶⁰ Practitioners should also carefully consider possible tradeoffs between interpretability and accuracy in developing AI models for forensic science.

d. Human Oversight of AI in Forensic Analysis

In forensic science, it is common to have a second person review a practitioner's work, sometimes referred to as a "technical review." Similarly, when using AI, it is important to maintain human oversight of analysis and results to ensure that the AI was consistently applied and identify possible irregularities. Human involvement is also important because, if forensic analysis involving AI will be the basis for evidence in court, a human expert must explain the AI use and interpret the results.⁶¹

Including a human review element in the forensic analysis process comes with risks, however. Human review can introduce human biases, such as confirmation bias with respect to a subject or automation bias with respect to the reliability of analysis.⁶² Training for forensic

⁵⁹ See Brandon L. Garrett & Cynthia Rudin, *Interpretable Algorithmic Forensics*, 120 PROCS. NAT'L ACAD. SCI. (Oct. 2023), <https://doi.org/10.1073/pnas.2301842120>.

⁶⁰ See Louise Kelly et al., *Explainable Artificial Intelligence for Digital Forensics: Opportunities, Challenges and a Drug Testing Case Study*, DIGITAL FORENSIC SCIENCE (B. Suresh Kumar Shetty & Pavanchand Shetty eds., 2020), <https://www.intechopen.com/chapters/73078>; Stuart W. Hall & Amin Sakzad, *Explainable Artificial Intelligence for Digital Forensics*, 4 WIREs FORENSIC SCI. (2021), <https://wires.onlinelibrary.wiley.com/doi/abs/10.1002/wfs2.1434>; Marthe S. Veldhuis et al., *Explainable Artificial Intelligence in Forensics: Realistic Explanations for Number of Contributor Predictions of DNA Profiles*, 56 FORENSIC SCI. INT'L: GENETICS 102632 (2022), <https://doi.org/10.1016/j.fsigen.2021.102632>. See generally Alejandro Barredo Arrieta et al., *Explainable Artificial Intelligence (XAI): Concepts, Taxonomies, Opportunities and Challenges Toward Responsible AI*, 58 INFO. FUSION 82 (2020), <https://doi.org/10.1016/j.inffus.2019.12.012>.

⁶¹ See Andrea Roth, *Machine Testimony*, 126 YALE L.J. 1972 (2017), https://www.yalelawjournal.org/pdf/RothFinal_c4o97on1.pdf.

⁶² See S. M. Kassir et al., *The Forensic Confirmation Bias: Problems, Perspectives, and Proposed Solutions*, 2 J. APPLIED RSCH. MEMORY & COGNITION 42 (2013), <https://doi.org/10.1016/j.jarmac.2013.01.001>. In forensic analysis of digital and multimedia evidence, for example, examiners often rely on tools that automate elements of analysis. The automation may not be robust against changes in the format of evidence, which could result from routine software updates. Automation bias could cause an examiner to miss relevant information if a tool is unsuccessful at

practitioners who use AI should address the risks of biases affecting analysis. Forensic analysis procedures should also minimize these risks, such as by minimizing the availability of irrelevant information to forensic examiners.⁶³

Recommendations

a. Policies for Use of AI in Forensic Analysis

Forensic science service providers (FSSPs), such as law enforcement laboratories, should have clear and documented policies regarding the use of AI in forensic analysis. These policies should address the types of AI that may be used, the circumstances in which AI may be used, governance requirements, and limitations. FSSPs should consider using the NIST AI Risk Management Framework, NIST's Trustworthy and Responsible Artificial Intelligence Resource Center, OMB Memoranda M-24-10 and M-24-18, and other appropriate AI guidance from federal agencies to develop guidance and establish governance programs.

Consistent with accepted standards for forensic analysis, human review and interpretation of AI outputs should remain standard procedure in forensic science applications.

The output of an AI system should not be the sole basis for conclusions in forensic analysis. A qualified examiner should interpret the output and apply their expert judgment to form conclusions.

FSSPs should consider potential risks from using AI in forensic analysis and should design and implement processes to mitigate those risks prior to using AI. AI impact assessments may be a valuable framework for these considerations.

Policies should establish rigorous validation requirements to ensure that AI uses are reliable, both as they are developed and as they are deployed by the FSSP. Appropriate validation will often involve pre-deployment testing for performance and demographic biases, using data and contexts that are representative of real-world use, as well as post-deployment monitoring.

FSSP policies should address AI interpretability and explainability and should set a general preference for interpretable models when they can meet operational needs.

Regular audits of AI use can ensure that examiners are following required procedures.

b. Procurement of AI Capabilities for Forensic Analysis

FSSPs should only procure tools that have a demonstrated acceptable level of accuracy. FSSPs should verify that the data used to build an AI model is high quality and representative of the FSSP's intended real-world use.

extracting or parsing evidence. A possible mitigation for this risk would be to ensure that, if a tool encounters an error, it is clearly conveyed to examiners.

⁶³ Appropriate steps to minimize risks of human biases for forensic analysis involving AI may be similar or identical to appropriate steps for minimizing bias risks for analysis that does not involve AI.

FSSPs should require information from vendors about intended uses, training data and methods, validation, potential limitations, and potential biases of products. Experts should carefully review vendor disclosures and relevant additional information, such as peer-reviewed publications and reviews by other laboratories, for alignment with forensic analysis objectives.

Partnerships with independent AI researchers may be beneficial for FSSPs in validating AI products. Where possible, FSSPs should avoid restrictive licensing agreements and other possible barriers to collaboration with third-party experts.

FSSPs should ensure that biases in AI systems and uses, including for sex, race, color, disability, and age, have been adequately evaluated and mitigated.

c. Datasets

When possible, FSSPs should choose AI capabilities that have been trained on large, high-quality, and representative datasets. Producing these training datasets may require cooperation among practitioners and forensic experts worldwide, accounting for differences across jurisdictions that could affect performance or introduce biases. If an AI use depends on data that can differ in relevant ways across jurisdictions, FSSPs should consider appropriately supplementing training data with datasets from their jurisdictions.

When evaluating an AI system, FSSPs should also evaluate the training data where possible to ensure that the data is accurate, complete, and representative of the intended deployment context.

FSSPs should ensure that the data used to validate an AI system is separate from the data used to build the system. Reuse of training data in testing can lead to misunderstanding of the system's performance and biases.

d. Training and Education

Forensic practitioners who use or interact with AI should be appropriately trained about the AI, including its design, intended use, performance, biases, and limitations. Practitioners should also receive appropriate training about applicable policies and required procedures, as well as how to mitigate the possible human biases associated with use of the AI.

FSSPs can benefit each other by sharing information about their experiences with AI and the policies and procedures that they have implemented.

Forensic leaders should stay current on emerging AI tools and monitor advancements in forensic applications to take advantage of the technology and make changes to existing use of AI tools as necessary.

IV. Predictive Policing

Introduction

Predictive policing is the use of quantitative analytical methods to identify times, places, and individuals likely to be associated with criminal activity.¹ Predictive policing models do not predict that specific crimes will occur, but rather estimate a general likelihood of crime. While law enforcement agencies use predictive policing tools for purposes like informing the allocation of officers, existing tools neither recommend nor evaluate responses to their output. Those responses could extend beyond traditional policing methods and determine the efficacy of strategies incorporating predictive policing.

Law enforcement agencies have used predictive policing tools since at least the 1990s.² Researchers studied “hot spots” policing with the Minneapolis Police Department and found crime reductions in patrolled areas,³ leading to further studies looking at approaches identifying crime “hot spots” to reduce drug and violent crime.⁴ The Department of Justice’s National Institute of Justice (NIJ) helped bring analysis tools to a wider range of U.S. law enforcement agencies by funding the development of CrimeStat software beginning in 1997.⁵ By 2008, when approximately 90% of law enforcement agencies were already using some form of “hot spots” policing,⁶ the NIJ joined with the Department of Justice’s Bureau of Justice Assistance to fund research into new predictive models that would turn “hot spots” mapping into a forward-looking tool for crime forecasting.⁷

¹ WALTER L. PERRY ET AL., RAND CORP., PREDICTIVE POLICING: THE ROLE OF CRIME FORECASTING IN LAW ENFORCEMENT OPERATIONS 1–2 (2013), https://www.rand.org/content/dam/rand/pubs/research_reports/RR200/RR233/RAND_RR233.pdf. The term “predictive policing” is something of a misnomer because it is not a policing strategy, like proactive policing, problem-oriented policing, or crime prevention through environmental design. Instead, the outputs of predictive policing inform the development and implementation of policing or other strategies to reduce crime. However, this chapter uses the term “predictive policing” given its widespread acceptance. The chapter also references predictive policing algorithms, models, and tools to emphasize that predictive analytics are separate from the responses, which may include various forms of policing as well as community approaches to reducing crime.

² Lawrence W. Sherman & David Weisburd, *General Deterrent Effects of Police Patrol in Crime “Hot Spots”: A Randomized, Controlled Trial*, 12 JUST. Q. 625, 643–46 (1995), <http://dx.doi.org/10.1080/07418829500096221>.

³ *Id.*

⁴ See, e.g., Anthony A. Braga et al., *The Effects of Hot Spots Policing on Crime: An Updated Systematic Review and Meta-Analysis*, 31 JUST. Q. 633, 634 (2014).

⁵ NED LEVINE, THE DEVELOPMENT OF A SPATIAL ANALYSIS TOOLKIT FOR USE IN A METROPOLITAN CRIME INCIDENT GEOGRAPHIC INFORMATION SYSTEM (1999) (Final Rep. to the Nat’l Inst. Just., Award No. 97-IJ-CX-0040), <https://www.ojp.gov/pdffiles1/nij/grants/179282.pdf>.

⁶ See POLICE EXEC. RSCH. F., VIOLENT CRIME IN AMERICA: WHAT WE KNOW ABOUT HOT SPOTS ENFORCEMENT, 3 (2008), https://www.policeforum.org/assets/docs/Critical_Issues_Series/violent%20crime%20in%20america%20-%20what%20we%20know%20about%20hot%20spots%20enforcement%202008.pdf (“nearly 9 out of 10 agencies use hot spots enforcement efforts directed either at larger hotspots areas like neighborhoods, smaller hot spot places like intersections, or both”).

⁷ Joel Hunt, *From Crime Mapping to Crime Forecasting: The Evolution of Place-Based Policing*, NIJ J., no. 281, Nov. 2019, <https://www.ojp.gov/pdffiles1/nij/252036.pdf>.

While even early versions of predictive policing fall near this report’s broad definition of artificial intelligence, predictive policing today uses more advanced statistical and machine learning methods and incorporates greater volumes and new types of data. Predictive policing tools can be grouped into two categories. “Location-based” or “place-based” tools attempt to forecast locations where crime is likely to cluster.⁸ Location-based tools also try to identify the types of crimes likely to occur and when they are likely to occur. “Person-based” predictive policing attempts to identify individuals who are more likely to commit crimes or be victims of crime in the future, and it has many similarities to risk assessment (*see* chapter V).⁹ The literature on predictive policing focuses overwhelmingly on “street crimes,” and this chapter follows that research.

Predictive policing tools can help law enforcement agencies, other public services, and community-based organizations promote public safety, efficiency, and transparency by informing decisions on how to allocate limited resources. Nevertheless, these tools do not currently prescribe the actions that should be taken at the places or with the people identified—actions such as additional enforcement, place-based problem-solving, diversion programs, job training, education, or environmental design interventions—or predict the effects of those actions.

Law enforcement use of predictive policing also raises significant risks, including the potential to create or entrench disparities. The data underlying predictive policing models may have significant gaps and errors, and it may reflect historical and human biases. Use of predictive policing models based on that data may also result in unintended, unjust outcomes, such as over- or under-policing of certain individuals and communities.

This chapter begins with descriptions of uses and risks of both place-based and person-based predictive policing tools, and it concludes with recommendations for law enforcement agencies to safely and effectively deploy these tools in a way that enhances efficiency and accuracy while protecting civil liberties, civil rights, and privacy.

Uses of Predictive Policing

a. Place-based

Place-based or “hot spots” policing involves “focusing limited resources on a small number of high-activity crime places.”¹⁰ Also called “spatial models,” promising analyses from Minneapolis and elsewhere suggested opportunities for these methods to assist in using policing resources more effectively and spurred extensive adoption by the 2000s.¹¹ During the 2010s, place-based predictive policing systems drew criticism about their accuracy, limitations, and reliance on historical crime.¹² As discussed below, predictive policing tools often rely heavily on historic crime

⁸ *See* PERRY ET AL., *supra* note 1, at 8–9 (describing a similar taxonomy of predictive policing methods).

⁹ *See id.*

¹⁰ Anthony A. Braga, *Effects of Hot Spots Policing on Crime*, CAMPBELL SYSTEMATIC REVIEWS 1, 4 (2007).

¹¹ SHERMAN & WEISBURD, *supra* note 2, at 643–46; POLICE EXEC. RSCH. F., *supra* note 6.

¹² *See, e.g.,* Lyria Bennett Moses & Janet Chan, *Algorithmic Prediction in Policing: Assumptions, Evaluation, and Accountability*, 28 POLICING & SOCIETY 806, 809–13 (2018) (examining underlying assumptions in predictive policing).

data, which can be problematic because the data may entrench historically discriminatory policing patterns. Criticisms also arose surrounding the immediate value and actionability of some systems' output for police departments.¹³

In 2016, the NIJ conducted the Real-Time Crime Forecasting Challenge to test the effectiveness and efficiency of predictive spatial models on crime data from Portland, Oregon.¹⁴ Although results from the Challenge should not be generalized beyond the context of that single city, comparing the performance of a wide range of algorithms there indicated that both simple and sophisticated spatial models can offer similar predictive accuracy.¹⁵ This finding supports the principle that for producing better outcomes, the specific tool or algorithm used for prediction may matter less than the selection and implementation of responses to high-risk areas.¹⁶ For example, community engagement in the development of responses (e.g., additional police presence versus social services) could determine whether those responses are effective.¹⁷

Since the 2010s, place-based predictive policing strategies have continued to evolve, as agencies incorporate evidence from past iterations and work with communities to respond with a wide range of interventions. Agencies are developing tools in-house, customizing them to fit into their workflows, and collaborating with other public services and community-based organizations.¹⁸ Strategies that currently make use of predictive policing include proactive policing and “hot spot” policing.¹⁹ These strategies include a range of specific approaches to prevent crime and focus limited resources geographically.²⁰

methods); DAVID ROBINSON & LOGAN KOEPKE, UPTURN, STUCK IN A PATTERN: EARLY EVIDENCE ON “PREDICTIVE POLICING” AND CIVIL RIGHTS 3–5 (2016), https://www.upturn.org/static/reports/2016/stuck-in-a-pattern/files/Upturn_-_Stuck_In_a_Pattern_v.1.01.pdf (discussing general limitations of predictive policing).

¹³ See Aaron Sankin & Surya Mattu, *Predictive Policing Software Terrible at Predicting Crimes*, THE MARKUP (Oct. 2, 2023, 10:00 AM), <https://themarkup.org/prediction-bias/2023/10/02/predictive-policing-software-terrible-at-predicting-crimes> (finding low success rate of a predictive policing tool).

¹⁴ *Real-Time Crime Forecasting Challenge*, NAT'L INST. JUST. (July 13, 2016), <https://nij.ojp.gov/funding/real-time-crime-forecasting-challenge> (archived content).

