

EVALUATION AS DUE PROCESS: CIVIL SERVICE IN AN AUTOMATED AGE

DANIEL E. HO*, OLIVIA MARTIN**,
AMY PEREZ*** & KIT RODOLFA****

Modern government is built on the civil service. The twentieth century consensus, emerging out of the New Deal, was based on a social safety net administered by civil servants to ensure fair administration and due process. Yet today, decades of denying agencies the legal flexibility and technical resources to modernize have created a crisis that fuels political backlash. This crisis is now driving hasty attempts at technological fixes: the “Department of Government Efficiency” seeks to dismantle the civil service through mass layoffs and rapid artificial intelligence (AI) deployment. At the same time, understaffed state agencies are experimenting with AI for benefits administration, with one claiming its system “taught itself all of the [state’s] eligibility rules” and eliminated a massive claims backlog in one month. Both approaches can rush toward AI solutions without adequate safeguards, presenting a false choice between dysfunctional status quo and reckless technological disruption, when thoughtful modernization could preserve administrative due process while meeting modern challenges.

* Daniel E. Ho is the William Benjamin Scott and Luna M. Scott Professor of Law, Professor of Political Science, Professor of Computer Science (by courtesy), Senior Fellow at the Stanford Institute for Economic Policy Research (SIEPR) and the Stanford Institute for Human-Centered AI (HAI), and Faculty Director of the Regulation, Evaluation, and Governance Lab (RegLab).

** Olivia H. Martin earned her J.D. at Stanford Law School and is a Ph.D. student in Economics at Stanford University and a Graduate Student Fellow at RegLab.

*** Amy Perez led unemployment insurance modernization work at the U.S. Department of Labor after years working at Colorado’s unemployment agency and is a Policy Fellow at RegLab.

**** Kit Rodolfa was a Principal Research Scientist at RegLab. We are grateful to Arnold Ventures and Stanford Impact Labs for supporting this work, leadership in various state unemployment offices for the conversations that helped inspire this article, Ellen Kwon for excellent research assistance, as well as to Phil Spesshardt, Andy Stettner, Nicole Zeichner, and seminar participants at Stanford’s RegLab for helpful feedback and comments.

This Article charts a path beyond this impasse by examining how administrative law should govern AI's evolving role in benefits adjudication. We trace how merit staffing requirements, designed to ensure procedural fairness and insulate decisions from politics, have also impeded responsible innovation. These regulations produce contradictory guidance that leaves states uncertain about AI's legal boundaries, forcing them to choose between violating timeliness mandates or merit staffing rules.

The pandemic intensified this dilemma, spurring widespread state experimentation—from fraud detection to chatbots to “auto-adjudication.” While AI tools can drastically reduce backlogs, increase outreach, and improve efficiency, they also introduce risks of error, opacity, and bias. Through a simple yet novel audit of state systems, we expose vulnerabilities in these largely unexamined experiments. Yet we also show, through the Internal Revenue Service’s (IRS’s) evaluation of its phone help line that integrated machine assistance, that rigorous assessment of human-AI systems is both feasible and essential.

*We argue that the path forward to vindicating constitutional and administrative law values lies in evaluation. While implicit in *Mathews v. Eldridge*, the failure to make this explicit has led to neglect. Evaluation as due process is not merely a bureaucratic exercise but a legal imperative—a means of operationalizing due process at scale and with flexibility that preserves the dignity, fairness, and accountability of administrative systems.*

INTRODUCTION.....	833
I. THE STATUS QUO: A STRAINED BENEFITS ADJUDICATION SYSTEM	838
A. <i>The Adjudication Process</i>	839
B. <i>The Demands of the Pandemic</i>	843
C. <i>Modernization</i>	848
II. GOVERNING BY MERIT STAFFING.....	851
A. <i>The Emergence of Merit Staffing Requirements</i>	851
B. <i>The Inherently Governmental Boundary to Merit Staffing</i>	853
C. <i>Interpretative Uncertainty in Defining Automation’s Role in Merit Staffing.</i> ..	855
III. EXPERIMENTATION UNDER LEGAL UNCERTAINTY.....	862
A. <i>Unfacilitated Experimentalism</i>	863
B. <i>The Three Waves of AI Tools</i>	864
1. <i>Wave One: Fraud Detection</i>	865
2. <i>Wave Two: Chatbots</i>	868
a. <i>Audit of All Identified Chatbots</i>	870
b. <i>Audit of Generative AI Chatbots</i>	874
c. <i>Legal and Policy Implications</i>	878
3. <i>Wave Three: Adjudication Assistance</i>	883
C. <i>The Federal Precursor to Modernization through Evaluation</i>	888
IV. THE LEGAL AND POLICY NECESSITY OF EVALUATION.....	891
A. <i>Evaluation of AI Systems Is Legally Mandated</i>	891

1. <i>Adopting AI Without Evaluation is Likely a Violation of Procedural Due Process</i>	891
2. <i>Adopting AI Without Evaluation Is Likely Arbitrary and Capricious</i>	895
3. <i>Adopting AI Without Evaluation Is Likely Contrary to Procurement Law and Policy</i>	898
B. <i>What Meaningful Evaluation Looks Like</i>	901
1. <i>Core Principles</i>	901
2. <i>Example: An Evidence-Based Adoption Cycle for AI Chatbots</i> ..	903
C. <i>Prospective Implementations of an Evaluation Mandate</i>	906
1. <i>Impact Litigation</i>	907
2. <i>Agency Guidance and Rulemaking</i>	908
3. <i>New Legislation</i>	909
4. <i>Procurement</i>	912
D. <i>Reinterpreting Merit Staffing for the AI Age</i>	914
1. <i>The Limitations of the Current Framework</i>	914
2. <i>Building Toward a New Role for Merit Staff</i>	915
CONCLUSION	916
APPENDIX TABLE A: RESPONSE SUMMARIES FOR CHATBOT AUDIT....	918
APPENDIX TABLE B: RESPONSE SUMMARIES FOR GENERATIVE AI CHATBOT AUDIT	921

INTRODUCTION

One of the central roles of government is the provision of the social safety net. As Justice Brennan wrote in *Goldberg v. Kelly*,¹ “Public assistance . . . is not mere charity, but a means to ‘promote the general Welfare, and secure the Blessings of Liberty to ourselves and our Posterity.’”² These significant constitutional interests counsel “uninterrupted provision” to those eligible to receive assistance; failure to do so risks violating citizens’ due process rights.³ The twentieth century consensus, emerging out of the New Deal, was that the safety net would be administered by a neutral civil service to ensure fair administration and due process. Yet the future role of safety net programs and the law governing them is at a critical juncture, shaped by three pressing challenges.

First, the infrastructure to administer these programs has become frozen in outdated frameworks due to both chronic resource constraints and rigid legal restrictions. This institutional paralysis—agencies operating with 1970s technology and 1930s bureaucratic constraints while processing twenty-first

1. 397 U.S. 254 (1970).
2. *Id.* at 265.
3. *Id.*

century caseloads⁴—has predictably fueled political backlash.⁵ The “Department of Government Efficiency” (DOGE) has seized on these failures to justify wholesale elimination, lambasting the civil service as an “antidemocratic” institution staffed by “millions of unelected, unappointed” officials.⁶ The solution is as sweeping as it is concerning: mass reductions in head count across federal agencies, with certain agencies “deleted outright,”⁷ accompanied by untested deployment of AI to fill the gaps.⁸ President Trump has personally acted on these critiques, making the civil service a focal point of the first days of his second term.⁹ This false choice between dysfunction and demolition

4. See generally JENNIFER PAHLKA, RECODING AMERICA 259–60 (2023) [hereinafter PAHLKA, RECODING AMERICA]; JENNIFER PAHLKA & ANDREW GREENWAY, NISKANEN CTR., THE *HOW WE NEED NOW: A CAPACITY AGENDA FOR 2025 AND BEYOND* 24 (2024), https://www.niskanencenter.org/wp-content/uploads/2024/12/Niskanen-State-Capacity-Paper_Jen-Pahlka-and-Andrew-Greenway-2.pdf [<https://perma.cc/PQR8-U96E>].

5. See Jennifer Pahlka, Opinion, *This Is How Democrats Can Counter Elon Musk*, N.Y. TIMES (Feb. 7, 2025), <https://www.nytimes.com/2025/02/07/opinion/democrats-elon-musk-doge.html> [<https://perma.cc/UN9Q-U66G>].

6. Elon Musk & Vivek Ramaswamy, Opinion, *The DOGE Plan to Reform Government*, WALL ST. J. (Nov. 20, 2024, 12:33 PM), <https://www.wsj.com/opinion/musk-and-ramaswamy-the-doge-plan-to-reform-government-supreme-court-guidance-end-executive-power-grab-fa51c020> [<https://perma.cc/F3PF-8TYJ>]. The Department of Government Efficiency (DOGE) is not (yet) a department, which would require an act of Congress. See U.S. CONST. art. II, § 2, cl. 2.

7. Eric Katz, *Trump’s ‘DOGE’ Commission Promises Mass Federal Layoffs, Ending Telework*, GOV’T EXEC. (Nov. 18, 2024), <https://www.govexec.com/workforce/2024/11/trumps-doge-commission-promises-mass-federal-layoffs-ending-telework/401111/> [<https://perma.cc/ME28-3G58>].

