

#### White Paper

# Responsible AI guidelines: A blueprint for ethical integration

The purpose of this paper is to provide key considerations in a simple, easyto-follow framework for governments looking to implement responsible AI.



Matt Peters Chief Technology Officer



Chris Zumberge VP, AI Solutions



Jessica Rodriguez Director, Business Solutions and Strategy

# Al is here to stay — here's what governments need to know

Of all the questions one could ask about the impact of artificial intelligence (AI) on technology and society, there is one which has a definite answer: Is AI here to stay? In short, yes. And it is rapidly changing the way we work, create, and interact daily.

The integration of AI into the public sector is accelerating, with the global AI market valued at \$371.7 billion in 2025 and projected to reach \$204 trillion by 2032, reflecting a compound annual growth rate (CAGR) of 30.6%.<sup>1</sup> According to a study done by the National Association of State Chief Information Officers (NASCIO), 72% of states surveyed have implemented enterprise policies and procedures for AI development and use.<sup>2</sup> These numbers underscore AI's potential to revolutionize public services, offering both efficiency gains and substantial economic benefits.

While researchers, developers, and startups alike are focused on advancing AI capabilities — working to build smarter large language models (LLMs), faster processing systems, and friendlier interfaces — there is a large portion of the population which has been thrust into figuring out exactly how to tame this technological wave. They have, willingly or unwillingly, become accountable for enforcing ethical and responsible AI. This is especially true of public sector organizations — with constituent satisfaction a top priority, governments and agencies must prioritize responsible AI use to protect themselves and their data.

The purpose of this paper is to provide a simple, easy-to-follow framework for governments and agencies looking to implement responsible AI. Our hope is that any organization, large or small, can leverage this framework and become empowered by AI's potential, rather than bogged down by its dangers.

# 8 key components of responsible AI guidelines

As AI technologies become increasingly integral to organizational operations, ethical considerations are essential in guiding their development and deployment. Ethical frameworks help safeguard privacy, <u>reinforce public trust</u>, and ensure that AI systems serve as tools for good rather than sources of worry.

To embrace AI without putting yourself and your organization at risk, it's important to have a set of guidelines to shape your AI journey. This section will explore some of the most important considerations for responsible AI use at all levels of your organization.

# 1. Transparency

Transparency in AI means that government agencies openly communicate when and how AI systems are used, making the systems' workings and decision criteria understandable to both experts and the public. It helps provide clarity about how AI systems work, the data they use, and the rationale behind their decisions to build trust and accountability. This is crucial — people must know how an algorithm influences decisions about them and have insight into why it made a recommendation.

Without transparency, AI decisions become "black boxes," undermining accountability and constituents' ability to challenge or appeal outcomes. To ensure transparency, consider building the following into your AI framework.



**Audit logs:** An AI audit log captures detailed information about the activities and decisions made by an AI system. It serves as a trail of the system's operations, and it typically includes the system's inputs and outputs, decision-making processes, user

interactions, system changes, and errors. Audit logs allow stakeholders to review and understand how decisions are being made, identifying areas for improvement, and reinforces transparency.

- Reasoning transparency: Closely document and monitor the paths an AI system takes to arrive at an outcome. This can explain the logic and processes involved in AI decision-making, inclusive of the data inputs, algorithmic steps, decision rules, and the final decision itself. Understanding AI reasoning is crucial for building trust with users, making it easier to validate decisions, identify errors, and refine your models for improved accuracy and fairness.
- Open-source intelligence gathering: Using open-source intelligence provides your Al systems with a wide range of publicly available data, helping to reduce biases and broaden the knowledge base. By using data that's open and accessible, Al models can be trained on more representative datasets, making their outcomes more reliable and transparent.

To ensure transparency in AI, some governments and agencies have published AI inventories and algorithmic impact assessments, and others have heightened transparency protocols for algorithmic decision-making. Incorporating the above considerations into your AI framework can make AI systems' logic and use cases more visible, helping your organization uphold opengovernment values and ensure decisions can be inspected and trusted by constituents.

# 2. Accountability

Accountability involves defining clear roles and responsibilities for AI outcomes, including mechanisms to enforce proper use of AI. Public sector organizations must ensure that if an AI system causes an error or harm, something or someone can be held responsible and take corrective action. This is especially vital in the public sector because AI tools can profoundly affect constituents' lives — from determining eligibility for services to guiding law enforcement — and someone must answer for their performance.