¹⁵ See YongJei Lee et al., *A Theory-Driven Algorithm for Real-Time Crime Hot Spot Forecasting*, 23 POLICE Q. 174, 194–96 (2020) (analyzing data from Portland and Cincinnati and arguing that authors' simple model in Microsoft Excel “has demonstrated similar levels of efficiency and accuracy, with lower economic or fiscal investments and greater transparency than other more expensive, commercial models”).

¹⁶ See, e.g., NAT'L ACAD. OF SCI., ENG'G, & MED., LAW ENFORCEMENT USES OF PREDICTIVE POLICING APPROACHES (Nov. 2024) (“NAS Proceeding 2024”), <https://nap.nationalacademies.org/catalog/28037/law-enforcement-use-of-person-based-predictive-policing-approaches-proceedings>.

¹⁷ NAS PROCEEDING 2024, *supra* note 16.

¹⁸ See Tim Lau, *Predictive Policing Explained*, BRENNAN CTR FOR JUST. (Apr. 1, 2020), <https://www.brennancenter.org/our-work/research-reports/predictive-policing-explained> (“[T]he NYPD developed its own in-house predictive policing algorithms and started to use them in 2013”).

¹⁹ See generally Albert Meijer & Martijn Wessels, *Predictive Policing: Review of Benefits and Drawbacks*, 42 INT'L J. PUB. ADMIN. 1031, 1033–34 (2019) (discussing predictive policing and conventional policing methods).

²⁰ See *Practice Profile: Hot Spots Policing*, NAT'L INST. JUST. <https://crimesolutions.ojp.gov/ratedpractices/hot-spots-policing#1-0> (“Hot spots policing strategies focus on small geographic areas or places, usually in urban settings, where crime is concentrated. Through hot spots policing strategies, law enforcement agencies can focus limited resources in areas where crime is most likely to occur.”)

Recent place-based predictive policing programs include Place-Based Investigations of Violent Offender Territories (PIVOT) in Cincinnati²¹ and Data-Informed Community Engagement (DICE) in cities such as Dallas, Kansas City, Newark, New Orleans, and St. Louis.²² PIVOT uses “hot spot” mapping combined with problem-oriented policing strategies to identify and mitigate factors that facilitate violence,²³ while DICE pairs risk terrain modeling with community engagement to address crime through place-based interventions.²⁴

b. Person-based

A person-based approach attempts to identify those most at risk for committing future crimes or being a victim of crime, either by using factors associated with individuals known to law enforcement to generate risk scores, or by identifying connections between individuals who may be linked to past crimes. This approach typically generates a list of the highest-risk individuals in the jurisdiction as a whole or within a given geographic area, such as a patrol zone. With respect to violent crime, research consistently shows that it is disproportionately concentrated among small numbers of individuals, groups, and locations at the highest risk for violence.²⁵

In the mid-1990s, the NIJ funded an early study of one form of predictive policing, focused deterrence, in which law enforcement efforts prioritized individuals disproportionately responsible for crime.²⁶ In the past decade, predictive policing programs like Chicago’s Strategic Subjects List built on earlier person-focused deterrence approaches and other initiatives, such as gang databases.²⁷ That generation of policing programs established the use of algorithms to analyze data on individuals’ criminal histories, social networks, and other risk factors in an attempt to identify those most likely to be involved in violent crime as perpetrators or victims. By the end of the

²¹ *Place-Based Investigations of Violent Offender Territories (PIVOT)*, CITY OF CINCINNATI, <https://www.cincinnati-oh.gov/police/community-involvement/pivot/>.

²² Learning Community, PUB. SAFETY COLLABORATIVE COUNCIL, <https://www.diceforpublicsafety.org/learning-community.html>.

²³ See Tamara D. Madensen et al., *Place-Based Investigations to Disrupt Crime Place Networks*, THE POLICE CHIEF MAG., Apr. 2017, at 14–15, https://www.policechiefmagazine.org/wp-content/uploads/PoliceChief_April2017_F_WEB.pdf (describing development and implementation of PIVOT).

²⁴ Sarah Minster, *The Data-Informed Community Engagement (DICE) Approach to Public Safety Turns Analytics into Action*, NAT’L LEAGUE OF CITIES (June 18, 2024), <https://www.nlc.org/article/2024/06/18/the-data-informed-community-engagement-dice-approach-to-public-safety-turns-analytics-into-action>.

²⁵ *Violent Crime Reduction Roadmap: Working Together to Build Safer Communities*, Action 2, BUREAU JUST. ASSISTANCE, <https://bj.a.ojp.gov/violent-crime-reduction-roadmap/action-2#0-3> (“Research consistently shows that violence disproportionately concentrates among small number of individuals, groups and locations at the highest risk for violence.”)

²⁶ Braga et al. 2001, *Problem-Oriented Policing, Deterrence, and Youth Violence: An Evaluation of Boston's Operation Ceasefire*, J. RSCH. CRIME & DELINQ. 195, 198 (2001), <https://journals.sagepub.com/doi/abs/10.1177/0022427801038003001> (discussing the Boston Gun Project and Operation Ceasefire).

²⁷ See generally Jessica Saunders et al., *Predictions Put into Practice: A Quasi-Experimental Evaluation of Chicago's Predictive Policing Pilot*, 12 J. EXPERIMENTAL CRIMINOLOGY 347 (2016), <https://link.springer.com/article/10.1007/s11292-016-9272-0>.

2010s, however, many of the programs from this era had halted following concerns regarding efficacy, bias, and civil rights.²⁸

Person-based strategies predicting risk for victimization rather than criminal offending are also being used in cases of child abuse²⁹ and gender-based violence.³⁰ These tools raise some similar concerns to offender-centered predictive policing tools.³¹

Risks of Predictive Policing

a. Data and Model Output Quality

Although traditional law enforcement interventions have relied on officers' knowledge, judgments, and personal experience in their communities, predictive modeling relies on collected quantitative data as input. Predictive models depend on historic crime data, such as calls for service, crime incidents, and arrests.³² As with all data, these datasets are imperfect.³³ They may underrepresent underreported crimes, such as intimate partner violence and sexual assault.³⁴ They may also reflect past policing patterns, which can affect the certainty of the model's output and disproportionately affect vulnerable communities.³⁵

Predictive models may also underrepresent the actual incidence of crime in locations where individuals are less likely to report crimes to police.³⁶ Some scholars have pointed to historical

²⁸ See, e.g., Annie Sweeney & Jeremy Gerner, *For Years Chicago Police Rated the Risk of Tens of Thousands Being Caught Up in Violence. That Controversial Effort Has Quietly Been Ended.*, CHI. TRIB. (Jan. 25, 2020, 2:55 AM), <https://www.chicagotribune.com/2020/01/24/for-years-chicago-police-rated-the-risk-of-tens-of-thousands-being-caught-up-in-violence-that-controversial-effort-has-quietly-been-ended/> (discussing a findings in a report issued by the City Inspector General's Office analyzing the use of Chicago's Strategic Subjects List).

²⁹ See, e.g., ALLEGHENY CNTY. DEP'T HUM. SERVS., SUMMARIZING RECENT RESEARCH ON PREDICTIVE RISK MODELS IN CHILD WELFARE 1 (2024), <https://www.alleghenycountyanalytics.us/wp-content/uploads/2024/05/24-ACDHS-04-Predictive-Risk-Algorithms.pdf> (discussing the Allegheny Family Screening Tool, "an algorithm designed to assist child welfare call screening caseworkers in their assessment of general protective service referrals regarding potential child maltreatment.").

³⁰ See, e.g., ETICAS FOUND., THE EXTERNAL AUDIT OF THE VIOGÉN SYSTEM (2022), <https://eticasfoundation.org/wp-content/uploads/2024/07/ETICAS-FND-The-External-Audit-of-the-VioGen-System-1.pdf> (examining VioGén, a gender-violence risk assessment tool used by the Spanish Ministry of Interior).

³¹ ALLEGHENY CNTY. DEP'T HUM. SERVS., *supra* note 29.

³² See, e.g., Jeffrey Brantingham et al., *Does Predictive Policing Lead to Biased Arrests? Results From a Randomized Controlled Trial*, STATISTICS AND PUBLIC POLICY 5(1), 1–6. <https://doi.org/10.1080/2330443X.2018.1438940> (evaluating bias of predictive algorithms used for police patrol using arrest data).

³³ *Id.* at 5 (discussing limitations of arrest data).

³⁴ See U.S. DEP'T JUST., FRAMEWORK FOR PROSECUTORS TO STRENGTHEN OUR NATIONAL RESPONSE TO SEXUAL ASSAULT AND DOMESTIC VIOLENCE INVOLVING ADULT VICTIMS 2 n. 8 (May 2024), <https://www.justice.gov/ovw/media/1352371/dl?inline> ("Sexual assaults and domestic violence are, in large part, underreported, under-investigated and under-prosecuted").

³⁵ Some observers have pointed to racial disparities in policing decisions in such instances. See, e.g., AM. C. L. UNION N.J., SELECTIVE POLICING RACIALLY DISPARATE ENFORCEMENT OF LOW-LEVEL OFFENSES IN NEW JERSEY 8 (Dec. 2015), https://www.aclu-nj.org/sites/default/files/field_documents/2015_12_21_aclunj_select_enf.pdf.

³⁶ The Bureau of Justice Statistics reported that in 2022, only 41.5% of violent crimes and 31.8% of property crimes were reported to police. ALEXANDRA THOMPSON & SUSANNAH N. TAPP, BUREAU OF JUST. STATS., U.S. DEP'T JUST., NCJ 307089, CRIMINAL VICTIMIZATION, 2022 6 tbl.4 (2023), <https://bjs.ojp.gov/document/cv22.pdf>.

public health statistics as another way to estimate the underlying incidence of crime³⁷ just like the Drug Enforcement Administration (DEA) has used public-facing CDC drug overdose data to evaluate and to more effectively address the drug threat in communities across the United States.³⁸ Looking further afield, other crime prediction models have incorporated a wide variety of data sourced from outside law enforcement and unconnected to past criminal conduct, including “elevation; zoning and water area coverage; density of hospitals, fire departments, transportation points such as bus stops and subway entrances, and schools,” as well as variables such as the day of the week and the weather.³⁹ Many datasets, including data on crime incidents, may also contain some uncertainty because data features such as the time of a crime or the location of a crime may have a large range or be inaccurately reported.

Police officers typically generate data each time they conduct a stop, write a report, or do anything else that leaves a data trail. Some police actions are subject to discretion, and the discretion of individual officers thus plays a role in shaping the input data.⁴⁰ As the next section discusses, the use of AI in predictive policing runs the risk of introducing or exacerbating disparities in this input data.

b. Civil Rights

Perhaps the most frequent criticism of predictive policing is that it has the potential to reproduce or even amplify biases embedded in historical crime data. The risk is that historical crime data will train models to predict crime in ways that magnify biases.⁴¹ For example, place-based models can drive increased police presence in an area with greater historical crime data. The increased presence can lead to more law enforcement activity in that area, which can result in even more officers assigned to the area. Ultimately, the models become “increasingly confident that the locations most likely to experience further criminal activity are exactly the locations they had previously believed to be high in crime.”⁴² The result may be disparate outcomes, by race or other

³⁷ See Kristian Lum & William Isaac, *To Predict and Serve?*, SIGNIFICANCE 15, 16–17 (Oct. 2016) <https://rss.onlinelibrary.wiley.com/doi/pdf/10.1111/j.1740-9713.2016.00960.x> (comparing police records of drug crimes with public health survey of illegal drug use, which the authors argue is a better “ground truth” for incidence of drug crimes).

³⁸ See, e.g., Press Release, Drug Enforcement Admin., DEA Launches New Initiative to Combat Drug-Related Violence and Overdoses in Communities Across America (Feb. 7, 2022), <https://www.dea.gov/press-releases/2022/02/07/dea-launches-new-initiative-combat-drug-related-violence-and-overdoses-0> (“DEA initiated a data-driven approach using national crime statistics and CDC data to identify hot spots of drug-related violence and overdose deaths across the country, in order to devote its law enforcement resources where they will have the most impact: the communities where criminal drug networks are causing most harm.”); DEA Administrator Anne Milgram Remarks as Delivered Press Conference (Dec. 16, 2021) <https://www.dea.gov/sites/default/files/2021-12/DEA%20Administrator%20Anne%20Milgram%20Remarks%20as%20Delivered-Dec%2016%202021.pdf> (discussing CDC data on drug overdose deaths and criminal drug network activities).

³⁹ Jerry H. Ratcliffe et al., *The Philadelphia Predictive Policing Experiment*, 17 J. EXPERIMENTAL CRIMINOLOGY 15, 21 (2021) (describing possible inputs for predictive policing model formerly known as HunchLab); see also Resourcerouter Frequently Asked Questions, SoundThinking, <https://www.soundthinking.com/faqs/resourcerouter-faqs/> (“We supplement data modeling with non-crime data . . . Typical examples include seasonality, time of month, day of the week, time of day, holidays, upcoming events, weather, and locations of liquor establishments.”).

⁴⁰ See Lum & Isaac, *supra* note 38, at 16 (describing “feedback loop” in predictive model trained on police records that overrepresent drug crimes in heavily patrolled areas).

⁴¹ *Id.* at 19.

⁴² *Id.* at 16.

demographic characteristics, for residents of that area even where rates of crime may be comparable to other neighborhoods.

Person-based models may also be trained on data reflecting underlying disparities,⁴³ which may adversely impact individuals who receive increased attention and law enforcement interactions. In one person-based program, “officers are instructed to focus their attention on the highest point-value individuals ... [who] are subjected to heightened surveillance, and, therefore, are more likely to be stopped, thus further increasing their point value.”⁴⁴

This situation can create feedback loops in which “targeting of certain areas or certain races creates the impression of higher crime rates in those areas, which then justifies continued police presence there.”⁴⁵ Beyond the impact on residents and individuals, feedback loops can erode public trust in police by increasing police presence, and thus enforcement activities, in communities that are already distrustful of law enforcement.

c. Privacy

Individuals that a predictive policing tool assesses to be more likely to engage in or become a victim of crime may be subject to additional law enforcement scrutiny, raising concerns about privacy, surveillance, and harassment. Even given stringent criteria for inclusion of an individual’s data in a predictive policing system, the absence of oversight and monitoring for adherence to these criteria can mean that some individuals experience law enforcement surveillance or interactions without substantiated links to criminal activity. Precedent exists for oversight concerns with criminal justice databases.⁴⁶

Developers of predictive policing models may also choose to augment police data by drawing from additional datasets.⁴⁷ However, merging other datasets may also raise privacy concerns, such as revealing sensitive details of individuals or communities.⁴⁸

⁴³ See Andrew Guthrie Ferguson, *Policing Predictive Policing*, 94 WASH. U. L. REV. 1109, 1148–49 (2017) (“The result has been to justify disproportionate minority contacts and the collection of minority names in databases. These actions then feed a confirmation feedback loop that equates those currently in the system with those who need to be policed by the system. Essentially, high-crime areas or high-value suspects might only be considered ‘high’ because police already have data about those areas or people.”).

⁴⁴ Sarah Brayne & Angèle Christin, *Technologies of Crime Prediction: The Reception of Algorithms in Policing and Criminal Courts*, *Social Problems* 1, 13 (2020) <https://doi.org/10.1093/socpro/spaa004>; see also Sarah Brayne, *Big Data Surveillance: The Case of Policing*, 977, 987 (2017), <https://doi.org/10.1177/0003122417725865> (“An individual having a high point value is predictive of future police contact and that police contact further increases the individual’s point value.”).