8. See Musk & Ramaswamy, *supra* note 6; Andrea Hsu, *Officially, 59,000 Federal Jobs Are Gone Under Trump. There’s More to the Picture*, NPR (June 6, 2025, 9:58 AM), <https://www.npr.org/2025/06/04/nx-s1-5421277/trump-federal-workers-layoffs-doge> [<https://perma.cc/NN5E-H958>]; Kyle Chayka, *Elon Musk’s A.I.-Fuelled War on Human Agency*, THE NEW YORKER (Feb. 12, 2025), <https://www.newyorker.com/culture/infinite-scroll/elon-musks-ai-fuelled-war-on-human-agency> [<https://perma.cc/3UFF-U89U>] (referring to DOGE’s implementation of artificial intelligence (AI), such as AI filters scanning Treasury Department grant proposals for “forbidden terms,” such as “techno-fascism by chatbot”); Coral Davenport, *Inside Trump’s Plan to Halt Hundreds of Regulations*, N.Y. TIMES (Apr. 16, 2025), <https://www.nytimes.com/2025/04/15/us/politics/trump-doge-regulations.html> [<https://perma.cc/87U6-EKM2>] (describing Elon Musk’s development of an AI tool “intended to comb through the 100,000-plus pages of the Code of Federal Regulations and identify rules that are outdated or legally vulnerable in the wake of the two Supreme Court decisions”).

9. See Exec. Order No. 14,171, 90 Fed. Reg. 8,625 (Jan. 20, 2025); Exec. Order No. 14,170, 90 Fed. Reg. 8,621 (Jan. 20, 2025); Memorandum on Restoring Accountability for Career Senior Executives, 90 Fed. Reg. 8,481 (Jan. 20, 2025); Memorandum on Hiring Freeze, 90 Fed. Reg. 8,247 (Jan. 20, 2025); Memorandum on Return to In-Person Work, 90 Fed. Reg. 8,251 (Jan. 20, 2025).

obscures a third path: equipping government with modern tools and appropriate oversight to fulfill its constitutional obligations.

Second, these systems are under tremendous duress precisely when citizens need them most, revealing how understaffing and outdated technology create artificial scarcity in public services. For example, during the pandemic, state unemployment agencies faced claim volumes as much as thirty-three times higher than normal.¹⁰ One state estimated that it would need at least three times as much staff to handle the increased workload; in actuality, staffing grew only 42% from 2020 to 2021.¹¹ As a result of these issues, over a third of applicants to unemployment insurance waited over a month for benefits for much of 2020.¹² Further, an estimated 11% of total unemployment insurance benefits paid during the pandemic—or \$100 billion—were identified as fraudulent.¹³ And these challenges extend beyond the pandemic: in 2023, nearly half of denials of food assistance claims were adjudicated improperly or without adequate notice.¹⁴ Such operational shortcomings do not just harm individual claimants; they fuel the narrative that government cannot deliver abundance and prosperity to its citizens.¹⁵

10. U.S. GOV'T ACCOUNTABILITY OFF., GAO-22-104251, UNEMPLOYMENT INSURANCE: PANDEMIC PROGRAMS POSED CHALLENGES, AND DOL COULD BETTER ADDRESS CUSTOMER SERVICE AND EMERGENCY PLANNING 12–13 (2022) [hereinafter 2022 GAO REPORT], <https://www.gao.gov/assets/gao-22-104251.pdf> [<https://perma.cc/8PNW-HDP4>].

11. *Id.* at 14; MINN. MGMT. & BUDGET, STATE OF MINNESOTA WORKFORCE PLANNING REPORT: FY 2021, at 8–9, 31 (2021), <https://mn.gov/mmb-stat/workforce-reports/2021.pdf> [<https://perma.cc/AMS4-YWF2>].

12. *Benefits: Timeliness and Quality Reports*, U.S. DEP'T OF LAB., <https://oui.doleta.gov/unemploy/btq.asp> [<https://perma.cc/K9R4-9F85>] (select “US Total,” “All First Payment Timeliness,” starting “2020 Jan” and ending “2020 Dec” and click submit; the values in column “28 Days” (roughly a month) are below 66% for June through October).

13. U.S. GOV'T ACCOUNTABILITY OFF., GAO-23-106696, UNEMPLOYMENT INSURANCE: ESTIMATED AMOUNT OF FRAUD DURING PANDEMIC LIKELY BETWEEN \$100 BILLION AND \$135 BILLION (2023) [hereinafter GAO PANDEMIC UNEMPLOYMENT FRAUD REPORT], <https://www.gao.gov/products/gao-23-106696> [<https://perma.cc/EQ64-MA7T>].

14. USDA FOOD & NUTRITION SERV., SUPPLEMENTAL NUTRITION ASSISTANCE PROGRAM: CASE AND PROCEDURAL ERROR RATES (CAPER) FISCAL YEAR 2023 (2024) [hereinafter CAPER RATES FY 2023], <https://fns-prod.azureedge.us/sites/default/files/resource-files/snap-qc-caper-fy23.pdf> [<https://perma.cc/U7SP-DK9B>]. This is unlikely to be due to increased strain during the pandemic, as applications to Supplemental Nutrition Assistance Program (SNAP) only increased by approximately 10% during the pandemic. USDA FOOD & NUTRITION SERV., SUPPLEMENTAL NUTRITION ASSISTANCE PROGRAM PARTICIPATION AND COSTS (2025), <https://fns-prod.azureedge.us/sites/default/files/resource-files/snap-annualsummary-8.pdf> [<https://perma.cc/U78T-QMGY>].

15. *See generally* EZRA KLEIN & DEREK THOMPSON, ABUNDANCE 83 (2025).

Third, many see AI and automation as solutions to these shortcomings. However, the government's track record with modernization efforts has been fraught, displaying either outright incompetence or design entrenched in scarcity thinking, using technology to deny rather than deliver benefits.¹⁶ The automated fraud detection system deployed by Michigan's labor department accused nearly 40,000 applicants of fraud with a staggering error rate of 93%, turning a tool meant to protect public resources into one that wrongfully denied them.¹⁷ Meanwhile, New York City's generative AI chatbot provided inaccurate legal advice, incorrectly informing landlords that they could reject tenants with Section 8 vouchers or that they could not evict tenants for refusing to pay rent.¹⁸ Without rigorous evaluation and human oversight, AI systems can entrench rather than eliminate the barriers that create artificial scarcity in public services.¹⁹

This Article addresses these central challenges to government by proposing a path forward that rejects both the dysfunctional status quo and the dismantling proposed by DOGE. Properly implemented, AI can expand access to benefits, reduce wait times from months to days, and ensure that eligible citizens receive the support they are entitled to. The choice is not between human judgment and machine efficiency, but rather how to combine both to achieve what neither can accomplish alone: a government that delivers on its constitutional promise of promoting the general welfare through accessible, accurate, and timely services. Our six key contributions advance debates about the future of government programs, the civil service, administrative law, and the digital state.

First, we provide an in-depth characterization of the operational reality of mass adjudication systems. Critiques of automation often fail to consider this empirical reality: state agencies face an acute staffing shortage and an onslaught of cases that cannot be handled by paper alone. We show this through case studies of two of the country's largest social safety net programs:

16. See generally PAHLKA, RECODING AMERICA, *supra* note 4.

17. Alejandro de la Garza, *States' Automated Systems Are Trapping Citizens in Bureaucratic Nightmares with Their Lives on the Line*, TIME (May 28, 2020, 2:24 PM), <https://time.com/5840609/algorithm-unemployment/> [<https://perma.cc/63LD-XLMK>]; Rachael Kohl, *Automated Statecraft: Faulty Programming and Improper Collections in Michigan's Unemployment Insurance Program*, 2024 WIS. L. REV. FORWARD 43, 44.

18. Kyle Orland, *NYC's Government Chatbot Is Lying About City Laws and Regulations*, ARS TECHNICA (Mar. 29, 2024, 4:22 PM), <https://arstechnica.com/ai/2024/03/nycs-government-chatbot-is-lying-about-city-laws-and-regulations/> [<https://perma.cc/TGJ3-RW75>].

19. Bryce Covert, *States Are Turning Their Public Benefits Systems over to AI. The Results Have Often Led to 'Immense Suffering'*, FAST CO. (Jan. 23, 2025), <https://www.fastcompany.com/91265363/states-are-turning-their-public-benefits-systems-over-to-ai-the-results-have-often-led-to-immense-suffering> [<https://perma.cc/S39Y-F5B8>].

unemployment insurance (UI), which provides transition payments for workers losing a job through no fault, and the Supplemental Nutrition Assistance Program (SNAP), which provides food benefits for low-income individuals. We explore how labor-intensive adjudication has struggled to deliver timely and accurate outcomes and how modernization efforts have remained piecemeal and constrained by legal and financial limitations.

Second, we trace the origins of state civil service or “merit staffing” requirements to the emergence of these benefits programs. We examine ambiguities and contradictions in federal policies that govern the boundaries between civil service and automation. Conflicting regulations mandate both that benefits determinations be made by merit staff and that technology be outsourced, creating profound uncertainty when automation is introduced into benefits adjudication.

Third, in light of legal uncertainty, we document a staggering diversity of state experimentation with AI-based tools in benefits adjudication. We categorize these efforts into three distinct waves of AI implementation—fraud detection, chatbots, and adjudication assistance—culminating in the use of Generative AI.