When it comes to ensuring AI accountability, consider the following:

- **Documentation:** Maintain transparency in your AI systems, clearly documenting decision-making processes and proactively adjusting practices based on your data.
- Audits and monitoring: Conduct regular audits and continuous monitoring of Al systems to ensure they operate as intended and adhere to established standards. This helps in identifying and rectifying issues promptly.
- Governance frameworks: Establish robust governance structures that outline clear
  roles, responsibilities, and procedures for managing AI systems. This includes defining accountability for decision-making and ensuring compliance with ethical standards.
- Training and education: Ensure that those involved in AI development and
  deployment understand the ethical, legal, and social implications of AI technologies and use them accordingly.

Government frameworks reinforce that AI accountability is critical. <u>Executive Order 13960</u> requires federal agencies to monitor, audit, and document AI system compliance with established safeguards. This means agencies are expected to train personnel in AI ethics, designate officials responsible for AI oversight, and be prepared to explain and justify AI-driven decisions. The National Institute of Standards and Technology (NIST) encourages similar steps, highlighting the importance of documentation and human review in bolstering accountability. These recommendations and frameworks provide a solid basis for state and local governments to build responsible AI frameworks themselves. By embedding accountability into AI governance, public agencies ensure there is always a human answerable for <u>how AI is used</u> — which is essential for public trust and compliance.

# 3. Reliability and safety

Reliability and safety refer to an AI system's ability to perform as expected, without causing unintended harm. In the public sector, these considerations are paramount — an unreliable algorithm in domains like healthcare, transportation, or criminal justice can lead to dangerous mistakes. Agencies need to have confidence that an AI tool will produce accurate and valid results across all populations and conditions. This means rigorous testing, validation, and ongoing monitoring of AI systems that includes flagging harmful content and ensuring training data is accurate and unbiased. And, importantly, it means communicating these standards across all levels of the organization and to the public.

# "In the public sector, these considerations are paramount — an unreliable algorithm in domains like healthcare, transportation, or criminal justice can lead to dangerous mistakes."

Neglecting reliability can have real consequences. For example, in 2025, the Washington Metropolitan Area Transit Authority (WMATA) faced challenges with its self-driving train system. The Washington Metrorail Safety Commission (WMSC) initially blocked the expansion of automated trains due to concerns over increased instances where trains overran stations. Although the situation didn't pose a direct safety threat, the service disruptions caused by the systems underscore the necessity for rigorous testing and validation of AI systems in public transportation to ensure reliability and user satisfaction.<sup>3</sup>

Safety is closely tied to reliability. AI must be resilient to errors, bias, and manipulation so it does not endanger constituents' wellbeing or rights. Building safe AI infrastructure can be achieved through responsible design and deployment, transparency and documentation, and rigorous training for both deployers and end-users. Considering safety throughout the development lifecycle can help you proactively deploy safe technologies, and consistent monitoring and refinement can improve safety with the AI tool's use.<sup>4</sup>

By prioritizing reliability and safety, government organizations fulfill their duty of care, deploying AI only when it is trustworthy and aligned with best practices, government recommendations, and industry standards.

# 4. Privacy and security

Privacy and security are critical when integrating AI into public sector operations, given the sensitive personal data that government systems often handle. AI can intensify privacy risks by processing large volumes of personal data or by drawing inferences about individuals, and the misuse of this data can have detrimental effects on public trust, organizational operations, financial bottom lines, and more. Public sector organizations must protect individuals' personal information from misuse and give people control over how their data is used.

This translates to practices like data minimization (using only data that is truly needed), deidentification of personal information, and obtaining consent or providing opt-outs where appropriate. Several states now have laws granting individuals the right to opt out of automated profiling or requiring algorithmic impact assessments to evaluate privacy impacts — these are US analogues to the European Union's (EU) General Data Protection Regulation (GDPR) emphasis on data rights.<sup>5</sup>

Security involves safeguarding AI systems and data against breaches, hacking, or other malicious threats. Equally vital, AI systems must be resilient against cyberattacks and data leaks. Government AI often involves critical infrastructure like power grids and public safety systems or sensitive constituent information like tax records and biometrics for identity verification. A breach or manipulation of such systems can undermine public trust and cause harm.