⁴⁵ Ferguson, *supra* note 43 at 1153.

⁴⁶ See, e.g., CAL. STATE AUDITOR, THE CALGANG CRIMINAL INTELLIGENCE SYSTEM 1, 31–32 (2016), <https://www.auditor.ca.gov/pdfs/reports/2015-130.pdf> (finding California gang database’s “oversight structure does not ensure that user agencies collect and maintain criminal intelligence in a manner that does not invade individuals’ privacy rights,” and that some agencies lacked adequate support for including individuals in database).

⁴⁷ See Ratcliffe et al., *supra* note 3938.

⁴⁸ E.g. Arvind Narayanan & Vitaly Shmatikov, *Robust De-anonymization of Large Sparse Datasets*, <https://ieeexplore.ieee.org/document/4531148> (connecting Netflix and IMDB data to learn potentially sensitive video-watching habits of individuals).

d. Interpreting and Explaining

As models become more complex, they gain the ability to represent more intricate relationships and sophisticated patterns in the data.⁴⁹ This increased flexibility, however, often comes with increased difficulty for humans to interpret and explain the internal workings of these complex models.⁵⁰

In predictive policing, more complex models are likely to be marginally more accurate than simpler ones,⁵¹ especially when predicting rarer crimes. However, the diminishing marginal improvements that come with more sophisticated algorithms may not justify their use over simpler ones, given their higher cost, greater complexity, and decreased interpretability. A less complex predictive policing tool that is more interpretable may better align with law enforcement and community needs and values than a highly complex “black box” tool. An interpretable tool can allow for greater accountability and oversight by providing concrete explanations for a model’s predictions and may even increase transparency into the factors that influenced decisions.

Recommendations

Thoughtful deployment of predictive policing tools can not only mitigate risks but can also ensure that AI use for predictive policing offers accuracy and efficiency benefits for law enforcement. Entities that have deployed or are considering deploying predictive policing tools should consider the following recommendations.

a. Assess Goals of the Predictive Policing Tool with the Community

When considering a predictive policing tool, a critical question is the particular community’s goals that have driven consideration of the tool. After identifying objectives, a community should outline how it will measure the success of the tool or alternatives in addressing those goals. These initial steps should engage public agencies, relevant community-based organizations, and other stakeholders in the community, such that agencies can determine how to align their use of tools to community needs and expectations, address concerns and potential risks, and establish consensus for the implementation process. Failure to effectively engage with impacted communities can undermine public trust. Law enforcement agencies weighing the adoption of predictive policing tools should therefore engage with community members and their representatives (e.g., county commissioners or city council members) regarding goals and measures of success.

b. Assess the Need for a Predictive Policing Tool and Possible Alternatives

Law enforcement agencies should assess the likelihood a predictive policing tool will address the community’s goals, the shortcomings the tool might have, alternatives to the tool

⁴⁹ See, e.g., PERRY ET AL., *supra* note 1, at 36 (comparing “simple methods,” such as regression analyses, to “black box models,” which can model “extremely complicated relationships”).

⁵⁰ Whereas simple models “are usually directly interpretable by a person,” more complex “black box models” are not. PERRY ET AL., *supra* note 1, at 36.

⁵¹ See, e.g., PERRY ET AL., *supra* note 1, at xix (“Although there is usually a correlation between the complexity of a model and its predictive power, increases in predictive power have tended to show diminishing returns.”).

(including non-AI alternatives), and ways the tool might address or exacerbate concerns. Along with these considerations, law enforcement agencies should evaluate whether they have the resources, personnel, and technical expertise necessary for the proper use and monitoring of these types of systems in the short and long term.

This assessment—and the full lifecycle of a predictive policing tool’s deployment—should include community engagement. Law enforcement agencies should seek to educate community members and their representatives on how a considered predictive policing tool works, the benefits of the tool, the short- and long-term costs of the tool, the risks associated with the tool (e.g., civil rights concerns), and plans for monitoring and mitigating risks. When consulting with stakeholders, law enforcement agencies should provide the public meaningful opportunities for participation by using plain language, considering language access needs, and communicating with those most likely to be negatively impacted by these tools.

c. Assess Which Data to Use to Train Models and Ensure Data Is As Accurate As Possible

Law enforcement agencies should carefully consider which data to include in a predictive policing model. To mitigate the risk of feedback loops, agencies should keep the model inputs focused on types of crime relevant to their goals. Additionally, because models should be continually updated with new data, agencies should consider filtering out data that reflect actions taken by police in response to previous model outputs in order to avoid feedback loops.

Agencies should also ensure a central role for humans in choosing which data to use as inputs, deciding which metrics to use for evaluating accuracy, screening the model’s outputs with a crime analyst’s eye for contextual knowledge and historical factors, and determining which practices to implement in response to predictions. Although humans have the potential to introduce their own biases, they are also essential for ensuring that predictive policing tools are being used according to law and regulation, as well as relevant policies. By being transparent with the public about these choices, especially with respect to the selection of input data for predictive tools, agencies can ensure accountability and build public trust.

Predictive models may have a veneer of neutrality that can be reinforced by automation bias.⁵² To counteract this, user interfaces for predictive policing systems should include constraints, cues, notifications, and embedded content that reinforce appropriate use of the model, especially with regard to critically evaluating model outputs. Training users to think critically about the limits to algorithmic objectivity and error rates must be a top priority in any law enforcement organization using predictive models.⁵³

⁵² “Automation bias refers to undue deference to automated systems by human actors that disregard contradictory information from other sources or do not (thoroughly) search for additional information.” Saar Alon-Barkat & Madaline Busuioc, *Human–AI Interactions in Public Sector Decision Making: “Automation Bias” and “Selective Adherence” to Algorithmic Advice*, 33 J. PUB. ADMIN. RSCH. & THEORY 153, 155 (2023).

⁵³ See ALEXANDER BABUTA & MARION OSWALD, ROYAL UNITED SERVS. INST. FOR DEF. & SEC. STUD., DATA ANALYTICS AND ALGORITHMIC BIAS IN POLICING 15 (2019), https://static.rusi.org/20190916_data_analytics_and_algorithmic_bias_in_policing_web_0.pdf (“Adequate training

Although elimination of all errors from a large dataset may be infeasible, law enforcement agencies have a responsibility to ensure that the data they collect on both incidents and individuals are as accurate as possible. Routine checks of police and other criminal justice data for accuracy and ongoing revisions to improve data quality can mitigate the risks associated with data errors for both the police and community members. To the extent feasible, organizations using predictive policing tools should establish channels for affected community members to determine what data is stored about them, learn about decisions made based on that data, and petition for omission or corrections of errors.

d. Test, Measure, Validate, and Reevaluate, Including Independently and in a Real-world Context

For validation purposes, agencies should use data collected before and after deployment. Pre-deployment data are recommended because they isolate the system's ability to forecast without any interference from changes in response to the predictions. Collecting post-deployment data is recommended because they might help in measuring the efficiency or accuracy of the model over time when responses are changing and potentially detecting, deterring, or preventing crimes in the forecasted areas.⁵⁴

After validating the model, agencies should also conduct operational and field testing to better understand the impact of the model combined with interventions. Even accurate and well-validated models could produce unintended consequences when implemented in practice, so agencies should have mechanisms in place to identify the broad impact of their prediction-based interventions—including the possibility of bias and discrimination—and make any appropriate changes based on continual evaluation.⁵⁵

The impact of a particular predictive policing approach on crime rates is a commonly applied metric, but other factors can yield a more complete picture of the impact, including: community-defined measures of success; citizen commendations or complaints; use of excessive or unwarranted force; traffic stops and field interviews; citations, fines, and fees; victim satisfaction surveys; citizen surveys about factors such as fear of crime and satisfaction with and trust in police; number of tips received; rates of victim and witness cooperation; officer health, wellbeing, retention, and job satisfaction; and geographic residency of new applicants to the police force.⁵⁶ To identify disparate treatment or impacts, evaluations of models should include data

focused on cognitive bias and fair decision-making would appear essential to ensure officers are able to consistently achieve the correct balance.”).

⁵⁴ “Deployment of prediction boxes to the field and delivery of policing dosage to those boxes is expected to suppress some fraction of crime in those locations. Predictive accuracy should therefore decline in response to directed patrol.” G.O. Mohler et al., *Randomized Controlled Field Trials of Predictive Policing*, 110 J. AM. STAT. ASS’N 1399, 1404 (2015).

⁵⁵ See, e.g., NAT’L INST. STANDARDS & TECH., ARTIFICIAL INTELLIGENCE RISK MANAGEMENT FRAMEWORK (AI RMF 1.0), at 33 tbl.4 (2023), <https://doi.org/10.6028/NIST.AI.100-1>.

⁵⁶ Research exploring the accuracy of place-based predictive policing algorithms commonly relies on the Prediction Accuracy Index (PAI) or the Prediction Efficiency Index* (PEI*). See Spencer Chainey et al., *The Utility of Hotspot Mapping for Predicting Spatial Patterns of Crime*, 21 SEC. J. 4 (2008) (proposing PAI); Veronica M. White et al., *A Discussion of Current Crime Forecasting Indices and an Improvement to the Prediction Efficiency Index for*

disaggregated by demographic factors such as race, national origin, age, gender, disability, and other characteristics.

Testing and validation are necessary regardless of whether an agency develops a predictive policing tool in-house or acquires the tool from a third party. Part of testing and evaluation should include completing an impact assessment for the predictive policing tool before putting it into use.

Agencies should adopt a mechanism for objective third-party auditing of models and source code to proactively address concerns about model accuracy, reliability, and potential for bias and discrimination. Agencies should routinely audit police databases to mitigate the costs of data errors for both users of the system and community members. Databases used in person-based predictive policing, for example, should be routinely audited to ensure that individuals meet the criteria for being included in the database.

Models should be reevaluated periodically based on the policing patterns and crime rates that result from their use, ensuring the models are fair; equitable, including minimizing feedback loops; and aligned with public safety priorities. If seeking to procure a predictive policing tool from a third party, agencies should require regular evaluation and auditing as part of the vendor contract.

Importantly, agencies should mitigate emerging risks to rights and safety, including by regularly updating the predictive policing tool to improve it and reduce its risks.

Metrics used to assess accuracy in person-based predictive policing have much in common with those used in risk assessment, because person-based tools estimate, among other things, the likelihood of a person being involved in a crime. See chapter V for a detailed discussion of risk assessment.

e. Broadly Consider Community Resources to Address Root Causes of Crime

How predictive policing tools are used matters.⁵⁷ In addition to (or in place of) a policing response, community-based interventions built around resources for public health and social services can help alleviate the burdens on officers and mitigate concerns about predictive models.⁵⁸

f. Adopt Policies

Law enforcement agencies should adopt policies that address the use of the predictive policing system pre- and post-deployment. Such policies should govern which predictive policing systems are approved for use, as well as how predictive models are selected, trained, implemented,

Applications, 37 SEC. J. 47 (2024); see also Ned Levine, *The “Hottest” Part of a Hotspot: Comments on “The Utility of Hotspot Mapping for Predicting Spatial Patterns of Crime”*, 21 SEC. J. 295 (2008) (discussing Recapture Rate Index (RRI)); Veronica M. White & Joel Hunt, *Measuring How Relatively “Good” a Hot-spot Map Is: A Summary of Current Metrics*, 1 PROCEEDINGS OF THE IISE ANNUAL CONFERENCE & EXPO 2022, at 97 (2022), available at <https://www.ojp.gov/pdffiles1/nij/305365.pdf> (evaluating PAI, RRI, PEI, and PEI*).

⁵⁷ Cf. EXEC. OFF. OF THE PRESIDENT, BLUEPRINT FOR AN AI BILL OF RIGHTS 53–54 (Oct. 2022), <https://www.whitehouse.gov/wp-content/uploads/2022/10/Blueprint-for-an-AI-Bill-of-Rights.pdf> (describing “predictive policing” as an application of AI that may impact civil rights, civil liberties, or privacy).

⁵⁸ ANDREW GUTHRIE FERGUSON, THE RISE OF BIG DATA POLICING: SURVEILLANCE, RACE, AND THE FUTURE OF LAW ENFORCEMENT 173–76 (2017).

interpreted, and evaluated. Each of these policies should consider the tradeoffs between public safety, community trust, and the potential harms of algorithmic error, human bias, automation bias, and system inequalities that may result in misclassification. For further discussion of governance measures, see the Conclusion and Best Practices chapter.

g. Ensure Adequate Human Training and Assessment

Law enforcement organizations using predictive policing tools should require regular training on the operation and appropriate use of the specific tools. The training should be appropriate for the tools' diverse users, and it should include modules to ensure users think critically about the limits to algorithmic objectivity, error rates, ways in which these tools are embedded in existing inequitable systems, and strategies to prevent and mitigate human biases and systemic inequities.

Users of predictive policing models should receive training on how to interpret the model's outputs. The required training should also address how to determine what type of intervention, from police or other services, would be best to prevent or address the types of crime predicted in a given community, especially if predictions are driven by public health factors like substance abuse, homelessness, and mental health disorders.

Law enforcement personnel should not be authorized to use predictive policing systems if their training is not current.

V. Risk Assessment

Risk assessment instruments use quantitative analysis to estimate the likelihood that a certain outcome will occur in the criminal justice system, such as whether an individual will recidivate or fail to appear for trial.¹ These instruments are commonly based on statistical models.² Agencies and courts may use the model output to inform decisions about individuals, including whether to require detention before trial, what sentence to impose, and what interventions to attempt.

Risk assessment in criminal justice dates back nearly a century,³ and these types of tools are currently widely used.⁴ Based on the data available, it appears that every state implements some form of risk assessment,⁵ with considerable local variation,⁶ and most states have a state law or court rule that addresses risk assessment.⁷ At the federal level, risk assessment instruments

¹ See P'SHIP ON AI, REPORT ON ALGORITHMIC RISK ASSESSMENT TOOLS IN THE U.S. CRIMINAL JUSTICE SYSTEM 7 (2019), <https://partnershiponai.org/paper/report-on-machine-learning-in-risk-assessment-tools-in-the-u-s-criminal-justice-system/> (discussing definitions of risk assessment).

² As noted in the introduction, this report follows the definition of “artificial intelligence” from OMB Memoranda M-24-10 and M-24-18, which does not depend on “the type of model” and includes “simple” models, such as regression models that use conventional statistics.

³ BERNARD E. HARCOURT, AGAINST PREDICTION: PROFILING, POLICING, AND PUNISHING IN AN ACTUARIAL AGE 47-76 (2013), <https://academic.oup.com/chicago-scholarship-online/book/21282>.