Fourth, to show the critical need for evaluation, we design and carry out a simple audit of state-level chatbot implementations, evaluating their capabilities, accuracy, and vulnerabilities in providing critical guidance to claimants. We show how these tools challenge the conventional line between general guidance and claimant-specific advice, with the latter conventionally seen as reserved for the civil service. Rigorous evaluation of the risks and benefits of such systems is sorely lacking.

Fifth, while evaluation, auditing, and assessing AI systems in benefits systems may seem like a daunting challenge, we show that this is eminently feasible. Indeed, we show that functionally, the federal government has already done so for the same kind of service that chatbots purport to offer: citizen assistance through call lines. The IRS’s modernization efforts to establish taxpayer help lines, which over time integrated machine assistance, highlight how a large, resource-constrained federal agency has implemented evaluation frameworks to navigate similar challenges, while balancing efficiency and accountability.

Finally, we argue that evaluation is the linchpin for reconciling automation with constitutional and administrative law values, notably due process. The current arsenal of legal frameworks—such as merit staffing and ex-post judicial review for due process—is inadequate to guide states in managing AI’s complex challenges. Modern benefits systems cannot rely alone on a select few plaintiffs, like those in *Goldberg* and *Mathews*, to ensure that they are meeting the demands of due process at the scale and speed inherent to AI. While scholars considering automation in the abstract have called for

transparency,²⁰ improved notice and re-designed hearings,²¹ and designing technical tools to ensure procedural regularity,²² evaluation uniquely accounts for the institutional realities of mass adjudication. Moreover, our proposal merges the dominant focus in AI and Machine Learning (ML) on evaluation standards with the day-to-day realities of public administration, demonstrating how to adapt these standards beyond mere lab benchmarks. Agencies must have clear benchmarks to determine whether a human or an automated system poses a greater “risk of an erroneous deprivation,”²³ particularly as automation and AI reshape benefits systems. Put differently, evaluation is due process.

The Article proceeds as follows. In Part I, we analyze the operational demands of benefits adjudication. Part II examines the history and law of merit staffing requirements. Part III examines the three waves of state experimentation with automation in benefits systems, offers our survey and audit results, and discusses the counterexample of evaluation of Internal Revenue Service (IRS) help lines. Part IV assesses existing legal sources of accountability and makes the case for evaluation as due process. Technology must assist and augment the civil service, and evaluation will be critical to designing government technology systems and modernizing government programs to complement and vindicate the values we have long expected of merit staff.

I. THE STATUS QUO: A STRAINED BENEFITS ADJUDICATION SYSTEM

Benefits administration—programs like UI and SNAP—represents one of the most frequent and consequential touchpoints between the public and government. Millions of Americans interact with these systems each year, relying on them during periods of economic hardship or personal crisis. As such, these programs play an outsized role in shaping public trust: when they function smoothly, they reinforce confidence in government’s capacity to deliver essential services. But when they fail—whether through delayed payments, administrative backlogs, or wrongful denials—they not only

20. See, e.g., Danielle Keats Citron, *Open Code Governance*, 2008 U. CHI. LEGAL F. 355, 355–56 (2008); Deirdre K. Mulligan & Kenneth A. Bamberger, *Saving Governance-By-Design*, 106 CAL. L. REV. 697, 770–83 (2018); Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, et al., *Model Cards for Model Reporting*, in FAT* ‘19: PROC. CONF. ON FAIRNESS, ACCOUNTABILITY, & TRANSPARENCY 220, 220 (2019).

21. See, e.g., Danielle Keats Citron, *Technological Due Process*, 85 WASH. U. L. REV. 1249, 1303–13 (2008); Kate Crawford & Jason Schultz, *Big Data and Due Process: Toward a Framework to Redress Predictive Privacy Harms*, 55 B.C. L. REV. 93, 93 (2014).

22. Joshua A. Kroll, Joanna Huey, Solon Barcoas, Edward W. Felten, Joel R. Reidenberg, David G. Robinson, et al., *Accountable Algorithms*, 165 U. PA. L. REV. 633, 662–72 (2017).

23. *Mathews v. Eldridge*, 424 U.S. 319, 335 (1976).

undermine the livelihoods of vulnerable claimants but also fuel broader frustrations with public institutions.

In 2020, roughly one in eight American adults received unemployment benefits.²⁴ That same year, nearly 40 million people—or one-eighth of the overall population, received food assistance through SNAP.²⁵ Both programs are emblematic examples of cooperative federalism²⁶: the federal government provides the majority of administrative funding and overarching guidelines, while state agencies handle administration and adjudication shaped by the specifics of state law. This structure allows for considerable flexibility in how state legislatures determine eligibility for benefits but often leaves little flexibility in how state agencies administer their programs on an operational level.²⁷

A. The Adjudication Process

Adjudicating benefits eligibility is a labor-intensive process for both agencies and applicants. In the case of UI, the process begins with a person submitting a claim to a state agency explaining why and how they separated from their job.²⁸ An adjudicator then decides their eligibility under state and federal law, generally evaluating the claimant on two criteria: whether they earned enough money during the applicable period and whether they left

24. It is estimated that somewhere between 23.6 million (according to U.S. Census Bureau records) and 45.4 million (according to Internal Revenue Service (IRS) records) unique individuals received unemployment insurance (UI) in 2020. Jeff Larrimore, Jacob Mortenson & David Splinter, *Unemployment Insurance in Survey and Administrative Data*, BD. OF GOVERNORS OF THE FED. RESV. SYS. (July 5, 2022), <https://www.federalreserve.gov/econres/notes/feds-notes/unemployment-insurance-in-survey-and-administrative-data-20220705.html> [<https://perma.cc/DQ3S-KWXJ>]. The adult population was approximately 258 million in 2020. Stella U. Ogunwole, Megan A. Rabe, Andrew W. Roberts & Zoe Caplan, *Population Under Age 18 Declined Last Decade*, U.S. CENSUS BUREAU (Aug. 12, 2021), <https://www.census.gov/library/stories/2021/08/united-states-adult-population-grew-faster-than-nations-total-population-from-2010-to-2020.html> [<https://perma.cc/LF8X-KVCJ>].

25. *Characteristics of SNAP Households: FY 2020 and Early Months of the COVID-19 Pandemic: Characteristics of SNAP Households*, USDA FOOD & NUTRITION SERV., <https://www.fns.usda.gov/research/snap/characteristics-households-fy20-early-covid19-characteristics/> [<https://perma.cc/Y9TD-FPRT>] (Jan. 24, 2025).

26. See Abbe R. Gluck, *Nationalism as the New Federalism (and Federalism as the New Nationalism): A Complementary Account (and Some Challenges) to the Nationalist School*, 59 ST. LOUIS U. L.J. 1045, 1048 (2015); Michael C. Dorf & Charles F. Sabel, *A Constitution of Democratic Experimentalism*, 98 COLUM. L. REV. 267, 340 (1998).

27. See *infra* Part II.

28. See *How Do I File for Unemployment Insurance?*, U.S. DEP'T OF LAB., <https://www.dol.gov/general/topic/unemployment-insurance> [<https://perma.cc/TP4B-6FSU>].

their job due to no fault of their own (e.g., laid off or needed to care for a sick family member).²⁹ The employer also is notified of the claim, is asked to provide details about the worker's separation, and may dispute the claim by alleging that the worker was terminated for cause or quit without reason.³⁰ If there are discrepancies between the claims, the adjudicator must conduct additional fact-finding with the claimant, the employer, or both.³¹ Fact-finding processes can take weeks, depending on the time taken to respond. Workers and employers can appeal the initial determination, resulting in a hearing before an administrative law judge.³²

If all goes smoothly, it takes two to three weeks from filing a claim to receiving the initial check.³³ Some claims are straightforward—for example, a layoff undisputed by an employer—but many require careful review, such as quitting a job for specific reasons allowed by law, like attending certain types of training.³⁴ The eligibility determination process is a significant filter: in 2022, only one quarter of unemployed individuals applied for UI, and only half of applicants received benefits due to eligibility issues.³⁵ Most of the unemployed who did not apply reported that they did not believe they were eligible to receive benefits.³⁶

SNAP applications begin similarly, with individuals completing an application either online or in person at a local welfare office. However, federal regulations mandate that all applicants also be interviewed.³⁷ During the

29. *Id.*

30. See EMP. & TRAINING ADMIN., EMPLOYMENT AND TRAINING HANDBOOK NO. 301, at 1–6, 10–13 (5th ed. 2005) [hereinafter ET HANDBOOK 301], https://oui.doleta.gov/dmstree/handbooks/301/5th/hb301_5.pdf [<https://perma.cc/CGY2-XN7K>]; see also, e.g., *Employer Fact-Finding Responses*, COLO. DEP'T OF LAB. & EMP., <https://cdle.colorado.gov/employers/myui-employer/resources/user-guide/employer-fact-finding-responses> [<https://perma.cc/93KN-E77U>].

31. See, e.g., *Employer Fact-Finding Responses*, *supra* note 30.

32. U.S. DEP'T OF LAB., UNEMPLOYMENT INSURANCE FACT SHEET 2, https://oui.doleta.gov/unemploy/docs/factsheet/UI_Program_FactSheet.pdf [<https://perma.cc/A8B2-YR25>].