Agencies are responding by strengthening <u>AI cybersecurity measures</u> as public expectations and legal standards evolve. Robust privacy protections like data encryption, differential privacy, and strict access controls, combined with strong security practices (regular audits, patching vulnerabilities, and adherence to foundations like <u>NIST's Cybersecurity Framework</u>), are nonnegotiable for responsible AI in government.

## 5. Fairness and non-discrimination

Fairness and non-discrimination ensure that AI systems do not produce biased outcomes or treat people unfairly, especially along the lines of race, gender, age, or other protected characteristics. Government agencies have a legal and ethical obligation to uphold civil rights — deploying AI must not become a way to unintentionally deny equal treatment. However, without careful design, AI models can amplify historical biases that present in data.

For example, studies have found that facial recognition algorithms often have significantly higher false match rates for women, the elderly, and people of color.<sup>6</sup> This bias has led to real injustices. In one instance, a black man in Detroit was wrongfully arrested and held for 30 hours in jail after facial recognition software incorrectly matched his driver's license photo with store surveillance footage of a shoplifter.<sup>7</sup>

Such cases underscore the imperative for fairness. <u>AI used in policing</u>, hiring, benefits allocation, or any public service must be rigorously evaluated for disparate impacts. If an AI tool is found to disproportionately harm a protected group (e.g., denying them loans, flagging them as high risk at higher rates), agencies need to adjust or abandon the tool to comply with anti-discrimination laws.

United States policymakers are actively addressing algorithmic fairness. The Blueprint for an AI Bill of Rights explicitly includes "Algorithmic Discrimination Protections," stating Americans should not face discrimination by algorithms and urging regular disparity testing and public reporting on AI impacts.<sup>8</sup> Disparity testing assesses how an AI model's predictions or decisions vary across different populations within the data. It can involve identifying which characteristics are relevant to fairness concerns, analyzing different outcomes for different data groups, using metrics to determine bias, and addressing the cause of any bias that does occur. This is a databacked, repeatable strategy for evaluating AI fairness.

"The Blueprint for an AI Bill of Rights explicitly includes 'Algorithmic Discrimination Protections,' stating Americans should not face discrimination by algorithms and urging regular disparity testing and public reporting on AI impacts."<sup>8</sup>

At the state and local level, new laws are emerging. Colorado's 2023 AI Act imposes a duty to prevent algorithmic discrimination in high-risk systems, and New York City's Local Law 144

mandates annual bias audits of AI hiring tools to detect and publicly disclose any disparate impact.<sup>9</sup> These steps all aim to identify and mitigate bias through techniques like diverse training data, algorithmic fairness metrics, and independent audits.

# 6. Human oversight

Human oversight means that human officials remain involved in the deployment and operation of AI, providing the ability to review, intervene, or override algorithmic decisions. In the context of government, this is important because it combines the efficiency of AI with the judgment and accountability of human decision-makers. The foundation of human oversight acts as a safety net, it can catch errors the AI might make and ensures that final decisions consider nuances or ethical factors that an algorithm might miss. Many US guidelines advocate a "human in the loop" for high-stakes AI. The AI Bill of Rights highlights this need and emphasizes that people should be able to opt out of AI-driven processes and get timely help from a human decision-maker.<sup>10</sup>

In practice, this could look like a human caseworker reviewing an AI's recommendation before denying someone public benefits, or a judge being required to review risk scores from an algorithm in the justice system. This principle was tested in Michigan's unemployment fraud detection incident a few years ago, where an automated system falsely accused thousands of fraud. The lack of effective human oversight or review in that system led to massive errors and was later condemned, prompting reforms to reintroduce human review for contested claims.<sup>11</sup>

United States agencies are embedding human oversight into AI governance policies. Many government AI applications now include oversight committees or review boards.

All these measures ensure that AI remains a tool under human control. Human oversight preserves the fundamental notion that automated systems serve the public interest, but do not replace human responsibility or judgment in government services.

# 7. Respect for human autonomy

Respect for human autonomy means AI systems should augment human decision-making, not undermine human freedom or agency. In the public sector, this principle translates to giving individuals affected by AI-driven decisions the ability to understand and, where appropriate, contest or refuse those decisions. Citizens should never feel coerced or powerless because an algorithm made a call — whether it's an eligibility determination or a policing action — without room for human consideration.