⁴ See *Risk Assessment Landscape*, PUB. SAFETY RISK ASSESSMENT CLEARINGHOUSE, <https://bja.ojp.gov/program/psrac/selection/risk-assessment-landscape> (surveying uses of risk assessment nationwide); *Risk Assessment Tool Database*, BERKMAN KLEIN CTR., <https://criminaljustice.tooltrack.org/> (compiling data from reports and jurisdictions about uses of risk assessment), *How Many Jurisdictions Use Each Tool?*, MAPPING PRETRIAL INJUSTICE, <https://pretrialrisk.com/national-landscape/how-many-jurisdictions-use-each-tool/> (summarizing survey data about uses of risk assessment); *Stanford Pretrial Risk Assessment Tools Factsheet Project*, STANFORD L. SCH. POL'Y LAB, <https://law.stanford.edu/pretrial-risk-assessment-tools-factsheet-project/> (collecting authoritative descriptions of common risk assessment instruments); NAT'L CTR. FOR STATE CTS., APPENDIX A: PROFILES OF ASSESSMENT INSTRUMENTS, https://www.ncsc.org/_data/assets/pdf_file/0014/27140/ran-appendix-a.pdf (same); *Pretrial Release: Risk Assessment Tools*, NAT'L CONF. OF STATE LEGISLATURES (June 30, 2022), <https://www.ncsl.org/civil-and-criminal-justice/pretrial-release-risk-assessment-tools> (surveying state laws and court rules that address the development and use of risk assessment instruments); *50-State Report on Public Safety Part 2, Strategy 2, Action Item 2*, THE COUNCIL OF STATE GOVERNMENTS, <https://50statespublicsafety.us/part-2/strategy-2/action-item-2/> (providing results of a 50-state survey on use of risk assessment in probation and parole); Cathy Hu et al., *National Scan of Policy and Practice in Risk Assessment Policy Brief Number 2017-01*, THE RISK ASSESSMENT CLEARINGHOUSE (July 20, 2018, 3:13 PM), <https://bja.ojp.gov/sites/g/files/xyckuh186/files/media/document/PB-Scan-of-Practice.pdf> (survey of risk assessment practices, policies, procurement, and perceptions in 43 states).

⁵ *50-State Report on Public Safety Part 2, Strategy 2, Action Item 2*, THE COUNCIL OF STATE GOVERNMENTS, <https://50statespublicsafety.us/part-2/strategy-2/action-item-2/> (providing results of a 50-state survey on use of risk assessment in probation and parole).

⁶ *Risk Assessment Landscape*, PUB. SAFETY RISK ASSESSMENT CLEARINGHOUSE, <https://bja.ojp.gov/program/psrac/selection/risk-assessment-landscape> (surveying uses of risk assessment nationwide).

⁷ NAT'L CONF. OF STATE LEGISLATURES (June 30, 2022), <https://www.ncsl.org/civil-and-criminal-justice/pretrial-release-risk-assessment-tools> (surveying state laws and court rules that address the development and use of risk assessment instruments).

inform decision-making by courts, court services, and prisons. The Model Penal Code also encourages the use of risk assessment.⁸

This chapter begins with an overview describing how risk assessment is currently used in critical stages of the criminal justice system, as well as the typical design of these instruments. Next, the chapter describes the potential benefits of quantitative risk assessment in making criminal justice more effective, equitable, and efficient, followed by a discussion of risks including inaccuracy, biases, errors in data, inadequate validation, and insufficient consideration of alternatives. The chapter closes with a set of recommendations for the development and use of risk assessment instruments.

Uses of Risk Assessment

Risk assessment instruments inform decisions throughout the criminal justice system, including the following phases of prosecution and detention.⁹ For each phase, this chapter discusses how risk assessment is used and provide examples of specific risk assessment tools currently in use. As discussed in greater detail below, risk assessment tools should be employed only after a thorough evaluation, including for predictive performance,¹⁰ potential bias,¹¹ and suitability for the subject population, as well as the implementation of appropriate precautions.

a. Pretrial Release

Pretrial service agencies and prosecutors make recommendations to judges and judicial officers who ultimately make decisions about the status of defendants before trial.¹² A defendant may be released with no or minimal conditions, released subject to supervision and conditions, required to post bail to be released, or detained. Risk assessment can inform these decisions by assisting decision-makers in estimating the likelihood that a defendant will fail to appear in court, commit another offense before trial, or pose a risk to public safety. Examples used at the state and local levels include the Virginia Pretrial Risk Assessment Instrument (VPRAI) and VPRAI

⁸ Am. Law Inst., MODEL PENAL CODE: SENTENCING § 6.03 reporter's note F (AM. L. INST., Proposed Final Draft 2017), available at https://robinainstitute.umn.edu/sites/robinainstitute.umn.edu/files/2022-02/mpcs_proposed_final_draft.pdf; *id.* at § 6.09, cmt. E.

⁹ See Melissa Hamilton, *Risk Assessment Tools in the Criminal Legal System – Theory and Practice: A Resource Guide*, Nat'l Ass'n of Crim. Def. Laws. 16 fig. 1 (Nov. 2020), <https://www.nacdl.org/getattachment/a92d7c30-32d4-4b49-9c57-6c14ed0b9894/riskassessmentreportnovember182020.pdf> (noting a range of criminal justice system decision points where risk assessment instruments are used).

¹⁰ As discussed further below, there are a range of possible predictive performance metrics for risk assessment models, including precision and the false positive rate.

¹¹ Similarly, as discussed further below, there are a range of possible bias metrics for risk assessment models, including calibration and predictive equality. See generally Alessandro Castelnovo et al., *A Clarification of the Nuances in the Fairness Metrics Landscape*, SCIENTIFIC REPORTS (2022); Simon Caton & Christian Haas, *Fairness in Machine Learning: A Survey*, ACM COMP. SURVEYS (2024).

¹² In the federal court system, federal laws establish a general presumption that a defendant should be released before trial unless the government establishes at a detention hearing that the defendant should be detained because they present a risk to the community or are unlikely to appear at subsequent proceedings. Many state court systems also have a presumption of release and allow for defendants to be released if they post cash bail.

Revised (VPRAI-R),¹³ Public Safety Assessment (PSA),¹⁴ Ohio Risk Assessment System Pretrial Assessment Tool (ORAS-PAT),¹⁵ and Correctional Offender Management Profiling for Alternative Sanctions Pretrial Release Risk (COMPAS PRRS-I and -II).¹⁶ In the federal system, U.S. Probation and Pretrial Services uses the Pretrial Risk Assessment (PTRA), an algorithmic tool developed by the Administrative Office of the U.S. Courts.¹⁷

b. Sentencing

Following a criminal conviction, a judge often has discretion to determine an appropriate sentence for the defendant, subject to applicable statutory ranges and sentencing guidelines.¹⁸ Risk assessment tools can support these decision-makers in estimating the likelihood that a defendant will recidivate, which is typically a factor in sentencing. Widely used risk assessments in sentencing include the Level of Service Revised (LSI-R) and Level of Service Case Management Inventory (LS/CMI),¹⁹ as well as the COMPAS General Recidivism Risk Scale (GRRS), Violent Recidivism Risk Scale (VRRS), and COMPAS-R Summative GRRS.²⁰ Some jurisdictions use

¹³ See MARIE VANNOSTRAND & KENNETH J. ROSE, VA. DEP'T OF CRIM. JUST. SERVS., PRETRIAL RISK ASSESSMENT IN VIRGINIA, VIRGINIA PRETRIAL RISK ASSESSMENT INSTRUMENT (2009), <https://www.dcjs.virginia.gov/sites/dcjs.virginia.gov/files/publications/corrections/virginia-pretrial-risk-assessment-report.pdf>; VA. DEP'T OF CRIM. JUST. SERVS., VIRGINIA PRETRIAL RISK ASSESSMENT INSTRUMENT - (VPRAI): INSTRUCTION MANUAL – VERSION 4.5 (2021), https://www.dcjs.virginia.gov/sites/dcjs.virginia.gov/files/publications/corrections/virginia-pretrial-risk-assessment-instrument-vprai_2.pdf; KENNETH ROSE, STANFORD L. SCH. POL'Y LAB, RISK ASSESSMENT FACT SHEET: VIRGINIA PRETRIAL RISK ASSESSMENT INSTRUMENT (VPRAI) (2019), <https://law.stanford.edu/wp-content/uploads/2019/06/VPRAI-Factsheet-FINAL-6-20.pdf>.

¹⁴ See ADVANCING PRETRIAL POL'Y & RSCH, *About the Public Safety Assessment*, <https://advancingpretrial.org/psa/about/>; KRISTIN BECHTEL (ARNOLD VENTURES), STANFORD L. SCH. POL'Y LAB, RISK ASSESSMENT FACT SHEET: PUBLIC SAFETY ASSESSMENT (PSA), <https://law.stanford.edu/wp-content/uploads/2019/05/PSA-Sheet-CC-Final-5.10-CC-Upload.pdf>.

¹⁵ Edward J. Latessa et al, *The Creation and Validation of the Ohio Risk Assessment System (ORAS)*, 74 Fed. Prob. J. 16 (2010), https://www.uscourts.gov/sites/default/files/74_1_2_0.pdf.

¹⁶ EQUIVANT, PRACTITIONER'S GUIDE TO COMPAS CORE (Apr. 2019) <https://web.archive.org/web/20190520172536/http://www.equivant.com/wp-content/uploads/Practitioners-Guide-to-COMPAS-Core-040419.pdf>.

¹⁷ *Pretrial Risk Assessment*, ADMIN. OFF. OF THE U.S. CTS., <https://www.uscourts.gov/services-forms/probation-and-pretrial-services/supervision/pretrial-risk-assessment>.

¹⁸ In several jurisdictions, a jury may determine the appropriate sentence. Research on how risk assessment interacts with human decision makers, discussed below, has focused on judges rather than juries.

¹⁹ Christopher T. Lowenkamp & Kristin Bechtel, *The Predictive Validity of the LSI-R on a Sample of Offenders Drawn from the Records of the Iowa Department of Corrections Data Management System*, 71 FED. PROB. J. 25 (2007) https://www.uscourts.gov/sites/default/files/71_3_4_0.pdf; JAMES AUSTIN ET AL., INST. ON CRIME, JUST. & CORR. AT GEO. WASH. UNIV., *Reliability and Validity Study of the LSI-R Risk Assessment Instrument, Final Report submitted to the Pennsylvania Board of Probation and Parole* (Jan. 9, 2003) <https://www.ojp.gov/ncjrs/virtual-library/abstracts/reliability-and-validity-study-lsi-r-risk-assessment-instrument>.

²⁰ Eugenie Jackson & Christina Mendoza, Equivant/Northpointe, *Setting the Record Straight: What the COMPAS Core Risk and Need Assessment Is and Is Not*, HARVARD DATA SCI. REV., Mar. 31, 2020 <https://hdsr.mitpress.mit.edu/pub/zhwo7ax4/release/7>; EQUIVANT, PRACTITIONER'S GUIDE TO COMPAS CORE (Apr. 2019) <https://web.archive.org/web/20190520172536/http://www.equivant.com/wp-content/uploads/Practitioners-Guide-to-COMPAS-Core-040419.pdf>; Antonio Cordella & Francesco Gualdi, *Algorithmic Formalization: Impacts on Administrative Processes*, PUB. ADMIN. (Aug. 27, 2024), <https://onlinelibrary.wiley.com/doi/10.1111/padm.13030>.

specialized instruments to evaluate the risk of recidivism in defendants convicted of sexual offenses.²¹

c. Prison Classification

If a defendant is in custody or has a custodial sentence imposed, a court, jail, prison, or supporting agency must make decisions about detention. These decisions can include the type of facility, housing unit assignment, placement in a general or special population, and availability of programs, services, and work. Risk assessment instruments can inform these decisions by calling attention to predictors of potential violence and other misconduct while in custody.²² The tools used for prison classification are predominantly designed to estimate the risk of recidivism after release,²³ though some have been modified or evaluated for predicting prison misconduct.²⁴ Tools

²¹ R. Karl Hanson et al., *Assessing the Risk and Needs of Supervised Sexual Offenders: A Prospective Study Using STABLE-2007, Static-99R, and Static-2002R*, 42 CRIM. JUST. & BEHAV., Dec. 2015, at 1205, 1205, <https://journals.sagepub.com/doi/full/10.1177/0093854815602094> (evaluating risk assessment tools applied to individuals convicted of sexual offenses).

²² See generally NAT'L INST. OF CORR., U.S. DEP'T JUST., OBJECTIVE PRISON CLASSIFICATION: A GUIDE FOR CORRECTIONAL AGENCIES (2nd ed. 2021), <https://nicic.gov/resources/nic-library/all-library-items/objective-prison-classification-guide-correctional-agencies>; Richard A. Berk et al., *A Randomized Experiment Testing Inmate Classification Systems*, 2 CRIM. & PUB. POL'Y 215 (2003) (describing a randomized study to evaluate a prison classification system); Patricia L. Hardyman et al., *Internal Prison Classification Systems: Case Studies in Their Development and Implementation*, NAT'L INST. OF CORR. (Jan. 2002), <https://nicic.gov/resources/nic-library/all-library-items/internal-prison-classification-systems-case-studies-their> (describing the prison classifications systems in several states); Joe Russo, Michael J.D. Vermeer, Dulani Woods & Brian A. Jackson, *Risk and Needs Assessments in Prisons: Identifying High-Priority Needs for Using Evidence-Based Practices*, RAND (Sept. 9, 2020), https://www.rand.org/pubs/research_reports/RRA108-5.html; Daryl G. Kroner & Jeremy F. Mills, *The Accuracy of Five Risk Appraisal Instruments in Predicting Institutional Misconduct and New Convictions*, 28 CRIM. JUST. & BEHAV., Aug. 2001, at 471, 471, <https://journals.sagepub.com/doi/10.1177/009385480102800405>; Thomas R. Kane, *The Validity of Prison Classification: An Introduction to Practical Considerations and Research Issues*, 32 CRIME & DELINQ., July 1986, at 367, 367, <https://journals.sagepub.com/doi/abs/10.1177/0011128786032003008> (describing methods for validating prison classification systems).

²³ James M. Byrne & Amy Dezemmer, *The Research Director Perspective on the Design, Implementation, and Impact of Risk Assessment and Offender Classification Systems in USA Prisons: A National Survey*, in HANDBOOK ON RISK AND NEED ASSESSMENT 48, 53 (Faye Taxman ed., 2016), <https://www.taylorfrancis.com/chapters/edit/10.4324/9781315682327-10/research-director-perspective-design-implementation-impact-risk-assessment-offender-classification-systems-usa-prisons-national-survey-james-byrneand-amy-dezemmer?context=ubx&refId=9495cf84-b71e-455f-9547-8e80237d7fa7> (reporting results of a multistate survey on prison classification systems and finding that most responding states used a risk assessment tool designed to predict recidivism).