33. *Id.* at 1.

34. See, e.g., CAL. CODE REGS. tit. 22, § 1256-5 (2025).

35. *Characteristics of Unemployment Insurance Applicants and Benefit Recipients Summary*, U.S. BUREAU OF LAB. STAT. (Mar. 29, 2023, 10:00 AM), <https://www.bls.gov/news.release/uisup.nr0.htm> [<https://perma.cc/X6KZ-PHGB>].

36. The UI take-up rate is around 77%, which is roughly in line with other social benefits programs. See, e.g., Stéphane Auray, David L. Fuller & Damba Lkhagvasuren, *Unemployment Insurance Take-Up Rates in an Equilibrium Search Model*, 112 EUR. ECON. REV. 1 (2019); Avraham Ebenstein & Kevin Stange, *Does Inconvenience Explain Low Take-Up? Evidence from Unemployment Insurance*, 29 J. POL'Y. ANALYSIS & MGMT. 111, 113 (2010).

37. 7 C.F.R. § 273.2 (2019).

interview, the eligibility worker must verify the applicant's information and review documents confirming their residence and income. To continue receiving benefits, recipients must recertify their income and eligibility at a frequency determined by the state (ranging from every six months to two years or more).³⁸ The interview, often available only during working hours, is commonly cited as the largest barrier to SNAP participation.³⁹

For both programs, federal agencies oversee state performance through statutorily mandated evaluation systems, though these systems only focus on the select areas emphasized in the law. The Social Security Act instructs the Secretary of Labor to only certify UI payments to states if their "methods of administration" are "reasonably calculated to insure full payment of unemployment compensation when due."⁴⁰ The U.S. Department of Labor (DOL) aims to meet this requirement by regularly evaluating the "Benefits Timeliness and Quality" (BTQ) of state agencies' eligibility determinations.⁴¹ For timeliness, DOL calculates the percentage of first payments made within two or three weeks, with the "acceptable" level of performance being 87%.⁴² For quality, DOL randomly audits 60–100 of a state's eligibility determinations

38. USDA FOOD & NUTRITION SERV., STATE OPTIONS REPORT 1 (2017), <https://fns-prod.azureedge.us/sites/default/files/snap/14-State-Options.pdf> [<https://perma.cc/3BYN-ESE6>]. During the pandemic, however, many states were able to waive the initial eligibility interview requirement for new applicants. See Memorandum from USDA Food & Nutrition Serv., to SNAP State Agencies, on Supplemental Nutrition Assistance Program (SNAP) – Adjusting Interview Requirements Due to Novel Coronavirus (COVID-19) – Blanket Approval 1 (Mar. 26, 2020), <https://fns-prod.azureedge.us/sites/default/files/resource-files/SNAP-COVID19-AdjustInterviewRequirements.pdf> [<https://perma.cc/2JNP-8A6H>]. However, the possibility for such waivers ended with the end of the official federal public health emergency. *COVID-19 Public Health Emergency*, U.S. DEP'T OF HEALTH & HUM. SERVS., <https://www.hhs.gov/coronavirus/covid-19-public-health-emergency/index.html> [<https://perma.cc/W3LG-NQW3>] (Dec. 15, 2023).

39. See ALL. TO TRANSFORM CALFRESH, ENROLLING MEDI-CAL PARTICIPANTS IN CALFRESH: WHAT WORKS?: LESSONS FROM COUNTY-LEVEL EXPERIMENTATION IN CALIFORNIA AND NATIONAL RESEARCH 24 (2019).

40. 42 U.S.C. § 503(a)(1).

41. ET HANDBOOK 301, *supra* note 30, at 1.

42. U.S. DEP'T OF LAB., UI PERFORMS CORE MEASURES, https://oui.doleta.gov/unemploy/pdf/Core_Measures.pdf [<https://perma.cc/K569-JNF8>]. The applicable standard is fourteen days for states with a "waiting week"—a non-compensable period that must be served before benefits commenced. This standard applies to most large states, including California, Texas, New York, and Florida. The twenty-one-day standard applies to states without a waiting week. 20 C.F.R. § 640.5 (2025).

each quarter.⁴³ Each issue is scored out of 100, with points deducted for inadequate fact-collection, improper application of law and policy, and insufficient reasoning in the written explanation.⁴⁴ A determination fails the audit if it scores below ninety-five, and the “acceptable” level of performance for a state is for 75% of determinations to pass.⁴⁵ From 2010 to 2019, states met this bar just 52% of the time; some states, like New Jersey, Florida, and Pennsylvania met this level less than 5% of the time.⁴⁶ States are also evaluated based on the extent of their improper payments—benefits both overpaid and underpaid.⁴⁷

SNAP similarly has a quality control system; however, it focuses almost exclusively on payment error rates—both over and underpayments.⁴⁸ However, the U.S. Department of Agriculture Food and Nutrition Service (FNS) also measures the case and procedural error rate (CAPER) at which households are improperly denied benefits due to an inaccurate or procedurally incorrect decision.⁴⁹ In 2019, 34% of denials were improperly made.⁵⁰ In 2023, this rate stood at 45%.⁵¹ FNS does not penalize states for “persistently high CAPERs.”⁵²

43. Large states draw a sample of 100 determinations, and small states draw a minimum sample of sixty determinations. The sample must be equally composed of separation (i.e., whether the claimant left the job due to no fault of their own) and non-separation issues (e.g., issues such as failing to actively search for work). ET HANDBOOK 301, *supra* note 30, at 1–2.

44. EMP. & TRAINING ADMIN., EMPLOYMENT AND TRAINING HANDBOOK NO. 401 (5th ed. 2017), https://www.dol.gov/sites/dolgov/files/ETA/handbooks/2017/ETHand401_5th.pdf [<https://perma.cc/86J9-5R3G>].

45. UPERFORMS CORE MEASURES, *supra* note 42.

46. Analysis of data from *Benefits: Timeliness and Quality Reports*, *supra* note 12. For these calculations, meeting the acceptable level of performance means scoring 75% or greater on both the non-separation and separations audits.

47. *See id.*

48. *SNAP Quality Control*, USDA FOOD & NUTRITION SERV., <https://www.fns.usda.gov/snap/qc> [<https://perma.cc/VHZ2-L9E4>] (Oct. 20, 2025). This is consistent with the language of the Food and Nutrition Act, which describes the quality control mandate as follows: “In carrying out the supplemental nutrition assistance program, the Secretary shall carry out a system that enhances payment accuracy and improves administration by establishing fiscal incentives that require State agencies with high payment error rates to share in the cost of payment error.” 7 U.S.C. § 2025(c).

49. *SNAP Case and Procedural Error Rates*, USDA FOOD & NUTRITION SERV., <https://www.fns.usda.gov/snap/qc/caper> [<https://perma.cc/DA9G-UUGK>] (June 30, 2025).

50. *Id.*

51. CAPER RATES FY 2023, *supra* note 14.

52. Dottie Rosenbaum & Katie Bergh, *SNAP Includes Extensive Payment Accuracy System*, CTR. ON BUDGET & POL’Y PRIORITIES (June 21, 2024), <https://www.cbpp.org/sites/default/files/6-28-23fa.pdf> [<https://perma.cc/6PCA-WY5V>].

There is a limited range of options for what the federal agency can do if states fail to meet performance benchmarks. In the case of UI, states are merely placed in corrective action plans if their annual performance is deemed unacceptable; in 2022, states submitted a total of 874 such action plans, one for each deficiency, such as failing to meet timeliness, quality, or appeals performance standards.⁵³ The only more severe corrective response available is the complete denial of administrative funding, the so-called “nuclear option.”⁵⁴

B. *The Demands of the Pandemic*

Applications to both UI and SNAP rose during the pandemic, as is expected for income-stabilizing benefits programs in times of economic crisis. For UI, this strain was exacerbated by the introduction of three UI expansion programs through the 2020 Coronavirus Aid, Relief, and Economic Security

53. U.S. DEP’T OF LAB., EMP. & TRAINING ADMIN., UNEMPLOYMENT INSURANCE PROGRAM LETTER (UIPL) NO. 17–22 (2022), <https://www.dol.gov/sites/dolgov/files/ETA/advisories/UIPL/2022/UIPL%2017-22/UIPL%2017-22%20-%20Blue%20.pdf> [<https://perma.cc/KJB5-LB5X>]; Agency Information Collection Activities; Comment Request; Unemployment Insurance (UI) State Quality Service Plan (SQSP) Planning and Reporting Guidelines, 88 Fed. Reg. 36,505 (June 5, 2023), <https://www.federalregister.gov/documents/2023/06/05/2023-11835/agency-information-collection-activities-comment-request-unemployment-insurance-ui-state-quality> [<https://perma.cc/HN8P-75S4>]. The corrective action plans may be quite brief; for example, many of Washington’s are just a few sentences each. *Unemployment Insurance Program*, WASHINGTON WORKFORCE INNOVATION AND OPPORTUNITY ACT STATE PLAN, <https://wioaplans.ed.gov/node/81486> [<https://perma.cc/UN2N-EAH2>].