Ensuring respect for autonomy involves measures like obtaining consent for AI data usage, providing opt-out choices, and avoiding "black box" systems that unilaterally determine outcomes with no human appeal. A black box AI is an AI system whose internal workings are a complete mystery to the users; while users can see the system's inputs and outputs, they can't see what happens within the AI tool to produce those outputs.<sup>12</sup> Such systems undermine human autonomy completely, and fly in the face of many <u>ethical AI frameworks</u> that list human autonomy as a core principle.

As AI use becomes more and more prevalent, state and local legislatures have been drafting governance to address the need for human autonomy when engaging with AI. For instance, some state AI laws echo GDPR-like rights. Colorado's AI Act, and proposals in states like Connecticut and Virginia, include provisions allowing people to opt out of data processing for the purpose of automated profiling decisions, especially those that produce legal, or equally significant, consequences for consumers.<sup>13</sup>

These safeguards protect individuals from being subject to an AI decision without recourse. By respecting human autonomy, public sector AI deployments adhere to democratic values. They

empower constituents with transparency, consent, and the assurance that a human element is always present, rather than subjecting people to a strict regime of "computer says so" with no human appeal. This principle ultimately maintains human dignity and agency in the face of automated processes.

# 8. Auditability

Auditability refers to the capacity to independently review and verify how an AI system operates and makes decisions, both before deployment and during use. For government agencies, making AI auditable is key to ensuring ongoing compliance with laws, detecting biases or errors, and improving systems over time.

An auditable AI system keeps records (such as logs of decisions, data inputs, and model versions) that allow internal or external examiners to trace outcomes back through the algorithm's process; this is sometimes called traceability. Auditability is important for accountability and trust. If a constituent challenges an automated decision (e.g., denial of a permit or a wrongful arrest based on AI), auditors should be able to reconstruct what the path was for AI's conclusions and how it arrived there.

The audit requirement forces algorithm providers to regularly test their systems for bias and gives regulators and the public a means to verify the tools' fairness. Real-world events highlight this need. After the Detroit facial recognition mismanagement, the settlement not only instituted oversight but also mandated a look-back audit (of all cases since 2017) where face recognition was used to obtain arrest warrants.<sup>14</sup> This retroactive audit is to uncover whether others were wrongfully implicated and is informing new policies, a clear example of auditability improving accountability.

# "The audit requirement forces algorithm providers to regularly test their systems for bias and gives regulators and the public a means to verify the tools' fairness."

The concept of auditability also extends to technical features: developers might include audit logs or explanation modules in AI software to facilitate later examinations. By making AI systems auditable, organizations create a feedback loop for transparency and improvement.

This also allows them to verify that AI decisions are lawful and fair after deployment, instead of blindly trusting that they will be free from bias or error. This practice, increasingly reinforced by law and policy, helps ensure that AI remains under rigorous oversight throughout its life cycle, and secures public trust in AI-driven government services.

# Implementing responsible AI practices

When implementing responsible AI, collaboration across public agency departments and organizations is crucial. By leveraging diverse expertise to create robust systems, teams can foster a shared vision for cohesive efforts. Pooling resources and knowledge ensures efficient use of time, technology, and funding, which is particularly important in resource-constrained environments. Ultimately, collaboration enhances the capacity to develop AI systems that are not only technologically advanced, but also ethically sound and socially beneficial.

Building a framework for AI accountability and ethical use doesn't have to be complicated, there are 6 primary actions to take.

#### Assess

- Conduct a needs assessment
- · Define objectives
- Evaluate risk

## Implement

- Pilot projects
- Protect data
- Implement system monitoring

#### Engage

- Develop ethical standards
- Formalize an AI policy
- Engage stakeholders

# Audit

- Establish audit processes and feedback mechanisms
- Maintain reports and document usage

## **Educate**

- Educate teams
- Build expertise
- Provide opportunities for continuous learning

# Review

- Evaluate outcomes against established ethical standards
- Adapt strategy and guidelines

# The importance of responsible AI adoption

The key to responsible AI adoption is being intentional and structured in the way AI tools are developed and used across organizations. Keeping both short- and long-term goals in mind as you build your strategy can help identify the most important areas to monitor. Maintaining open communication amongst those overseeing AI systems and those using them is critical.

Building a framework of guidelines will hold your teams accountable for responsible use, helping to protect data and your brand from unwanted interference. With a solid understanding of responsible use, you can harness the transformative power of AI all while staying protected against rising threats.

As a people-first service provider, CAI has partnered with state and local government agencies to automate processes and integrate <u>responsible AI practices</u> into their ecosystems with the employee and constituent experience front of mind.