²⁴ Edward Latessa et al., *Creation and Validation of the Ohio Risk Assessment System: Final Report*, UNIV. OF CINCINNATI CTR. FOR CRIM. JUST. RSCH. (July 2009), https://www.uc.edu/content/dam/uc/ccjr/docs/reports/project_reports/ORAS_Final_Report.pdf (describing the ORAS Prison Intake Tool, which is designed for prison intake use and validated as an estimate of post-release recidivism); Joshua S. Long, *Appropriate Classification of Prisoners: Balancing Prison Safety with the Least Restrictive Placements of Ohio Inmates* (Jun. 22, 2020) (Ph.D. dissertation, University of Cincinnati), <https://cech.uc.edu/content/dam/refresh/cech-62/school-of-criminal-justice/research/2020/Joshua%20S.%20Long%206-22-20.pdf>; Matthew Makarios & Edward J. Latessa, *Developing a Risk and Needs Assessment Instrument for Prison Inmates: The Issue of Outcome*, 40 CRIM. JUST. & BEHAV., Dec. 2013, at 1449, 1449, <https://journals.sagepub.com/doi/full/10.1177/0093854813496240> (comparing risk assessment tools in the prison classification context, one designed to estimate misconduct in custody and the other designed to estimate post-release recidivism, and concluding that are significant differences and a “one size fits all” approach to risk assessment may undermine validity).

that have been specifically designed and validated for predicting prison misconduct are less common.²⁵ In the federal system, the Federal Bureau of Prisons currently uses the Prisoner Assessment Tool Targeting Estimated Risk and Needs (PATTERN)²⁶ to track dynamic changes in risk and the Standardized Prisoner Assessment for Reduction in Criminality (SPARC-13)²⁷ to identify programmatic and treatment needs of inmates.

d. Probation, parole, and supervision

Risk assessment tools are commonly used to inform decisions about the appropriate level of supervision for convicted persons who are not in custody, conditions for release and reentry plans from custody, and eligibility for earned release from custody.²⁸ Assessments that attempt to calculate likelihood of recidivism are common in these settings, too, such as LSI-R and COMPAS GRRS. In the federal system, the Federal Bureau of Prisons uses PATTERN to inform eligibility for earned time.²⁹ The Administrative Office of the U.S. Courts also developed the Post Conviction Risk Assessment (PCRA), which aims to predict general and violent recidivism and is used by federal probation officers.³⁰

Risk Assessment Design

There are, broadly, two types of risk assessment instruments: actuarial models, which quantitatively combine factors to estimate the likelihood that a risk will occur, and structured

²⁵ E.g., Grant Duwe, *The Development and Validation of a Classification System Predicting Severe and Frequent Prison Misconduct*, 100 THE PRISON J., Mar. 2020, at 173 <https://journals.sagepub.com/doi/10.1177/0032885519894587> (describing the design and validation of a risk assessment tool for estimating the likelihood of misconduct in custody); Mark D. Cunningham et al., *An Actuarial Model for Assessment of Prison Violence Risk Among Maximum Security Inmates*, 12 ASSESSMENT, Mar. 2005, at 40 <https://journals.sagepub.com/doi/abs/10.1177/1073191104272815> (similar).

²⁶ NAT'L INST. OF JUST., 2023 REVIEW AND REVALIDATION OF THE FIRST STEP ACT RISK ASSESSMENT TOOL (Aug. 2024), <https://www.ojp.gov/pdffiles1/nij/309264.pdf>; Zachary Hamilton et al., *Tailoring to a Mandate: The Development and Validation of the Prisoner Assessment Tool Targeting Estimated Risk and Needs (PATTERN)*, 39 JUST. Q., Apr. 2021, at 1129, 1129, <https://www.tandfonline.com/doi/full/10.1080/07418825.2021.1906930>.

²⁷ NAT'L INST. OF JUST., NCJ 309349, 2023 REVIEW AND VALIDATION OF THE FEDERAL BUREAU OF PRISONS NEEDS ASSESSMENT SYSTEM (Sept. 2024), <https://www.ojp.gov/pdffiles1/nij/309349.pdf>.

²⁸ See, e.g., John Monahan & Jennifer L. Skeem, *Risk Assessment in Criminal Sentencing*, 12 ANN. REV. CRIM. PSYCH. 489, 494, 496 (2016) (describing risk assessments to shorten a sentence on the “back end”); Sheldon X. Zhang et al., *An Analysis of Prisoner Reentry and Parole Risk Using COMPAS and Traditional Criminal History Measures*, 6 CRIME & DELINQUENCY 167 (2014) (describing use of COMPAS in parole supervision).

²⁹ See U.S. GOV'T ACCOUNTABILITY OFF., PUB. NO. GAO-23-105139, BUREAU OF PRISONS SHOULD IMPROVE EFFORTS TO IMPLEMENT ITS RISK AND NEEDS ASSESSMENT SYSTEM, <https://www.gao.gov/assets/gao-23-105139.pdf> (describing how PATTERN is used for earned time eligibility).

³⁰ *Post Conviction Risk Assessment*, ADMIN. OFF. OF THE U.S. CTS., <https://www.uscourts.gov/services-forms/probation-and-pretrial-services/supervision/post-conviction-risk-assessment>; see also Seena Fazel et al., *The Predictive Performance of Criminal Risk Assessment Tools Used at Sentencing: Systematic Review of Validation Studies*, 81 J. CRIM. JUST., July–Aug. 2022 (analyzing sentencing and post-conviction tools).

judgments,³¹ which provide frameworks for applying experience and intuition.³² The simplest actuarial designs compute a sum of factors, without any weighting, to generate an overall risk score. More complex approaches use conventional statistical methods, typically linear or logistic regression. Recent models use advanced statistics and machine learning methods, such as boosted regression³³ or gradient-boosted decision trees.³⁴ Typically, the quantitative output of a model is converted into a category, such as low, moderate, or high risk, by applying predefined thresholds.³⁵

The features that risk assessments take into consideration also vary significantly. Some have over 100 distinct inputs, while others consider fewer than 10.³⁶ The general long-term trend is toward fewer features that have a more readily understandable relationship with outcomes in the criminal justice system.³⁷ Early risk assessments tended to emphasize subjective characterizations of defendants, derived from professional interviews.³⁸ Some early risk assessment tools also relied on impermissible factors, such as race.³⁹

Modern risk assessment instruments tend to employ more objective factors that are more closely related to criminal justice outcomes. The factors may be fixed (e.g., criminal history) or changeable over time (e.g., time since last infraction or participation in drug treatment).⁴⁰ Some recent risk assessment tools also account for stages of the criminal justice system and periodic

³¹ Structured judgments are not a focus of this report, since they fall outside the definition of artificial intelligence in OMB Memoranda M-24-10 and M-24-18. This chapter discusses them below as an important comparison for actuarial models, which can be AI within the OMB definition.

³² Sarah L. Desmarais & Jay P. Singh, *Risk Assessment Instruments Validated and Implemented in Correctional Settings in the United States*, THE COUNCIL OF STATE GOV'TS JUST. CTR. (Mar. 27, 2013), <https://csgjusticecenter.org/wp-content/uploads/2020/02/Risk-Assessment-Instruments-Validated-and-Implemented-in-Correctional-Settings-in-the-United-States.pdf> (surveying the design and validation of risk assessment tools); Jennifer L. Skeem & John Monahan, *Current Directions in Violence Risk Assessment*, VA. PUB. L. & LEGAL THEORY RSCH. PAPER NO. 2011-13 (Mar. 2011), https://gspp.berkeley.edu/assets/uploads/research/pdf/03-2011_Current_Directions_in_Violence_Risk_Assessment.pdf.

³³ Zachary Hamilton et al., *Tailoring to a Mandate: The Development and Validation of the Prisoner Assessment Tool Targeting Estimated Risk and Needs (PATTERN)*, 39 JUST. Q., Apr. 2021, at 1129, 1129, <https://www.tandfonline.com/doi/full/10.1080/07418825.2021.1906930>.

³⁴ Jon Kleinberg et al., *Human Decisions and Machine Predictions*, 133 Q.J. ECON., Aug. 2017, at 237, 237, <https://academic.oup.com/qje/article/133/1/237/4095198>.

³⁵ The thresholds are typically set when developing, validating, or revalidating a risk assessment tool. The thresholds and their quantitative and qualitative meanings vary by risk assessment tool.

³⁶ See the references accompanying the Uses of Risk Assessment section for factors that widely used risk assessment tools analyze.

³⁷ Bernard E. Harcourt, *Risk as a Proxy for Race: The Dangers of Risk Assessment*, 27 FED. SENT'G REP. 237 (2015), https://scholarship.law.columbia.edu/cgi/viewcontent.cgi?article=3568&context=faculty_scholarship (describing the long-term historical trend in risk assessment factors).

³⁸ James Bonta & D.A. Andrews, *Risk-Need-Responsivity Model for Offender Assessment and Rehabilitation*, PUBLIC SAFETY CANADA (Jan. 2007), <https://www.publicsafety.gc.ca/cnt/rsrscs/pblctns/rsk-nd-rspnsvty/rsk-nd-rspnsvty-eng.pdf>; D.A. Andrews, James Bonta & J. Stephen Wormith, *The Recent Past and Near Future of Risk and/or Need Assessment*, 52 CRIME & DELINQ., Jan. 2006, at 7, 7, <https://journals.sagepub.com/doi/10.1177/0011128705281756>; Desmarais & Singh *supra* note 32.

³⁹ Early risk assessment instruments considered race as a factor. Modern models do not. See Bernard E. Harcourt, *Risk as a Proxy for Race: The Dangers of Risk Assessment*, 27 FED. SENT'G REP. 237 (2015), https://scholarship.law.columbia.edu/cgi/viewcontent.cgi?article=3568&context=faculty_scholarship.

⁴⁰ These different categories of factors are sometimes delineated as “static” and “dynamic.”

reassessment and are designed to interact with treatment and supervision conditions that can affect risk calculations.⁴¹

Future risk assessment approaches could incorporate more frequent model validation and updates, as well as more frequent reassessments based on evolving factors.⁴² These changes could improve predictive performance and reduce biases. For example, researchers are currently exploring how to evaluate changes in risk based on real-time location and provide prompt support to individuals based on their unique risks and needs.⁴³ These types of ongoing assessments may create privacy risks for individuals, and future approaches may have to further account for both maximizing predictive performance and respecting privacy.⁴⁴

Potential Benefits for a More Effective, Equitable, and Efficient Criminal Justice System

The central promise of risk assessment is that empirical evaluation of the risk of future harmful behavior could be more accurate, transparent, and equitable than subjective human judgments alone.

More accurate assessments of risk have the potential for significant benefits. They could enable better alignment between justice system functions including rehabilitation, deterrence, and incapacitation, and the prevention of future offenses.⁴⁵ Better predictions could also more accurately identify people who are unlikely to reoffend, channeling them toward lesser pretrial restrictions, lesser sentences, and less restrictive conditions of release.⁴⁶

Efficiency in resource allocation is another possible upside of more accurate estimates of risk likelihood.⁴⁷ Costlier aspects of criminal justice monitoring or detention could be better

⁴¹ The 2022 annual review and validation of PATTERN, for example, found that individuals in federal prison could and often did see their risk level change from first to last assessment, independent of simply getting older (which is associated with lower risk). NAT'L INST. OF JUST., NCJ 305720, 2022 REVIEW AND REVALIDATION OF THE FIRST STEP ACT RISK ASSESSMENT TOOL 18, 35 (Mar. 2022), <https://www.ojp.gov/pdffiles1/nij/305720.pdf>.

⁴² See D. Michael Applegarth et al., *Imperfect Tools: A Research Note on Developing, Applying, and Increasing Understanding of Criminal Justice Risk Assessments*, 34 CRIM. JUST. POL'Y REV. 319, 323 (2023), <https://journals.sagepub.com/doi/epub/10.1177/08874034231180505> (noting support of dynamic factors by focus group of winners of the National Institute of Justice's Recidivism Forecasting Challenge).

⁴³ The National Institute of Justice recently funded research to develop a real-time, cellphone-based intelligent tracking system to monitor people who are on community supervision with the goal of flagging potentially risky behavior and providing support to avert such actions. See MARCUS ROGERS, NCJ 308693, AI ENABLED COMMUNITY SUPERVISION FOR CRIMINAL JUSTICE SERVICES 2–3 (2024) (Final Rep. to the Nat'l Inst. of Just., Award No. 2019-75-CX-K001), <https://www.ojp.gov/pdffiles1/nij/grants/308693.pdf>.

⁴⁴ See *id.* at 12–13, 36.

⁴⁵ See KLEINBERG ET AL., *supra* note 33 at 237; Erin Collins, *Punishing Risk*, 107 GEO. L.J. 57, 72-73 (2018), <https://www.law.georgetown.edu/georgetown-law-journal/wp-content/uploads/sites/26/2018/12/Punishing-Risk-2.pdf>; John Monahan & Jennifer L. Skeem, *Risk Assessment in Criminal Sentencing*, 12 ANN. REV. CLINICAL PSYCH., 489, 489, <https://doi.org/10.1146/annurev-clinpsy-021815-092945>.

⁴⁶ See Sarah L. Desmarais, John Monahan & James Austin, *The Empirical Case for Pretrial Risk Assessment Instruments*, CRIM. J. & BEHAV. (2021).

⁴⁷ Erin Collins, *Punishing Risk*, 107 GEO. L.J. 57, 76-77 (2018), <https://www.law.georgetown.edu/georgetown-law-journal/wp-content/uploads/sites/26/2018/12/Punishing-Risk-2.pdf>.

directed toward situations where there is greater predicted benefit.⁴⁸ For a given level of public resources, it may be possible to obtain a greater level of public safety.⁴⁹

Risk assessment instruments also have the potential to increase transparency of human judgments.⁵⁰ A risk assessment tool could be made publicly accessible, along with its design documentation, validation studies, and guidance to practitioners.⁵¹ The data used for risk assessment could provide the basis for review of decisions—and correction, if necessary—by affected individuals and their counsel. A subjective decision-maker, by contrast, might not (intentionally or otherwise) fully explain the information that they considered and how they arrived at a decision.

Equity is another important motivation for using risk assessment tools. Models can be designed and validated to minimize disparities in predictive performance across demographic groups.⁵² Models can also be designed and validated to ensure that individuals with similar salient characteristics, such as the crime charged and criminal history, receive similar estimates of risk likelihood.⁵³

When risk assessments are appropriately designed, validated, and used, actuarial models can be predictive of criminal justice outcomes and can outperform human judgments alone.⁵⁴

⁴⁸ There are, to be sure, other factors to consider in allocating criminal justice resources. Crime prevention is an important goal, but not the only goal.

⁴⁹ KLEINBERG ET AL. *supra* note 33.

⁵⁰ Alex Chohlas-Wood, *Understanding Risk Assessment Instruments in Criminal Justice*, BROOKINGS (June 19, 2020), <https://www.brookings.edu/articles/understanding-risk-assessment-instruments-in-criminal-justice/>.

⁵¹ The Public Safety Assessment (PSA), for example, is available to the public. So are the PSA's design documentation, validation studies, and guidance to practitioners.

⁵² This conception of equity is, in AI research and practice, often referred to as “group fairness.” See the design documentation and validation studies referenced in the first section of this chapter for examples of how risk assessment tool developers address equity considerations. As discussed further below, there are differing possible equity metrics, and risk assessment tools may be considered equitable by some metrics and not by others. The factors used for risk assessment can also be closely related to demographics (e.g., residential and income data can be related to race), adding further complexity to the challenge of measuring and mitigating biases.

⁵³ In AI research and practice, this conception of equity is often referred to as “individual fairness.” As discussed further below, the use of “cut points” with risk assessment tools can introduce or exacerbate risks to individual fairness. See Jane R. Bambauer, Tal Zarsky & Jonathan Mayer, *When a Small Change Makes a Big Difference: Algorithmic Fairness Among Similar Individuals*, 55 U.C. DAVIS L. REV. 2337 (2022).