54. Social Security Act, Pub. L. No. 74-271, sec. 303, § 503, 49 Stat. 620 (1935). Such extreme sanctions appear to never have been used, though some Senate Republicans urged their use against California during the pandemic in response to high levels of fraud in the state. Press Release, U.S. Senate Comm. on Health, Educ., Lab. & Pensions, Ranking Members Cassidy, Crapo Demand Transparency on DOL’s New Policy Forgiving Julie Su for Losing \$32 Billion to Unemployment Insurance Fraud (May 8, 2024), <https://www.help.senate.gov/ranking/newsroom/press/ranking-members-cassidy-crapo-demand-transparency-on-dols-new-policy-forgiving-julie-su-for-losing-32-billion-to-unemployment-insurance-fraud> [<https://perma.cc/5UG7-9MUQ>]. The inadequacy of this “nuclear option” is also lamented in Project 2025’s chapter on the U.S. Department of Labor (DOL), which calls on Congress to “provide DOL with more reasonable enforcement tools for the UI system.” Jonathan Berry, *Department of Labor and Related Agencies*, in MANDATE FOR LEADERSHIP: THE CONSERVATIVE PROMISE, PROJECT 2025: PRESIDENTIAL TRANSITION PROJECT 598 (2023), https://static.heritage.org/project2025/2025_MandateForLeadership_FULL.pdf [<https://perma.cc/BA25-7WVX>].

(CARES) Act.⁵⁵ Not even including the new CARES Act programs, regular initial UI claims submitted early in the pandemic ranged from eleven to thirty-three times higher than volume in the preceding three months.⁵⁶ States were particularly unprepared because administrative funding from the federal government is tied to claims-related caseloads, which were at historic lows immediately prior to the pandemic.⁵⁷

As a result, states struggled to adequately increase staff to accommodate this surge in claims; some agencies borrowed staff from other state departments or hired contractors to take on certain roles in call centers.⁵⁸ While some agencies were able to increase their capacity in these ways—for example, Florida increased staff by about 93%⁵⁹—the staffing increases paled in comparison to the increase in claim volume. Minnesota officials reported that “based on the model DOL uses to allocate funding [to states] . . . the Minnesota UI agency would have needed about 4,000 total staff”—nearly three times the level of their 1,500 employees⁶⁰—“a staff level they characterized as impractical.”⁶¹ Staffing up quickly was further complicated by the fact that claim adjudication demands attaining a high level of statutory understanding and institutional knowledge. Official training periods are typically six to eight weeks,⁶² but many workers take far longer to feel fully “up to speed.”⁶³ In Former U.S. Deputy Chief Technology Officer Jennifer

55. OFF. OF INSPECTOR GEN., U.S. DEP’T OF LAB., COVID-19: STATES STRUGGLED TO IMPLEMENT CARES ACT UNEMPLOYMENT INSURANCE PROGRAMS (2021), <https://www.oig.dol.gov/public/reports/oa/2021/19-21-004-03-315.pdf> [<https://perma.cc/6P3X-FLG8>]. These extensions, however, were not COVID-19 specific. The UI program was also expanded during the Great Recession.

56. 2022 GAO REPORT, *supra* note 10, at 12.

57. *Id.* at 7. Funding available for state administration declined by about 32% in real terms from 2010 to 2019. *Id.* at 7–8.

58. *Id.* at 13; *see also, e.g.*, Letter from Julie A. Su, Sec., Cal. Lab. & Workforce Dev. Agency, to Members of the Legislature 4 (Aug. 19, 2020), <https://abgt.assembly.ca.gov/sites/abgt.assembly.ca.gov/files/Sec.%20Su%20Response%20to%20Asm.%20Chiu%20Letter%208.19.20%20and%20Referenced%20Letter.pdf> [<https://perma.cc/53C2-RXDV>] (indicating that California used over 500 “surge staff” from Deloitte to staff call centers during the pandemic).

59. 2022 GAO REPORT, *supra* note 10, at 14.

60. MINN. MGMT. & BUDGET, STATE OF MINNESOTA WORKFORCE PLANNING REPORT: FY 2021 33 (2021), <https://mn.gov/mmb-stat/workforce-reports/2021.pdf> [<https://perma.cc/UQY8-T34T>].

61. 2022 GAO REPORT, *supra* note 10, at 14.

62. *See, e.g., id.* at 14; *Unemployment Benefit Specialist (Adjudicator) FAQs*, WIS. DEP’T OF WORKFORCE DEV., <https://dwd.wisconsin.gov/jobs/faqs/250725-ubs-adjudicator.htm> [<https://perma.cc/7FLW-BGM4>].

63. PAHLKA, RECODING AMERICA, *supra* note 4, at 40.

Pahlka's study of California's UI agency, an employee referred to himself as "the new guy" after seventeen years on the job—his tenure dwarfed by colleagues with over twenty-five years of experience.⁶⁴ In such settings, adding staff can initially lower productivity as senior staff are diverted to training new hires; for example, after hiring additional staff early in the pandemic, California's agency required two to five times longer to complete certain tasks.⁶⁵

The results of this staggering caseload coupled with inadequate staffing were predictable: increased wait times for applicants, a significant drop in decision quality, and heightened fraud rates.⁶⁶ Put differently, the UI system failed precisely at a time when American workers needed it most. Immediately prior to the pandemic, roughly 90% of first UI payments were made within three weeks, but as depicted in the bottom panel of Figure 1, this rate plummeted to nearly 50% with COVID-19, and did not rise above 80% until 2023.⁶⁷ At the beginning of the pandemic, around 10% of initial applicants—or around 250,000 people—waited over seventy days to receive their initial benefit, with some waiting as long as four months or more.⁶⁸ Quality measures also dropped sharply (top panel of Figure 1), from a pass rate of over 80% at baseline to below 60% in the first period after the pandemic. As staff were overwhelmed learning to operate new systems, many fraudulent claims were approved. In 2023, the Government Accountability Office estimated that states made out somewhere between \$100 and \$135 billion—about 11% to 15% of total benefits paid—in fraudulent payments during the pandemic, sparking congressional outrage.⁶⁹

64. *Id.*

65. YOLANDA RICHARDSON & JENNIFER PAHLKA, EMPLOYMENT DEVELOPMENT DEPARTMENT STRIKE TEAM DETAILED ASSESSMENT AND RECOMMENDATIONS 15 (2020), <https://www.govops.ca.gov/wp-content/uploads/sites/11/2020/09/Assessment.pdf> [<https://perma.cc/Q68Q-NNPN>].

66. See GAO PANDEMIC UNEMPLOYMENT FRAUD REPORT, *supra* note 13.

67. See *Benefits: Timeliness and Quality Reports*, *supra* note 12; U.S. DEP'T OF LAB., MANPOWER ADMIN., UIPL NO. 1126 (1971), https://oui.doleta.gov/dmstree/uipl/uipl_pre75/uipl_1126.htm [<https://perma.cc/H357-NEVG>].

68. Greg Iacurci, *He Was Stuck in Unemployment 'Limbo' for Four Months. Then Came \$23,000 in Benefits*, CNBC (Aug. 15, 2020, 11:52 AM), <https://www.cnbc.com/2020/08/15/coronavirus-mans-23000-in-unemployment-pay-delayed-four-months.html> [<https://perma.cc/GW-H3-AWH4>] (four months in New York); Kate Davidson, *Months Later, Thousands Still Waiting for Unemployment Benefits in Oregon*, OR. PUB. BROAD. (Aug. 5, 2020, 7:11 PM), <https://www.opb.org/article/2020/08/05/months-later-still-waiting-for-unemployment-benefits-in-oregon/> [<https://perma.cc/N7LZ-VG6J>] (five months in Oregon).

69. GAO PANDEMIC UNEMPLOYMENT FRAUD REPORT, *supra* note 13; Richard Lardner, Jennifer McDermott & Aaron Kessler, *How Billions in COVID-19 Pandemic Relief Aid Was Stolen or*

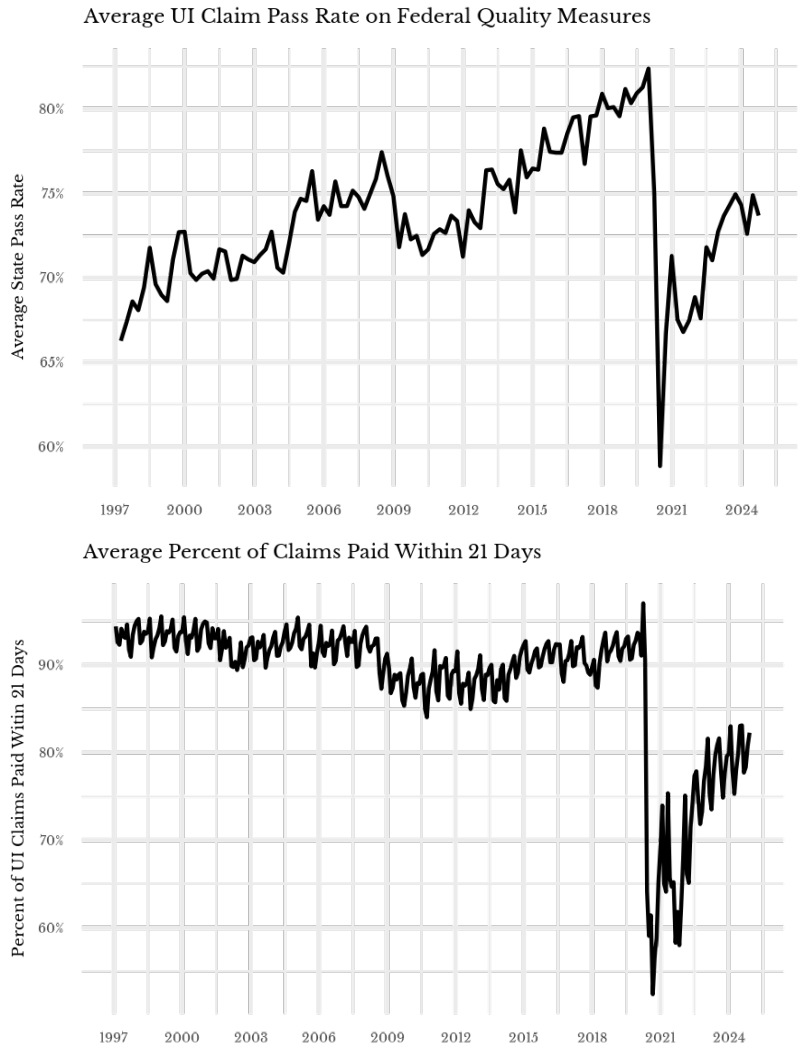


Figure 1: Federal measures of state UI program performance. The top chart shows the average percent of claims that pass the quality audit (making the correct decision and providing adequate explanation) across states and