We leverage our deep understanding of how policies and funding impact technology decisions in the public sector, as well as our 10+ years of experience in building and implementing intelligent systems, to help our clients move beyond the 'chatbot' toward meaningful, responsible AI.

Take the next step in your organization's responsible AI adoption journey by visiting <u>CAI.io/</u> <u>services/data-and-artificial-intelligence</u>.

#### References

- 1. "Artificial Intelligence (AI) Market." Markets and Markets. April 2025. <u>https://www.marketsandmarkets.</u> com/Market-Reports/artificial-intelligence-market-74851580.html
- 2. "How Are States Using Generative Artificial Intelligence?" Government Technology. November 12, 2024. https://www.govtech.com/artificial-intelligence/how-are-states-using-generative-artificial-intelligence
- 3. Rachel Weiner. "Accused of making Metro less safe, watchdog relents on self-driving trains." The Washington Post. May 21, 2025. <u>https://www.washingtonpost.com/dc-md-va/2025/05/21/green-yellow-metro-automated</u>
- "3 AI Risks and Trustworthiness." National Institute of Standards and Technology, U.S. Department of Commerce. March 2025. <u>https://airc.nist.gov/airmf-resources/airmf/3-sec-characteristics/#3-1-valid-and-reliable</u>
- 5. Hope Anderson, Nick Reem, Juliann Susas. "Automated Decision Making Emerges as an Early Target of State AI Deregulation." White & Case. March 7, 2025. <u>https://www.whitecase.com/insight-alert/</u> <u>automated-decision-making-emerges-early-target-state-ai-regulation</u>
- 6. Dr. Charles H. Romine. "Facial Recognition Technology (FRT)." NIST. February 6, 2020. <u>https://www.nist.gov/speech-testimony/facial-recognition-technology-frt-0</u>
- 7. ACLU. "Williams v. City of Detroit." January 29, 2024. <u>https://www.aclu.org/cases/williams-v-city-of-detroit-face-recognition-false-arrest</u>
- Adam S. Forman, Nathaniel M. Glasser. "The White House Releases "Blueprint for an AI Bill of Rights: Making Automated Systems Work for the American People." Epstein Becker Green. October 12, 2022. <u>https://www.workforcebulletin.com/the-white-house-releases-blueprint-for-an-ai-bill-of-rights-making-automated-systems-work-for-the-american-people</u>
- 9. Hope Anderson, et al. "Automated Decision Making..." <u>https://www.whitecase.com/insight-alert/</u> <u>automated-decision-making-emerges-early-target-state-ai-regulation</u>
- 10. Adam S. Forman, Nathaniel M. Glasser. "The White House Releases..." <u>https://www.workforcebulletin.</u> <u>com/the-white-house-releases-blueprint-for-an-ai-bill-of-rights-making-automated-systems-work-for-</u> <u>the-american-people</u>
- 11. Gretchen Carr. "Case Over the Michigan Unemployment Insurance Agency's Faulty Automated System Finally Settled." University of Michigan; Science, Technology, and Public Policy. August 2024. <u>https://stpp.fordschool.umich.edu/sites/stpp/files/2024-08/stpp-midas-explainer.pdf</u>
- 12. Matthew Kosinski. "What is black box artificial intelligence (AI)?" IBM. October 29, 2024. <u>https://www.ibm.com/think/topics/black-box-ai</u>
- 13. Hope Anderson, et al. "Automated Decision Making..." <u>https://www.whitecase.com/insight-alert/</u> <u>automated-decision-making-emerges-early-target-state-ai-regulation</u>
- 14. Hilary Golston, David Komer. "Facial recognition false arrest of man by Detroit police wins settlement." Fox 2. October 16, 2024. <u>https://www.fox2detroit.com/news/facial-recognition-false-arrest-man-detroit-police-wins-settlement</u>

#### About CAI

CAI is a global services firm with over 9,000 associates worldwide and a yearly revenue of \$1.3 billion+. We have over 40 years of excellence in uniting talent and technology to power the possible for our clients, colleagues, and communities. As a privately held company, we have the freedom and focus to do what's right — whatever it takes. Our tailor-made solutions create lasting results across the public and commercial sectors, and we are trailblazers in bringing neurodiversity to the enterprise.

Learn how CAI powers the possible at <u>www.cai.io</u>