⁵⁴ SHAMENA ANWAR ET AL., RAND CORP., RR-A3299-1, WHAT HAPPENS WHEN JUDGES FOLLOW THE RECOMMENDATIONS OF PRETRIAL DETENTION RISK ASSESSMENT INSTRUMENTS MORE OFTEN? 8 (2024) https://www.rand.org/pubs/research_reports/RR-A3299-1.html; Desmarais et al., *supra* note 31; Jodi L. Viljoen et al., *Are risk assessment tools more accurate than unstructured judgments in predicting violent, any, and sexual offending? A meta-analysis of direct comparison studies*, BEHAV. SCI. & LAW (2024), <https://onlinelibrary.wiley.com/doi/full/10.1002/bsl.2698>; Zhiyuan Lin et al., *The limits of human predictions of recidivism*, 6 SCI. ADV. (2020), <https://www.science.org/doi/10.1126/sciadv.aaz0652>; R. Karl Hanson & Kelly E. Morton-Bourgon, *The Accuracy of Recidivism Risk Assessments for Sexual Offenders: A Meta-Analysis*, PUBLIC SAFETY AND EMERGENCY PREPAREDNESS CANADA (2007), <https://publications.gc.ca/collections/Collection/PS3-1-2007-1E.pdf>; D. A. Andrews et al., *The Recent Past and Near Future of Risk and/or Need Assessment*, 52 CRIME & DELINQ. 7 (2006), <https://journals.sagepub.com/doi/10.1177/0011128705281756>; D. Mossman, *Assessing Predictions of Violence: Being Accurate about Accuracy*, 62 J. CONSULT. CLIN. PSYCHOL. 783 (1994), <https://pubmed.ncbi.nlm.nih.gov/7962882/>; Don M. Gottfredson, *Effects of Judges' Sentencing Decisions on Criminal Careers*, National Institute of Justice: Research in Brief (November, 1999),

Modern models that are widely used have been evaluated for differing predictive performance by race or gender, with mixed results.⁵⁵ As discussed further below, studies on the predictive performance and disparities of risk assessment tools have significant limitations, and results substantially differ by testing methods, prediction and bias metrics, populations of individuals studied, and translation from quantitative to qualitative findings.

It is possible that use of more advanced machine learning methods could meaningfully improve the performance and bias characteristics of risk assessment tools.⁵⁶ This area of research is nascent, however, and advanced models may be more difficult to understand and analyze.

The Risks of Risk Assessment

While there are potential benefits to risk assessment, those benefits may not be fully realized. Use of these tools can also potentially reinforce bias, inequality, and other problems in the criminal justice system. When considering implementation of a risk assessment tool, careful analysis of possible downsides is important.

Accuracy, which is a leading rationale for risk assessment in criminal justice, is also a leading concern. Applying standards widely used in scientific research, validation studies indicate that risk assessment models can better predict criminal justice outcomes than random chance, and possibly in some circumstances better than human judgment alone, but a risk materializing (or not) remains far from certain.⁵⁷ As an example, in one state's recent comparative evaluation of several

<https://www.ojp.gov/pdffiles1/nij/178889.pdf>; see generally Michael J. White et al., *The Meta-Analysis of Clinical Judgment Project: Fifty-Six Years of Accumulated Research on Clinical Versus Statistical Prediction*, 34 THE COUNSELING PSYCHOL. 341 (2006), <https://journals.sagepub.com/doi/10.1177/0011000005285875>; William M. Gove et al., *Clinical Versus Mechanical Prediction: A Meta-Analysis*, 12 PSYCHOL. ASSESS. 19 (2000), <http://zaldlab.psy.vanderbilt.edu/resources/wmg00pa.pdf>; Robyn M. Dawes et al., *Clinical Versus Actuarial Judgment*, 243 SCIENCE 1668 (1989), <https://meehl.umn.edu/sites/meehl.umn.edu/files/files/138cstixdawesfaustmeehl.pdf>.

⁵⁵ See the validation studies for risk assessment tools referenced in the first section of this chapter. There has been limited meta-analysis of relevant research. See Sarah L. Desmarais et al., *Predictive Validity of Pretrial Risk Assessments: A Systematic Review of the Literature*, 48 CRIM. JUST. & BEHAV. 398 (2021) (surveying validation studies for pretrial risk assessment tools and noting in supplementary material the demographics in studies); Seena Fazel et al., *The Predictive Performance of Criminal Risk Assessment Tools Used at Sentencing: Systematic Review of Validation Studies*, 81 J. CRIM. JUST., July–Aug. 2022 (surveying validation studies for risk assessment tools used at sentencing and noting population demographics). There has also been limited research directly comparing risk assessment tools using similar populations and methods. See JUD. COUNCIL CAL., PRETRIAL RISK ASSESSMENT TOOL VALIDATION (2022), <https://www.courts.ca.gov/documents/Pretrial-Pilot-Program-Risk-Assesment-Tool-Validation-2022.pdf> (reporting results, including on race and gender disparities, from a coordinated multicounty validation study that examined several risk assessment tools).

⁵⁶ Compare KLEINBERG ET AL., supra note 33 at 259-260 (finding that gradient-boosted decision trees substantially outperform logistic regression in predicting failure to appear) with Jongbin Jung et al., *Simple Rules for Complex Decisions*, April 2017, <https://doi.org/10.48550/arXiv.1702.04690> (finding that simple rules can have equivalent performance to random forest models in predicting failure to appear).

⁵⁷ JUD. COUNCIL CAL., PRETRIAL RISK ASSESSMENT TOOL VALIDATION (2022), <https://www.courts.ca.gov/documents/Pretrial-Pilot-Program-Risk-Assesment-Tool-Validation-2022.pdf>; Sarah L. Desmarais et al., *Performance of Recidivism Risk Assessment Instruments in U.S. Correctional Settings*, in HANDBOOK OF RECIDIVISM RISK/NEEDS ASSESSMENT TOOLS 15 (J.P. Singh, D.G. Kroner, J.S. Wormith, S.L. Desmarais & Z. Hamilton eds., 2018), <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781119184256.ch1>;

pretrial risk assessment instruments, among individuals with the highest possible risk score, about a sixth to a quarter (depending on the tool) were subsequently arrested for a violent offense.⁵⁸ As even this high arrest rate reflects, users of the tool should bear in mind that risk is not inevitability.

Beyond the field of criminal justice, research has demonstrated that even with exceptionally high-quality and long-term data, and even with sophisticated machine learning models, accurately predicting life outcomes such as an eviction or job layoff can be beyond reach.⁵⁹ Outcomes in the criminal justice system are similarly products of complex and interrelated individual, environmental, and societal factors.

Perpetuation of bias and inequality is also a significant risk.⁶⁰ Studies of risk assessment tools have, in some instances, found notable differences in predictive performance by race and gender.⁶¹ Moreover, a model that has similar predictive performance across groups by some metrics may have substantial differences by others. For example, a risk assessment tool that has comparable precision in predicting recidivism across demographic groups may have very different false positive rates across groups.⁶² Researchers have demonstrated that, under realistic assumptions, it can be mathematically impossible to achieve equality across multiple different bias

DESMARAIS ET AL (2021) *supra* note 55 at 398; Seena Fazel et al., *The predictive performance of criminal risk assessment tools used at sentencing: Systematic review of validation studies*, 81 J. CRIM. JUST. (2022), <https://pmc.ncbi.nlm.nih.gov/articles/PMC9755051/pdf/main.pdf>; Anne A. H. de Hond et al., *Interpreting area under the receiver operating characteristic curve*, 4 THE LANCET DIGITAL HEALTH 853 (2022), [https://www.thelancet.com/journals/landig/article/PIIS2589-7500\(22\)00188-1/fulltext](https://www.thelancet.com/journals/landig/article/PIIS2589-7500(22)00188-1/fulltext); Seth J. Prins & Adam Reich, *Criminogenic risk assessment: A meta-review and critical analysis*, 23 PUNISH. & SOC'Y 578 (2021) <https://www.ncbi.nlm.nih.gov/pmc/articles/PMC9385164/>; Sonja B. Starr, *Evidence-Based Sentencing and the Scientific Rationalization of Discrimination*, 66 Stan. L. Rev. 803, 842-62 (2014).

⁵⁸ JUD. COUNCIL CAL., PRETRIAL RISK ASSESSMENT TOOL VALIDATION (2022), <https://www.courts.ca.gov/documents/Pretrial-Pilot-Program-Risk-Assesment-Tool-Validation-2022.pdf>.

⁵⁹ E.g., Matthew J. Salganik, *Measuring the predictability of life outcomes with a scientific mass collaboration*, 117 PROCS. NAT'L ACAD. SCIS. 8399 (2020), <https://www.pnas.org/doi/10.1073/pnas.1915006117>.

⁶⁰ This discussion focuses on race and gender disparities, which have been a primary focus of relevant research. Other types of disparities may exist and have received limited study, such as on the basis of disability or language. For example, a risk assessment tool that uses employment as a risk factor without considering whether the person receives Social Security Disability Insurance may result in discrimination on the basis of disability. There is also limited research on disparities across subgroups combining race and gender.

⁶¹ See JUD. COUNCIL CAL., PRETRIAL RISK ASSESSMENT TOOL VALIDATION (2022), <https://www.courts.ca.gov/documents/Pretrial-Pilot-Program-Risk-Assesment-Tool-Validation-2022.pdf>; Desmarais et al (2021) *supra* note 55 at 398; Howard Henderson et al., *Determining Racial Equity in Pretrial Risk Assessment*, 86 FED. PROB. 26 (2022), https://www.uscourts.gov/sites/default/files/86_3_5.pdf; Matthew DeMichele et al., *The Public Safety Assessment: A Re-Validation and Assessment of Predictive Utility and Differential Prediction by Race and Gender in Kentucky* (Apr. 25, 2018), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3168452.

⁶² E.g., Julia Angwin et al., *Machine Bias*, ProPublica (May 23, 2016), <https://www.propublica.org/article/machine-bias-risk-assessments-in-criminal-sentencing> (demonstrating these properties in an evaluation of the COMPAS risk assessment tool).

metrics like these.⁶³ Navigating the subtle and complex tradeoffs across predictive performance and bias measures is a challenge inherent in risk assessment.⁶⁴

Data is one potential source of bias for risk assessment models. Models are based on data from the criminal justice system, which can encode existing disparities. Consider, as an example, a risk assessment model that uses housing stability as an input and that outputs a score for risk of future arrest. Housing stability may depend on biases in the housing market, and arrests may be affected by biases in policing. The model may incorporate and perpetuate these biases, much like it incorporates other trends in data. Some scholars question the extent to which risk assessment can mitigate disparities in the criminal justice system and society; risk assessment involves prediction based on past events, which may themselves reflect inequality.⁶⁵ Some scholars have further suggested that risk assessment may create circular processes, in which individuals or groups are affected by the criminal justice system, then flagged as higher risk, leading to further involvement in the system.⁶⁶

The process of designing, implementing, and validating risk assessment tools may also be affected by disparities. Risk assessment tools may perpetuate inequality when they are developed in a manner that does not incorporate, for example, the perspectives of individuals within affected communities.⁶⁷ Tools may also be developed using data that is readily available, rather than alternative types of data that may be more predictive of criminal justice outcomes and less prone to encoding historical biases (e.g., using arrests rather than convictions).⁶⁸

Data quality is another area of potential concern with risk assessment tools. Discrepancies in how data is collected and categorized can result in misleading, incomplete, and inaccurate records that do not accurately reflect factors relevant to risk assessment.⁶⁹ Some inputs to risk

⁶³ Jon Kleinberg et al., *Inherent Trade-Offs in the Fair Determination of Risk Scores*, Nov. 2016, <https://arxiv.org/abs/1609.05807>; Alexandra Chouldechova, *Fair Prediction with Disparate Impact: A Study of Bias in Recidivism Prediction Instruments*, Feb. 2017, at 2, 4, <https://arxiv.org/abs/1703.00056>; Richard Berk et al., *Fairness in Criminal Justice Risk Assessments: the State of the Art*, 50 *Sociol. Methods & Res.* 3 (2018), <https://journals.sagepub.com/doi/pdf/10.1177/0049124118782533>.

⁶⁴ See, e.g., Andrew Bell et al., *The Possibility of Fairness: Revisiting the Impossibility Theorem in Practice* (FacT '23: Proceedings of the 2023 ACM Conference on Fairness, Accountability, and Transparency, 2023), <https://dl.acm.org/doi/abs/10.1145/3593013.3594007>.

⁶⁵ See, e.g., Sandra G. Mayson, *Bias In, Bias Out*, 129 *YALE L.J.* 2122 (2019), <https://www.yalelawjournal.org/article/bias-in-bias-out>; Ben Green, *The False Promise of Risk Assessments: Epistemic Reform and the Limits of Fairness*, ACM FAT* (2020), <https://dl.acm.org/doi/10.1145/3351095.3372869>.

⁶⁶ See, e.g., Shawn D. Bushway, “Nothing Is More Opaque Than Absolute Transparency”: *The Use of Prior History to Guide Sentencing*, 2.1 *HARVARD DATA SCI. REV.* (2020), <https://hdsr.mitpress.mit.edu/pub/dudgcmk3/release/7>.

⁶⁷ Ngozi Okidegbe, *The Democratizing Potential of Algorithms?*, 53 *CONN. L. REV.* 739 (2022), https://scholarship.law.bu.edu/cgi/viewcontent.cgi?article=4138&context=faculty_scholarship.

⁶⁸ Jessica M. Eaglin, *Constructing Recidivism Risk*, 67 *EMORY L.J.* 59, 101-104 (2017), <https://scholarlycommons.law.emory.edu/cgi/viewcontent.cgi?article=1046&context=elj>.

⁶⁹ See, e.g., Sarah Lageson, *Criminally Bad Data: Inaccurate Criminal Records, Data Brokers, and Algorithmic Injustice*, 2023 *U. ILL. L. REV.* 1771, 1775–81, 1786 (describing sources of criminal record errors and the impact on automated decision-making). For example, inconsistent collection and coding of individuals’ race and ethnicity can complicate efforts to understand how justice outcomes for people of color differ from those of white individuals. See KELLY ROBERTS FREEMAN, CATHY HU & JESSE JANNETTA, *RACIAL EQUITY AND CRIMINAL JUSTICE RISK ASSESSMENT*, *URBAN INST.* 3-4 (2021), <https://www.urban.org/sites/default/files/publication/103864/racial-equity-and-criminal-justice-risk-assessment.pdf>.

assessment models can be subjective, such as the quality of an individual's relationship with their family, introducing possible inaccuracy and bias.⁷⁰ There has been little evaluation about how consistent practitioners and individuals being assessed are in determining these inputs.⁷¹

Transparency is also a concern. Individuals who are subject to a risk assessment tool (and their representatives) may not know that the tool was used or have sufficient information to understand how it works and how it performs. Affected individuals also may not be aware of the inputs provided to the tool or have an opportunity to correct mistakes. While some models are entirely open, providing free public access to design documentation, implementing materials, and validation studies, others are more restrictive. Commercial licensing requirements, nondisclosure agreements, and trade secret protections, for example, can inhibit evaluation and understanding.⁷² The White House Blueprint for an AI Bill of Rights recommends that AI models used in sensitive domains, such as criminal justice, should provide “meaningful access to examine the system.”⁷³ Access to models is essential for enabling independent evaluation to quantify performance and identify issues.

The use of risk assessment models in contexts where they have not been properly validated is another source of concern. A model that is trained on one population, at one time, for one purpose may perform very differently on other populations, at later times, or when used for other purposes.⁷⁴ There are significant differences in criminal trends and criminal law across jurisdictions, calling into question how well risk assessment models generalize. Risk assessment models are often used without any recent validation on the local population.⁷⁵ It is a best practice

⁷⁰ See Beth Karp, *What Even Is a Criminal Attitude? –And Other Problems with Attitude and Associational Factors in Criminal Risk Assessment*, 75 STAN. L. REV. 1431 (2023), <https://review.law.stanford.edu/wp-content/uploads/sites/3/2023/06/Karp-75-Stan.-L.-Rev.-1431.pdf>.

⁷¹ Sarah L. Desmarais et al., *Performance of Recidivism Risk Assessment Instruments in U.S. Correctional Settings*, in HANDBOOK OF RECIDIVISM RISK/NEEDS ASSESSMENT TOOLS (J.P. Singh, D.G. Kroner, J.S. Wormith, S.L. Desmarais & Z. Hamilton eds., 2018), <https://onlinelibrary.wiley.com/doi/abs/10.1002/9781119184256.ch1>.