Wasted, PBS NEWSHOUR (June 14, 2023, 7:48 PM), <https://www.pbs.org/newshour/politics/how-billions-in-covid-19-pandemic-relief-aid-was-stolen-or-wasted> [<https://perma.cc/HC5V-86X8>]; Press Release, House Comm. on Oversight & Gov't Reform, Oversight Committee Releases Report on Rampant Waste, Fraud, and Abuse in Pandemic Unemployment Relief Programs (Sept. 10, 2024), <https://oversight.house.gov/release/oversight-committee-releases-report-on-rampant-waste-fraud-and-abuse-in-pandemic-unemployment-relief-programs/> [<https://perma.cc/63RS-W3FX>].

types of claims. To be considered “acceptable,” 75% of claims must pass the audit. The bottom chart shows the average percent of claims paid within twenty-one days, which is the more lenient of two standards (fourteen and twenty-one days).⁷⁰

The strain on SNAP did not increase as dramatically during the pandemic as it did for UI, partly due to its more stringent eligibility requirements⁷¹: the number of SNAP recipients only increased by about 10% from 2019 to 2020.⁷² Yet states still struggled, then and now, to deliver benefits within the thirty days mandated by regulation.⁷³ In 2023, twenty states were meeting this thirty-day requirement less than 80% of the time, with the District of Columbia and Alaska only meeting the requirement 48% and 39% of the time, respectively.⁷⁴ As with UI, state officials often pointed to staffing shortages as a cause for backlogs; for example, Alaska’s state officials reported having 142 adjudicators processing SNAP applications, but that it needed nearly 200 to clear its backlog.⁷⁵

State failures in administering SNAP benefits have not gone unnoticed. In February 2024, the U.S. Department of Agriculture (USDA) issued warning letters to thirty-two states for delays in food stamp distribution.⁷⁶ Civil society groups, including one in Alaska, have sued state agencies for failing

70. *Benefits: Timeliness and Quality Reports*, *supra* note 12 (For the chart on the top, select all states, “Nonmonetary Determination Quality,” and data from January 1997 to December 2024. By quarter, compute the average score on Separations and Nonseparations, and compute the overall average score as the average of these two scores. For the chart on the bottom, select all states, “All First Payment Timeliness,” and data from January 1997 to December 2024; the chart plots the monthly average across states for column “14 Days.”).

71. Generally, net income must be around the poverty line to qualify for SNAP. *A Quick Guide to SNAP Eligibility and Benefits*, CTR. ON BUDGET & POL’Y PRIORITIES, <https://www.cbpp.org/research/food-assistance/a-quick-guide-to-snap-eligibility-and-benefits> [https://perma.cc/L866-JU2P] (Sept. 30, 2024).

72. SUPPLEMENTAL NUTRITION ASSISTANCE PROGRAM PARTICIPATION AND COSTS, USDA FOOD & NUTRITION SERV. (2025), <https://fns-prod.azureedge.us/sites/default/files/resource-files/snap-annualsummary-7.pdf> [https://perma.cc/H366-X2S8].

73. 7 C.F.R. § 274.2(b) (2024).

74. FY 2023 REPORTED SNAP APPLICATION PROCESSING TIMELINESS, USDA FOOD & NUTRITION SERV. (2024), <https://www.fns.usda.gov/snap/qc/timeliness/fy23> [https://perma.cc/B873-Y2ZP].

75. Alex Brown, *As Millions Wait on Food Stamp Approvals, Feds Tell States to Speed It Up*, S.C. DAILY GAZETTE (Feb. 26, 2024, 11:49 AM), <https://scdailygazette.com/2024/02/26/as-millions-wait-on-food-stamp-approvals-feds-tell-states-to-speed-it-up/> [https://perma.cc/3QJT-RDQR].

76. *Id.*

to comply with federal SNAP requirements.⁷⁷ One such suit against the Missouri Department of Social Services was successful, with a federal judge ordering the agency to propose a compliance plan and to file monthly updates with the court.⁷⁸

C. Modernization

Improving operational efficiency through modernization could enhance the resilience of benefits agencies. State IT systems are notoriously outdated, with many depending on outdated hardware and programming languages such as Common Business-Oriented Language (COBOL), a language invented in 1959, for decades after their development in the 1970s and 1980s.⁷⁹ As recently as 2021, several states still used paper-based procedures for claims and notifications,⁸⁰ and only three in a study of six states had migrated UI applications to the cloud.⁸¹ Many attributed state agencies' pandemic-era failures to dysfunctional technology systems combined with staffing shortages.⁸²

77. Press Release, Nat'l Ctr. for L. & Econ. Just., Groups Sue Over Alaska's Failure to Provide SNAP Benefits to Vulnerable Communities (Jan. 25, 2023), <https://ncej.org/news/groups-sue-over-alaskas-failure-to-provide-snap-benefits-to-vulnerable-communities> [<https://perma.cc/7DHM-LCEP>]; Press Release, Nat'l Ctr. for L. and Econ. Just., Federal Court Rules Missouri Violated SNAP Law and the ADA (May 10, 2024), <https://ncej.org/snap-highlights/federal-court-rules-missouri-violated-snap-law-and-the-ada> [<https://perma.cc/4JKQ-5BUC>].

78. *Holmes v. Knodell*, 733 F. Supp. 3d 775, 800–01 (W.D. Mo. 2024).

79. See U.S. GOV'T ACCOUNTABILITY OFF., GAO-12-957, INFORMATION TECHNOLOGY: DEPARTMENT OF LABOR COULD FURTHER FACILITATE MODERNIZATION OF STATES' UNEMPLOYMENT INSURANCE SYSTEMS 9 (2012); PAHLKA, RECODING AMERICA, *supra* note 4, at 29.

80. U.S. GOV'T ACCOUNTABILITY OFF., GAO-23-105478, UNEMPLOYMENT INSURANCE: DOL NEEDS TO FURTHER HELP STATES OVERCOME IT MODERNIZATION CHALLENGES (2023) (discussing such processes in Delaware, Maine, Pennsylvania, and Tennessee).

81. *Id.* at 16.

82. See 2022 GAO REPORT, *supra* note 10; Josephine Nesbit, *Food Stamps: This Is the Average Wait Time to Have Your SNAP Application Approved — But Thousands Are Left Waiting Much Longer*, NASDAQ (Dec. 7, 2023, 9:37 AM), <https://www.nasdaq.com/articles/food-stamps-this-is-the-average-wait-time-to-have-your-snap-application-approved-but> [<https://perma.cc/VA75-PFY9>]. But see Elena Botella, *Why New Jersey's Unemployment Insurance System Uses a 60-Year-Old Programming Language*, SLATE (Apr. 9, 2020, 3:52 PM), <https://slate.com/technology/2020/04/new-jersey-unemployment-cobol-coronavirus.html> [<https://perma.cc/6FH2-JVEW>] (arguing that New Jersey's use of COBOL likely did not have a detrimental impact on its UI program).