⁷² Hannah Bloch-Wehba, *Access to Algorithms*, 88 FORDHAM L. REV. 1265 (2020), <https://ir.lawnet.fordham.edu/cgi/viewcontent.cgi?article=5649&context=flr>; Rebecca Wexler, *Life, Liberty, and Trade Secrets: Intellectual Property in the Criminal Justice System*, 70 STAN. L. REV. 1343 (2018), <https://review.law.stanford.edu/wp-content/uploads/sites/3/2018/06/70-Stan.-L.-Rev.-1343.pdf>; Cynthia Rudin et al., *The Age of Secrecy and Unfairness in Recidivism Prediction*, 2.1 HARV. DATA SCI. REV. (2020), <https://hdsr.mitpress.mit.edu/pub/7z10o269/release/7>; Cynthia Rudin et al., *Broader Issues Surrounding Model Transparency in Criminal Justice Risk Scoring*, 2.1 HARV. DATA SCI. REV. (2020), <https://hdsr.mitpress.mit.edu/pub/8jy98s9q/release/3>; Greg Ridgeway, *Transparency, Statistics, and Justice System Knowledge is Essential for Science of Risk Assessment*, 2.1 HARV. DATA SCI. REV. (2020), <https://hdsr.mitpress.mit.edu/pub/vu6rc1yv/release/7>.

⁷³ EXEC. OFF. OF THE PRESIDENT, BLUEPRINT FOR AN AI BILL OF RIGHTS 51 (2022), <https://www.whitehouse.gov/ostp/ai-bill-of-rights/>.

⁷⁴ ALEXANDRA CHOULDECHOVA & KRISTIAN LUM, SAFETY & JUST. CHALLENGE, THE PRESENT AND FUTURE OF AI IN PRE-TRIAL RISK ASSESSMENT INSTRUMENTS 3, 5 (2020), <https://safetyandjusticechallenge.org/resources/the-present-and-future-of-ai-in-pre-trial-risk-assessment-instruments/>; Erika Montana et al., *Cohort bias in predictive risk assessments of future criminal justice system involvement*, 120 PROCS. NAT'L ACAD. SCIS. (2023), <https://www.pnas.org/doi/10.1073/pnas.2301990120>.

⁷⁵ 50-State Report on Public Safety Part 2, Strategy 2, Action Item 2, THE COUNCIL OF STATE GOVERNMENTS, <https://50statespublicsafety.us/part-2/strategy-2/action-item-2/>.

to build a predictive model on the most representative data available, or failing that, to evaluate a model on representative data to determine if it is suitable for use.⁷⁶

There is limited research on the effects of risk assessment tools, both on the criminal justice systems and officials who use them and on the communities where they are implemented.⁷⁷ More research is needed to understand how practitioners integrate risk assessment information into decision-making, including to examine the risks of automation and confirmation biases, why officials sometimes choose to override recommendations based on risk assessment, and how risk assessment affects the accuracy and disparities of decisions.⁷⁸ Similarly, more research is needed on how risk assessment can impact communities. Studies of this type typically lack adequate controls to disentangle implementation of risk assessment from other trends, and the time horizon for studies may be too short to capture benefits and downsides.

Converting model output into categories (e.g., “high risk”), rather than presenting probabilities and confidence intervals, may omit important information for decision-makers.⁷⁹ An individual may be just over or under the threshold for a risk category, for example. While risk categories may have value in helping to explain results and counter precision bias, omitting more detailed information can deprive decision-makers of essential context.

A final challenge is that evaluating risk is fundamentally different from determining the appropriate treatment for an individual from a range of possible treatments.⁸⁰ Focusing on prediction of negative outcomes such as recidivism or failure to appear could divert resources from, and limit consideration of, interventions that may reduce the risk of those negative outcomes

⁷⁶ For example, the Minnesota Department of Corrections compared a nationally available tool to a state-specific one and concluded that “there is a home-field advantage to risk assessment,” because the Minnesota tool outperformed an off-the-shelf tool. Grant Duwe, Mn. Dep’t of Corr., *Evaluating Bias, Shrinkage and the Home-Field Advantage: Results from a Revalidation of the MnSTARR 2.0*, at 29–31 (2021).

⁷⁷ Jodi L. Viljoen et al., *Impact of Risk Assessment Instruments on Rates of Pretrial Detention, Postconviction Placements, and Release: A Systematic Review and Meta-Analysis*, 43 L. HUM. BEHAV. 397 (2019), <https://psycnet.apa.org/fulltext/2019-46921-001.pdf>; Megan Stevenson, *Assessing Risk Assessment in Action*, 103 MINN. L. REV. 303 (2018), <https://scholarship.law.umn.edu/cgi/viewcontent.cgi?article=1057&context=mlr>. Shamena Anwar et al., RAND, RR-A3299-1, *WHAT HAPPENS WHEN JUDGES FOLLOW THE RECOMMENDATIONS OF PRETRIAL DETENTION RISK ASSESSMENT INSTRUMENTS MORE OFTEN?* 8 (2024), https://www.rand.org/pubs/research_reports/RR-A3299-1.html; Matthew DeMichele et al., *What Do Criminal Justice Professionals Think About Risk Assessment at Pretrial?*, 83 FED. PROB. 32 (2019), https://www.uscourts.gov/sites/default/files/83_1_5_0.pdf; Sarah Riley, *Overriding (In)justice: Pretrial Risk Assessment Administration on the Frontlines*, FaccT ’24: Proceedings of the 2024 ACM Conference on Fairness, Accountability, and Transparency, 2024), <https://facctconference.org/static/papers24/facct24-35.pdf>; Brandon Garrett & John Monahan, *Judging Risk*, 108 CAL. L. REV. 439 (2020), <https://www.californialawreview.org/print/judging-risk>; Megan T. Stevenson & Jennifer L. Doleac, *Algorithmic Risk Assessment in the Hands of Humans*, AM. ECON. J.: ECON. POL’Y. (Nov. 2024), <https://www.aeaweb.org/articles?id=10.1257/pol.20220620>.

⁷⁹ Melissa Hamilton, *Risk Assessment Tools in the Criminal Legal System – Theory and Practice: A Resource Guide*, NAT’L ASS’N OF CRIM. DEF. LAWS. 44–48 (Nov. 2020), <https://www.nacdl.org/getattachment/a92d7c30-32d4-4b49-9c57-6c14ed0b9894/riskassessmentreportnovember182020.pdf>; Starr (2014), *supra* note 57.

⁸⁰ Chelsea Barabas et al., *Interventions over Predictions: Reframing the Ethical Debate for Actuarial Risk Assessment* (Conference on Fairness, Accountability, and Transparency, Proceedings of Machine Learning Research, 2018), <https://proceedings.mlr.press/v81/barabas18a/barabas18a.pdf>; Erin Collins, *Punishing Risk*, 107 GEO. L.J. 57 (2018), <https://www.law.georgetown.edu/georgetown-law-journal/wp-content/uploads/sites/26/2018/12/Punishing-Risk-2.pdf>; Starr (2014), *supra* note 57.

and otherwise improve individual rehabilitation and public safety. For example, emphasizing criminal justice risks may deemphasize interventions outside the criminal justice system, such as substance abuse treatment or job training, that could have positive effects.

Recommendations

a. Risk Assessment Tool Design, Implementation, and Validation

Before moving forward with a risk assessment tool, it is important to document objectives, performance and bias criteria, expected benefits and risks, and alternatives to using the tool. Risk assessment tools should only be used if the expected benefits clearly outweigh the risks, substantiated by adequate evidence. In some circumstances, it may be preferable to use trained professionals and structured judgments rather than an actuarial risk assessment model.

Throughout the process of deciding whether and how to use risk assessment, engagement with affected communities and stakeholders is essential. Their perspectives are vital for setting goals, identifying possible issues, and integrating mitigation strategies. Because risk assessment tools are built and validated on data, a critical early step is careful evaluation of the quality of available data, identifying gaps, inconsistencies, or biases that could affect risk assessment. Agencies should consider setting policies and procedures in place to collect, label, and store data efficiently and accurately, ensuring the integrity, completeness, and provenance of the data. Transparency about available data is a valuable part of maximizing the quality of the data used, as stakeholders may be well-positioned to identify opportunities and shortcomings with that data. If a tool involves risk factors within an individual's control, transparency can also enable affected individuals to better understand how to reduce their risk score.

If data under consideration for a risk assessment tool includes information covered by antidiscrimination law (e.g., race, gender, disability, or age), or foreseeable proxies for those protected characteristics (e.g., ZIP code as a proxy for race), designers and implementers should be especially cautious and ensure legal compliance. They should also bear in mind that large criminal justice databases are almost certain to contain errors. Criminal justice data may also reflect inconsistencies in recordkeeping and enforcement. Developers of risk assessment tools should account for these issues, trying to quantify them and build tools that are robust against errors and inconsistencies. At the same time, if there are types of data available that can serve a similar purpose in risk assessment, they should carefully consider tradeoffs. Relying on charges filed or convictions rather than arrests may reduce the possibility of bias, for example, by accounting for disparities in policing and involving additional components of the criminal justice system.

Developers of risk assessment tools should also be mindful of the value of a fit between the data they are utilizing and the jurisdiction in which the tool will be used. The data used for the output of the risk assessment tool should also closely resemble the decision that the tool will support. For example, a tool for predicting misconduct in custody would ideally be built with data about misconduct in custody, rather than data about recidivism after release. Also, where there is sufficient data in a jurisdiction, it is preferable to build a model specific to that jurisdiction rather than rely on a model built with data that may not be representative. And if a risk assessment tool

uses data from multiple jurisdictions, the tool should account for possible differences in data definitions across jurisdictions.

Risk assessment tools should also be evaluated in a real-world context before deployment. First, this evaluation should include performance in predicting real-world outcomes, on the specific population whose risk is being assessed, as well as measurement of predictive performance across demographic groups, to identify possible biases. There should also be a baseline included in the evaluation, such as predictions made by a human or a previous model, to comparatively assess a tool's performance. Second, deployment of new risk assessment tools should be phased, facilitating observation and mitigation of risks. This approach also provides data for comparatively assessing the community impacts of using risk assessment tool. Third, the design, implementation, and validation of a risk assessment tool should be reviewed by an evaluator who is independent of the tool's developer. Finally, if significant risks are identified before deployment, mitigations for those risks should be designed, implemented, and validated before the deployment continues.

b. Continuous Monitoring and Human Oversight

Continuous monitoring of risk assessment tools is essential for ensuring they continue to meet objectives and rapidly identifying emerging issues. These tools should be continuously monitored to identify possible changes in predictive performance, biases, or data. This monitoring should be automated, to the extent possible, to allow quick responses to changes. Tools should also be periodically revalidated, to ensure that trends in criminal justice and changes in laws, policies, and practices have not undermined performance. Similarly, tools should be retrained or updated on a periodic basis, to account for changing context. This process should consider prior feedback and outcomes, which may suggest ways of improving the tool's design.

Risk assessment tools should be paired with policies and procedures, and ideally automated monitoring, to ensure the accuracy of input data. They should not displace human decision-making in the criminal justice system. They should be accompanied by guidance and training for decision-makers about the tool and its limitations.

Individuals affected by a risk assessment tool should receive notice and explanation and should have an opportunity to respond. When a decision-maker in the criminal justice system uses a risk assessment tool, to the extent possible, the individual being assessed (and their representatives) should receive notice of the tool and explanation of how it is being used. Where feasible, an individual should be able to seek correction of errors in any input data or submit additional relevant context.

c. Training and Education

Individuals who use risk assessments in criminal justice should be trained on how to interpret a tool's scores, as well as the limits of prediction. With respect to limits, training curricula should pay special attention to the potential for bias and creating or perpetuating disparities.

Users of a risk assessment tool, including judges, prosecutors, defense attorneys, pretrial or probation officers, correctional staff, community supervision staff, and others, should be

informed about how the tool works and its limitations. Training programs and other educational offerings should evolve with changes in research and the tool.

d. Policies

Agencies should issue guidance to risk assessment users on appropriate and prohibited uses of the risk assessment tool. Agencies should also issue guidance to risk assessment users on when decisions from the risk assessment tool may require additional human oversight and intervention—for example, when there is high uncertainty in a risk assessment tool’s predictions, or when a decision has particularly high impact.

Agencies should provide public notice and documentation on the use of the system, as well as information on its design (including input data) and the testing conducted to mitigate risks of the system.

Agencies should ensure that there is meaningful access to the risk assessment tool, such that researchers and stakeholders can conduct their own evaluations.

Agencies should only use risk assessment tools where the prediction model can be made transparent to the public. Agencies should not use risk assessment tools where this core functionality is obscured by trade secret protections or other barriers.

Agencies should not use tools with prediction thresholds that cannot be evaluated and changed by agency policymakers.

Agencies should set a model’s prediction thresholds according to a methodology for achieving specific objectives, and agencies should provide decision-makers with individualized context about a risk assessment beyond a categorization.

e. Research

Future research should examine risk assessment tools to better understand predictive performance, biases, how they compare to alternatives, and how they affect decision-makers and communities.

VI. Conclusion & Best Practices

The preceding chapters offer insight into uses of AI within the criminal justice system. These uses potentially offer important benefits to the institutions and individuals in the criminal justice system—including judges, pretrial and probation officers, prosecutors, law enforcement officers, forensic professionals, criminal defendants, and defense attorneys—and to the broader public. The preceding chapters also discuss significant risks as well as important policy and technological considerations that users should understand and mitigate. Each chapter underscores that responsible use of AI is critical, especially in the criminal justice context, where the deployment of this technology is quickly evolving and where public safety and individual rights are on the line. Each chapter also recommends practices to mitigate risks and preserve privacy, civil liberties, and civil rights when using AI in the criminal justice context.

Central to all these recommendations is AI governance. As a general matter, governance “provides a framework for decision-making by establishing standards and procedures and clarifying roles and responsibilities.”¹ As applied to AI, governance includes the broad range of decisions surrounding the technology’s development, use, and safeguards. Further, in the context of the criminal justice system, AI governance must be adaptable to account for change, but it must also be grounded in enduring values. Indeed, AI governance in this space must account for civil rights and civil liberties just as much as technical considerations such as data quality and data security. Governance is central to the Department of Justice’s responsible use of AI, and this report emphasizes the importance of AI governance at the diverse law enforcement agencies and other criminal justice institutions using AI or contemplating its use.

AI governance will take different shapes, depending on the size, scope, mission, and resources of each actor or institution operating within the criminal justice system. As with other forms of governance, smaller or resource-constrained criminal justice actors—including those at the state, local, and municipal level—may not be in a position to fully implement each recommendation exactly as written in this report. However, a key takeaway is that AI presents a new set of considerations that actors and institutions in the criminal justice system should take into account in adapting their overall governance structures.

The recommendations included below, and detailed in each of the preceding chapters, draw from the work of a community of scholars across the United States and abroad, as well as the work of our colleagues elsewhere in the Federal Government on publications such as:

- The White House Office of Science and Technology Policy’s “The Blueprint for an AI Bill of Rights: Making Automated Systems Work for the American People”²;

¹ Governance, Digital.gov, <https://digital.gov/topics/governance/>.