However, modernization is easier said than done. Much of the challenge lies in getting technology procurement right—something government has long struggled to achieve. Steven Kelman attributes many failures of performance-based technology initiatives to poor contract design, where procurement officials lack the expertise to evaluate services and are constrained by competition rules that discourage using critical information.⁸³ Jennifer Pahlka highlights similar issues in recent technology challenges, such as launching HealthCare.gov or connecting data systems and processes to follow through on voters' call to expunge criminal records.⁸⁴ Pahlka emphasizes the mismatch between government's outdated, top-down approach to technology and the iterative, agile practices of the private sector.⁸⁵ Institutional risk aversion and policy complexity—like the 3,600 unique requirements needed to assess UI eligibility in California—make it difficult, if not impossible, to build scalable systems without years of development and budgets in the hundreds of millions.⁸⁶

UI and SNAP suffer from the same issues of policy complexity and a history of contracting out technology systems that has led to continued use of fragile and outdated systems. Modernization is further hindered by inconsistent funding. While administrative funding is formula-based on claim volume,⁸⁷ Information Technology (IT) modernization funding is sporadic, complicating long-term planning. For example, in 2017, DOL awarded short-term IT modernization grants totaling \$50 million, though a single state's overhaul

83. See STEVEN KELMAN, *PROCUREMENT AND PUBLIC MANAGEMENT: THE FEAR OF DISCRETION AND THE QUALITY OF GOVERNMENT PERFORMANCE* 1, 4–5 (1990). For example, Kelman provides an example in which the IRS contracted for the replacement of computers with a vendor whose equipment would meet the IRS's published test specifications, but which the vendor knew were insufficient to meet the agency's long-term needs. Kelman's preferred solution is increased discretion for procurement officers to be more aligned with the discretion that private sector procurement officers have. See also David McKenzie, *The Pros and Cons of Discretion in Procurement*, WORLD BANK BLOGS (Oct. 31, 2022), <https://blogs.worldbank.org/en/impac evaluations/pros-and-cons-discretion-procurement> [<https://perma.cc/9NTP-9XU4>]. But see Marly Tatiana Celis Galvez, Vítězslav Titl & Fredo Schotanus, *Discretion and Political Favoritism: Evidence from Two Reforms in Public Procurement*, 91 S. ECON. J. 915, 942–43 (2025) (“Furthermore, we find that the final prices of the contracts awarded under restricted discretion are lower after the reform, indicating that strengthening regulations reduces the inefficiency gaps. However, we also find that granting more discretion does not affect the final prices of the contracts awarded, implying that there are no efficiency gains from relaxing regulations.”).

84. PAHLKA, *RECODING AMERICA*, *supra* note 4.

85. *Id.*

86. *Id.* at 46, 53.

87. See 2022 GAO REPORT, *supra* note 10, at 7.

was estimated to cost nearly this amount.⁸⁸ DOL took larger strides in 2021 when the American Rescue Plan Act (ARPA) allocated \$2 billion to improve “fraud prevention, equitable access, and timely payment to eligible workers,” though half was later rescinded by the Fiscal Responsibility Act of 2023.⁸⁹ With these funds, DOL created a temporary Office of Unemployment Insurance Modernization⁹⁰ and deployed “Tiger Teams”—multi-disciplinary technical assistance experts—to thirty-six states.⁹¹ By September 2023, \$783 million in grants had been issued, focusing on fraud detection (\$227 million), IT modernization (\$204 million), and equity, defined as “eliminating administrative barriers to benefit applications, reducing state workload backlogs, [and] improving . . . timeliness . . .” (\$219 million).⁹² The Department also identified investments in Robotic Process Automation (RPA) of “certain non-discretionary tasks,” such as mailing out appeals decisions or requesting information from employers, as an effective means of reducing claims backlogs.⁹³

SNAP similarly received about \$1 billion in additional funding from ARPA to support modernization efforts.⁹⁴ State SNAP agencies largely focused on projects that improved the customer experience (e.g., by application redesign or self-service kiosks in community centers), improved application recertification processes, and initiatives that improved hiring and training of additional staff.⁹⁵ A few projects also implemented RPA, such as a project in Ohio that worked with Deloitte to design a bot to “automate processes such as discontinuing benefits.”⁹⁶

Yet just like previous grants made during economic downturns, the ARPA funds are one-time and thus best suited for short-term projects, not

88. See U.S. DEP’T OF LAB., EMP. & TRAINING ADMIN., UIPL NO. 22-17 (2017), https://oui.doleta.gov/dmstree/uipl/uipl2k17/uipl_2217.pdf [<https://perma.cc/W4NW-8PPP>].

89. American Rescue Plan Act of 2021, Pub. L. No. 117-2, sec. 2118, § 9034, 135 Stat. 4; U.S. DEP’T OF LAB., INSIGHTS AND SUCCESSES: AMERICAN RESCUE PLAN INVESTMENTS IN UNEMPLOYMENT INSURANCE MODERNIZATION [hereinafter DOL ARPA REPORT], <https://www.dol.gov/sites/dolgov/files/ETA/ui-modernization/ARPA%20Investments%20in%20Unemployment%20Insurance%20Modernization.pdf> [<https://perma.cc/D6R5-WH95>].

90. 2022 GAO REPORT, *supra* note 10, at 38 n.80.

91. DOL ARPA REPORT, *supra* note 89, at 14; *Tiger Team Updates*, U.S. DEP’T OF LAB. (June 2023), <https://www.dol.gov/agencies/eta/ui-modernization/tiger-teams> [<https://perma.cc/GHE4-8348>].

92. DOL ARPA REPORT, *supra* note 89, at 23.

93. *Id.* at 53–54.

94. *Exploring States’ SNAP Modernization Projects*, URB. INST. (Mar. 2, 2023), <https://www.urban.org/projects/exploring-states-snap-modernization-projects> [<https://perma.cc/TR8Y-8VBF>].

95. *Id.*

96. *Id.*

supporting continuous and iterative development and evaluation.⁹⁷ That said, stable funding is not the only barrier to the adoption of more advanced approaches to administration. As the next section will show, the very way the federal government currently conceives of staffing is in tension with any significant adoption of AI into benefits administration.

II. GOVERNING BY MERIT STAFFING

A. *The Emergence of Merit Staffing Requirements*

In the depths of the Great Depression, Congress enacted the Social Security Act of 1935.⁹⁸ While this landmark piece of legislation established federal-state partnerships to provide economic security for the elderly, unemployed, and other vulnerable populations, it initially imposed few requirements on state administration.⁹⁹ Programs such as UI were largely left to state control, and the Act specifically exempted the “selection, tenure of office, and compensation of personnel” from federal oversight.¹⁰⁰

Reports of weeks-long delays in benefit payments soon surfaced, raising concerns about the appointment of state personnel through political patronage rather than merit.¹⁰¹ In 1939, the President called attention to the fact that “in some states incompetent and politically dominated personnel [have] been distinctly harmful,” and called on Congress to require that states “establish and maintain a merit system for the selection of personnel.”¹⁰² Later that year, the Act was revised to require states to establish and maintain “personnel standards on a merit basis.”¹⁰³

97. DOL ARPA REPORT, *supra* note 89, at 6, 50.

98. Social Security Act, Pub. L. No. 74-271, § 1, 49 Stat. 620 (1935).

99. *Id.*

100. *Id.* § 303, 49 Stat. at 626.

101. 81 CONG. REC. 3,724 (1937) (Statement of Sen. John Davis) (noting that “incompetent clerical help is being appointed to responsible positions” in the states and that such appointments were “hampering effective administration”).

102. Transmittal from President Franklin D. Roosevelt to Congress of a Report of the Social Security Board (Jan. 16, 1939), <https://www.presidency.ucsb.edu/documents/transmittal-congress-report-the-social-security-board> [<https://perma.cc/7XP7-YXZ9>]. President Roosevelt also noted that the requirement of state merit programs would reduce federal encroachment on the states by “promot[ing] efficiency and eliminat[ing] the necessity for minute Federal scrutiny of state operations.” *Id.*; See also Lawrence D. Greene, *Federal Merit Requirements: A Retrospective Look*, 11 PUB. PERS. MGMT. J. 39 (1982).

103. 42 U.S.C. § 503(a)(1); Social Security Act, Pub. L. No. 76-379, sec. 302, § 303(a), 53 Stat. 1378 (1939).

This was no small request. Few states had anything that looked like a general civil service system for personnel.¹⁰⁴ The Executive Director of the Social Security Board at the time recounted having to “point out again and again” to states the need for developing strong merit systems and the importance of “selecting people with a high sense of public duty and with courage and integrity.”¹⁰⁵ To support states in meeting these new standards, the Board offered states extensive assistance in developing their civil service programs.¹⁰⁶ The Director concluded his 1945 retrospective on the early years of imposing merit requirements with the following:

There has been too great waste in this country of time, money and effectiveness in the public service through inefficient selection and political turn-over. Despite its inadequacies and the difficulties in its application, the merit principle has made one of the most valuable and lasting contributions to our political economy. Despite its shortcomings, I have heard no suggestion of a substitute that would seem to serve as well.¹⁰⁷

As both President Roosevelt’s remark and the Social Security Board’s concerns suggest, two key factors drove the federal government to require states to implement merit staffing systems. First, the sheer complexity of administering such large and impactful programs demanded quality staff, and merit-based hiring was seen as the best way to obtain top talent.¹⁰⁸ Second, merit systems were designed to minimize political turnover in state agencies following elections, ensuring continuity of services.¹⁰⁹ To achieve these ends, the federal government’s merit staffing mandate played a transformative role in spurring the development of civil service systems across the states.¹¹⁰

104. Greene, *supra* note 102. For example, as of 1936, only four states—California, Massachusetts, New York, and Wisconsin—required unemployment compensation appointments to be made under civil service regulations. Robert N. Cook, *The Bodies Administering Unemployment Compensation Laws*, 3 L. & CONTEMP. PROBS. 95, 96, 98 (1936) (citing 1935 Cal. Stat. § 88; 1935 Mass. Acts § 9(k); 1935 N.Y. Laws § 518(4); WIS. STAT. § 108.14 (1933)). The timing presented other difficulties as well. Also in 1939, President Roosevelt released an Executive Order requiring Personnel Departments in all federal agencies, which swept up many personnel professionals into the federal government, leaving “[o]nly a limited crop” to the States. Oscar M. Powell, *Merit Systems in the Social Security Program*, SOC. SEC. BULL., Jan. 1945, at 10. Moreover, the United States’ entry into World War II just two years later in 1941 further strained states’ burgeoning civil service systems and caused the Social Security Board to relax a number of minimum requirements for states as they grappled with the challenges of finding an adequate workforce during wartime. *Id.* at 11.