² EXEC. OFF. OF THE PRESIDENT, BLUEPRINT FOR AN AI BILL OF RIGHTS 5 (Oct. 2022), <https://www.whitehouse.gov/wp-content/uploads/2022/10/Blueprint-for-an-AI-Bill-of-Rights.pdf>.

- NIST’s “Artificial Intelligence Risk Management Framework,”³ and Special Publication 1270, “Towards a Standard for Identifying and Managing Bias in Artificial Intelligence”⁴;
- OMB’s Memorandum M-24-10, “Advancing Governance, Innovation, and Risk Management for Agency Use of Artificial Intelligence,”⁵ and Memorandum M-24-18, “Advancing the Responsible Acquisition of Artificial Intelligence in Government”⁶;
- “National Security Memorandum on Advancing the United States’ Leadership in Artificial Intelligence; Harnessing Artificial Intelligence to Fulfill National Security Objectives; and Fostering the Safety, Security, and Trustworthiness of Artificial Intelligence”⁷ and the associated “Framework to Advance AI Governance and Risk Management in National Security.”⁸

The recommendations in this chapter are also goals that the Department of Justice is working toward for its own uses of AI, as described in the Department’s “Compliance Plan for OMB Memorandum M-24-10.”⁹

Foundations for AI Governance

To establish a durable and comprehensive AI governance program, criminal justice agencies should identify the problem to solve and the reasons why the use of AI is preferable to alternatives; establish clear organizational and reporting structures to provide oversight, monitoring, and evaluation; hire, train, and retain a workforce with adequate resources to devise, enact, and enforce policies; build, operate, and routinely monitor the AI systems; regularly evaluate the performance of systems for accuracy and any unintended biases or disparities; and mitigate their associated risks.

³ NAT’L INST. STANDARDS & TECH., ARTIFICIAL INTELLIGENCE RISK MANAGEMENT FRAMEWORK: GENERATIVE ARTIFICIAL INTELLIGENCE PROFILE (July 2024), available at <https://doi.org/10.6028/NIST.AI.600-1>.

⁴ Reva Schwartz et al., NAT. INST. STANDARDS & TECH. SPECIAL PUBL’N 1270, TOWARDS A STANDARD FOR IDENTIFYING AND MANAGING BIAS IN ARTIFICIAL INTELLIGENCE (2022), <https://nvlpubs.nist.gov/nistpubs/SpecialPublications/NIST.SP.1270.pdf>.

⁵ MEMORANDUM FROM SHALANDA D. YOUNG, DIR., OFF. MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, TO HEADS OF EXEC. DEP’T’S & AGENCIES, (Mar. 28, 2024), available at <https://www.whitehouse.gov/wp-content/uploads/2024/03/M-24-10-Advancing-Governance-Innovation-and-Risk-Management-for-Agency-Use-of-Artificial-Intelligence.pdf>.

⁶ MEMORANDUM FROM SHALANDA D. YOUNG, DIR., OFF. MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, TO HEADS OF EXEC. DEP’T’S & AGENCIES, (Sep. 24, 2024), available at <https://www.whitehouse.gov/wp-content/uploads/2024/10/M-24-18-AI-Acquisition-Memorandum.pdf>.

⁷ MEMORANDUM ON ADVANCING THE UNITED STATES’ LEADERSHIP IN ARTIFICIAL INTELLIGENCE; HARNESSING ARTIFICIAL INTELLIGENCE TO FULFILL NATIONAL SECURITY OBJECTIVES; AND FOSTERING THE SAFETY, SECURITY, AND TRUSTWORTHINESS OF ARTIFICIAL INTELLIGENCE (Oct. 24, 2024), available at <https://www.govinfo.gov/app/details/DCPD-202400945>.

⁸ FRAMEWORK TO ADVANCE AI GOVERNANCE AND RISK MANAGEMENT IN NATIONAL SECURITY (OCT. 24, 2024), available at <https://ai.gov/wp-content/uploads/2024/10/NSM-Framework-to-Advance-AI-Governance-and-Risk-Management-in-National-Security.pdf>.

⁹ U.S. DEP’T OF JUST., COMPLIANCE PLAN FOR OMB MEMORANDUM M-24-10 (Oct. 2024), available at <https://www.justice.gov/media/1373026/dl>.

a. Conduct a Needs and Alternatives Assessment

Criminal justice agencies should know what problem or function they are trying to solve or address and how both AI and non-AI alternatives might present a solution.¹⁰ For example, judges, pretrial and probation officers, prosecutors, law enforcements officers, and defense attorneys may be seeking to more accurately calculate the likelihood that an individual will have a particular outcome in the criminal justice system. That calculation can in turn help judges make transparent, non-discriminatory, and equitable decisions. Similarly, agencies seek to drive law enforcement efficiencies and best direct limited resources. Agencies could consider risk assessment tools as compared to subjective assessments, or predictive policing tools as compared to traditional resource allocation measures.

b. Establish a Clear Organizational Structure

Criminal justice agencies, like other institutions working within the criminal justice system, vary in size, scope, resources, and mission. Regardless of institutional size, however, a tiered human-centered structure that assigns to individuals distinct operational, decision-making, and oversight responsibilities commensurate with their function, training, and experience will be key to maintaining an effective AI governance program. For example, at the Department of Justice, the Attorney General has designated a Chief AI Officer (CAIO) and an Emerging Technology Board (ETB), as required in OMB M-24-10. The CAIO has primary responsibility for coordination of the Department's use of AI, efforts to promote AI innovation, and management of risks from the use of AI. The CAIO coordinates with components across DOJ and reports to the Deputy Attorney General. The ETB serves as the Department's AI Governance Board.

c. Catalog Uses of AI

Agencies should catalog existing and new uses of AI and facilitate organization-wide visibility into these uses. Awareness of the AI uses occurring within the agency is the foundation for other governance steps. As feasible, agencies should also aim for greater transparency by making their use cases public. For example, the Department of Justice's AI governance process begins with identifying existing, new, and planned uses of AI across the Department. Components will report AI use cases, per the Advancing American AI Act and OMB M-24-10, and the Department will also review procurement, privacy governance, and IT governance records to ensure comprehensiveness.

d. Hire, Train, and Retain AI Workforce

Agencies interested in using AI should hire, train, and retain a workforce that understands the technical aspects, operational considerations, and tradeoffs associated with the use of AI (and remains knowledgeable as AI rapidly evolves). An adequately trained workforce should be able to understand, among other things, the source and limitations of the data on which a system is trained, the potential sources of any discrimination and bias, how these sources can affect outcomes in the

¹⁰ For institutions already using AI, this might take the form of an inventory that addresses the functionality of the tool (what does it do?), alternative tools (are there other options to achieve that result?), and a justification for using the tool that accounts for risks and mitigation strategies (why was the tool chosen among the existing alternatives?).

criminal justice system, the steps necessary to mitigate bias and discrimination, and means of implementing effective monitoring and auditing of AI systems.

Pre-Deployment Measures

Prior to deploying AI systems, law enforcement agencies and other agencies involved in criminal justice should also undertake the following measures.

e. Implement Policies and Procedures

Use of AI in the criminal justice system must be governed by robust human-centered policies and procedures that describe the permitted and prohibited uses of the AI system, the data used to train the system, the metrics and processes used to validate or evaluate the accuracy of the output, the frequency of monitoring, and the risks of using a particular AI system. The use of AI in the criminal justice system must also comply with the Constitution; statutes; rules of evidence, discovery, and procedure; discovery obligations; and ethical commitments. These policies and procedures should be developed by agencies in collaboration with relevant stakeholders and should anticipate and incorporate civil rights and civil liberties concerns. In addition, these policies should be regularly re-evaluated to ensure accuracy and consistency with current technology, best practices, and legal requirements.

Keeping “Humans in the Loop:” Individuals and institutions operating in the criminal justice system that use AI systems must ensure that human judgment drives the system’s design, implementation, and use. Users should be mindful that AI cannot displace human decision-making (and, indeed, doing so would be counter to established laws, norms, and principles designating different roles to, for instance, judges, officers, and prosecutors), and AI system outputs should be reviewed and verified by a human being. Courts, agencies, and decision-makers should create decision-making processes and guidelines for the use of AI to ensure a human reviews any outputs for accuracy. For high-impact decisions (such as cause for arrest or length of a sentence), the output of an AI system should not be the sole basis for a decision. For example, Department of Justice decisions regarding investigative steps, detention, prosecution, and evaluation or analysis of forensic evidence may take into account AI outputs to assist in collection and analysis of information, but trained professionals make the decisions.

Community Outreach and Public Notice: Consistent with applicable law and government guidance, and to the extent feasible, criminal justice entities should actively engage the public and relevant stakeholders regarding the intended use of potentially rights-or-safety-impacting AI technologies.¹¹ Disclosures should not include sensitive operational and law enforcement information but should generally include a plain language description of the system; the data used to train the system (e.g., driver license photos in the case of a facial recognition AI system); who

¹¹ OMB Memorandum M-24-10 offers useful guidance to agencies assessing whether AI may be rights-impacting or safety-impacting. In particular, M-24-10 defines these terms (note that rights-impacting AI includes AI that, among other factors, may have a significant effect on civil rights, civil liberties, or privacy) and specifies AI uses presumed to be rights- or safety-impacting. An agency’s engagement with the community and stakeholders may identify additional, potentially community-specific factors to consider in assessing the possible impact of an AI technology.

will have access to the system; the circumstances for the system's deployment; and governance mechanisms in place, such as whether the system will have an opt-out mechanism once deployed. Disclosure and engagement regarding the general purpose of an AI system might be feasible and useful even if disclosing details of capabilities and implementation is infeasible. Engagement can help build trust and ensure public support for adopted technologies.

Defining and Categorizing Risk: Criminal justice entities using AI should define and measure AI risks to rights and safety; categorize AI uses based on those risks; adopt commensurately greater safeguards for uses with rights- and safety-impacting risks; and determine their risk tolerance as an organization. The potential impact of AI uses varies dramatically, and the rigor of governance processes to manage risk may vary accordingly.

Regular Updates: Policies and procedures will inevitably require updates and should contemplate the circumstances that may warrant revisions. Agencies may consider, for instance, reviewing their AI governance policies on a yearly basis to account for technological and legal developments in this area. Agencies should consult with their legal counsel as needed to make appropriate revisions to reflect intervening legal decisions or changes in the applicable rules or statutes in the civil or criminal context.

f. Identify and Analyze Training Data Quality and Sources

Prior to deploying any AI system, individuals and institutions operating in the criminal justice system should closely track and identify the quality and sources of the data on which the AI system was trained and assess the relevance of the data the system was trained on to their particular use and jurisdiction. This is a crucial consideration for all AI systems but may be particularly relevant to institutions intending to contract AI systems designed, administered, or maintained by third-party vendors who may use proprietary training data sets or may otherwise be unwilling to disclose such information to their clients. To that end, criminal justice actors should obtain legal advice and involve procurement professionals to understand and negotiate contracts that allow those actors to examine the training data and address other important civil rights, civil liberties, or privacy concerns.

g. Test Systems Under Deployment Conditions

Criminal justice entities should conduct rigorous pre-deployment testing of any AI or automated systems. To the extent possible, such testing should be performed by independent third parties following domain-specific best practices and under conditions that mirror actual end-to-end deployment conditions. This testing must account for and seek to reflect the complex operational and societal contexts in which these systems are used. Testing a system in isolation may not accurately or fully reflect the system's impact. Agencies should also use pilots and limited releases in advance of full deployment to identify and mitigate against potential risks. The performance of the AI system should be evaluated against the performance of the status quo—such as technological tools or human processes that AI will be replacing or supplementing—and alternatives, and performance factors should include testing for bias, discrimination, and disparities. To instill public confidence, agencies may also consider facilitating external testing as feasible.

h. Evaluate Risks and Implement Risk Mitigation Strategies

After testing, criminal justice agencies should identify rights- or safety-impacting use cases and possible mitigation measures. Mitigation measures should include reconsidering use of rights- or safety-impacting use cases for which risks cannot be effectively mitigated. The decision not to proceed with a use case must be a feasible option.

i. Create Technology-Specific Policies

Some AI use cases, such as facial recognition, create unique or significant considerations that may justify dedicated policies. Agencies should consider whether specific policies are appropriate for particular use cases.

Post-Deployment Measures

After deploying AI systems, criminal justice agencies should also undertake the following measures.

j. Monitor the AI System

Following deployment, entities should conduct regular audits and monitoring of their uses of AI. Entities should adopt a mechanism for objective, independent auditing of models and source code to proactively address concerns about model accuracy, reliability, outdated training or test data, and potential for bias and discrimination.

k. Evaluate New Uses

Agencies should also consider whether and how their use of the AI system might have evolved over time, or how the context in which the system was deployed may have changed. If agencies identify new uses or new data sources that were not considered or evaluated pre-deployment, or if the context changes substantially, those agencies should implement pre-deployment strategies to the new uses and new data and evaluate whether those new uses and new data present new risks to the rights or safety of individuals, whether the new use is appropriate in light of those risks, and any mitigation strategies to be executed.

l. Continue Engaging with the Public

Post-deployment, agencies should continue to actively consult the public, and in particular, communities they judge will be most likely to be affected by the implementation of AI in the criminal justice system. Community engagement should be proactive, and feedback should be solicited on a regular, ongoing basis. Entities should establish channels for affected community members to seek recourse and alternatives to AI systems where practicable.

While this report considers many critical uses of AI in the criminal justice system as specified by EO 14110, existing and possible uses extend well beyond these cases. Use cases associated with future and emerging AI technologies—such as generative AI—might also have a transformative impact on the criminal justice system. Establishing robust AI governance programs

at criminal justice agencies can position those agencies to address risks of these use cases and better capture their promise.

Legal Disclaimer

This report is intended to support the development of policies and practices that protect civil rights, civil liberties, privacy, and equity, and that promote democratic values in the building, deployment, and governance of automated systems. It reflects and describes academic studies, policy proposals, and advocacy positions that might not be adopted or endorsed by the U.S. government.

This report is non-binding and does not constitute U.S. government policy. It does not supersede, modify, or direct an interpretation of any existing statute, regulation, policy, or international instrument. It does not constitute binding guidance for the public or federal agencies and therefore does not require compliance with the principles described herein. It also is not determinative of what the U.S. government's position will be in any international negotiation. The Department's inclusion of recommendations does not indicate a determination that adoption of those recommendations would necessarily satisfy requirements set forth in existing statutes, regulations, policies, or international instruments, or the requirements of the federal agencies that enforce them. Recommendations are not intended to, and do not, prohibit or limit any lawful activity of a government agency, including law enforcement, national security, or intelligence activities. The appropriate application of the recommendations set forth in this report depends significantly on the context in which automated systems are being utilized. In some circumstances, application of these recommendations in whole or in part may not be appropriate given the intended use of automated systems to achieve government agency missions. Future sector-specific guidance will likely be necessary and important for guiding the use of automated systems in certain settings.

This report recognizes that law enforcement activities require a balancing of equities, for example, between the protection of sensitive law enforcement information and public access to information; as such, public access to information may not be appropriate, or may need to be adjusted to protect sources, methods, and other law enforcement equities. In all circumstances, federal departments and agencies remain subject to judicial and congressional oversight, as well as oversight by private and public civil liberties and privacy-focused organizations, as well as existing laws, policies, and safeguards that govern automated systems.

This report is not intended to, and does not, create any legal right, benefit, or defense, substantive or procedural, enforceable at law or in equity by any party against the United States, its departments, agencies, or entities, its officers, employees, or agents, or any other person, nor does it constitute a waiver of sovereign immunity.