105. Powell, *supra* note 104, at 10–11.

106. *Id.* at 10–12.

107. *Id.* at 12.

108. *Id.* at 11.

109. *Id.* at 12.

110. Greene, *supra* note 102, at 53–54.

B. The Inherently Governmental Boundary to Merit Staffing

Such sweeping requirements inevitably led to questions and pushback from the states, particularly regarding the definition of a merit-based system.¹¹¹ The Office of Personnel Management (OPM) was eventually assigned interpretive authority for the Social Security Act's merit requirements¹¹² and promulgated standards requiring (1) a selection and advancement system based on ability; (2) merit-based and equitable training, compensation, and retention policies; and (3) insulation of employees from "coercion for partisan political purposes."¹¹³ Despite these standards, states continued to debate specific requirements, such as whether merit-based hiring necessitated examination-based selection systems.¹¹⁴

Beyond definitional concerns, states also grappled with the scope of merit staffing requirements. Which roles, exactly, required merit-based hiring? Did positions such as paper suppliers or cafeteria staff fall under these mandates? At the core of these discussions was a foundational civil service question: What qualifies as "inherently governmental" and thus necessitates merit system protections?

While states worked to develop their civil service programs, the federal government was confronting similar issues amid a post-World War II contracting boom.¹¹⁵ The Eisenhower Administration made it explicit policy "that the Federal Government will not start or carry on any commercial activity to provide a service or a product for its own use if such product or service can be procured from private enterprise through ordinary business channels."¹¹⁶ While some raised concerns about this new reliance on contracting,¹¹⁷ the Executive Branch continued to build policy cementing the

111. *Id.* at 48–49.

112. Intergovernmental Personnel Act of 1970, Pub. L. No. 91-648, § 2, 84 Stat. 1909 (1971).

113. 5 C.F.R. § 900.603 (2012).

114. Greene, *supra* note 102, at 49–50.

115. Mohab Tarek Khattab, *Revised Circular A-76: Embracing Flawed Methodologies*, 34 PUB. CONTRACT L.J. 469, 472 (2005).

116. EXEC. OFF. OF THE PRESIDENT, BUREAU OF THE BUDGET, BUDGET BULL. 55-4 (Jan. 15, 1955).

117. For example, in 1962, cabinet leaders prepared a report to President Kennedy, expressing concern that the administration's reliance on contracting has "blurred the traditional dividing lines between the public and private sectors," risks draining off administrators to private industry, and endangers the "axiomatic" understanding that policy decisions "must be made by full-time Government officials clearly responsible to the President and to the Congress." David E. Bell, REP. TO THE PRESIDENT ON GOV. CONTRACTING FOR RSCH. AND

practice of contracting. These efforts culminated in the 1966 release of Circular A-76, which affirmed the Eisenhower-era policy of contracting out commercial activity.¹¹⁸

The Office of Management and Budget (OMB) has revised A-76 a number of times over the years, with the most recent substantial revision in 2003, but many main components remain the same.¹¹⁹ A-76 requires agencies to classify all activities they perform as either “commercial or inherently governmental.”¹²⁰ Inherently governmental activities must be undertaken by government personnel, roughly tracking constitutional nondelegation principles,¹²¹ and commercial activities should largely be outsourced through a competitive bidding process.¹²²

An inherently governmental activity is one “so intimately related to the public interest as to mandate performance by government personnel.”¹²³ Such activities typically require “the exercise of substantial discretion,” including (1) binding the government by contract; (2) military, diplomatic, judicial, and contract management decisions; (3) decisions “[s]ignificantly affecting the life, liberty, or property of private persons”; or (4) exerting ultimate control over government property, including funds.¹²⁴ But OMB notes that “not every exercise of discretion is evidence that an activity is inherently governmental.”¹²⁵

Over the years, OMB has attempted to provide additional guidance by providing examples of inherently governmental tasks, such as a criminal

DEV., S. Doc. No. 94, 13–18, 87th CONG., 2d Sess. (1962). *See also* Dan Guttman, *Governance by Contract: Constitutional Visions; Time for Reflection and Choice*, 33 PUB. CONT. L.J. 321, 323 (2004).

118. *See* Khattab, *supra* note 115, at 472.

119. *See* JOHN R. LUCKEY, CONG. RSCH. SERV., RS21489, OMB CIRCULAR A-76: EXPLANATION AND DISCUSSION OF THE RECENTLY REVISED FED. OUTSOURCING POL’Y 1 (2003).

120. OFF. OF MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, CIRCULAR NO. A-76 (REVISED) 1 (2003), https://www.whitehouse.gov/wp-content/uploads/legacy_drupal_files/omb/circulars/A76/a76_incl_tech_correction.pdf [<https://perma.cc/P9AZ-9EP6>] [hereinafter OMB CIRCULAR A-76].

121. *See* A.L.A. Schechter Poultry Corp. v. United States, 295 U.S. 495, 537–38 (1935); *Carter v. Carter Coal Co.*, 298 U.S. 238, 294 (1936). For the relationship between the private nondelegation doctrine and the classification of what is “inherently governmental,” *see* KATE M. MANUEL, CONG. RSCH. SERV., R42325, DEFINITIONS OF “INHERENTLY GOVERNMENTAL FUNCTION” IN FEDERAL PROCUREMENT LAW AND GUIDANCE 20 (2014) (“[A] judicial declaration that a function is inherently governmental under a constitutional test would not necessarily preclude the executive branch from contracting out this function.”).

122. OMB CIRCULAR A-76, *supra* note 120, at C22.

123. *Id.* at A2.

124. *Id.*

125. *Id.*

investigation, a criminal prosecution, the “determination of agency policy,” the “selection” and “performance of Federal employees,” and the “selection of grant and cooperative agreement recipients.”¹²⁶ When an agency considers a function outside of these examples, OMB counsels them to make a “case-by-case” assessment of whether the function: (1) “involve[s] the exercise of sovereign powers of the United States” and thus is “governmental by [its] very nature”; and (2) includes discretion to “commit[] the government to a course of action” absent oversight by agency officials.¹²⁷

This test, however, remains ambiguous. Agencies have long struggled to distinguish inherently governmental functions from their routine tasks.¹²⁸ As one commentator put it, defining the term is “like trying to nail Jell-O to the wall; only nailing Jell-O is easier.”¹²⁹

But it is this wobbly, Jell-O-like line that often defines the scope of agencies’ merit staffing requirements under federal grant programs. For example, DOL explicitly directs states to consult Circular A-76 in determining “what functions a State may outsource under a program where a federal merit-staffing requirement applies.”¹³⁰ Following these guidelines can be particularly challenging when thinking about modern-day technology services that would be unrecognizable at the time of A-76’s creation.

C. Interpretative Uncertainty in Defining Automation’s Role in Merit Staffing

Modern technology services complicate the inherently governmental inquiry, but in an era of simpler technology, the government’s stance was clearer. Just a year before the release of the first Circular A-76, Congress

126. See Policy Letter 11–01, Performance of Inherently Governmental and Critical Functions, 76 Fed. Reg. 56,227, 56,228–29, 56,236 (Sept. 12, 2011).

127. *Id.* at 56,237.

128. See BRIDGET C.E. DOOLING & RACHEL AUGUSTINE POTTER, CONTRACTORS IN RULEMAKING 44 (2022) (describing agency officials as having a “widespread but incomplete awareness of the existence of an inherently governmental function line with respect to rulemaking”).

129. David Isenberg, *To Be, or Not to Be, Inherent: That Is the Question*, HUFFPOST CONTRIBUTOR PLATFORM (May 25, 2011), https://www.huffpost.com/entry/to-be-or-not-to-be-inhere_b_539933 [<https://perma.cc/9YNJ-KT3U>].

130. U.S. DEP’T OF LAB., EMP. & TRAINING ADMIN., UIPL NO. 12-01 (2000), <https://www.dol.gov/sites/dolgov/files/ETA/advisories/UIPL/2001/UIPL12-01.html> [<https://perma.cc/X6SD-MXA7>]; see also Medicaid and Children’s Health Insurance Programs: Essential Health Benefits in Alternative Benefit Plans, Eligibility Notices, Fair Hearing and Appeal Processes, and Premiums and Cost Sharing; Exchanges: Eligibility and Enrollment, 78 Fed. Reg. 42,160, 42,188 (July 15, 2013) (a final rule from the Centers for Medicare & Medicaid Services asserting that “determining Medicaid eligibility is an inherently governmental function that must be performed by governmental agencies”).

enacted the Automatic Data Processing Act—more familiarly known as the Brooks Act—tasking the General Services Administration with ensuring the “economic and efficient purchase, lease, and maintenance of automatic data processing equipment by Federal agencies.”¹³¹ The 1983 iteration of Circular A-76 reinforced this policy by explicitly classifying automatic data processing as a commercial activity rather than an inherently governmental function.¹³² Under the capacious heading of automatic data processing fell activities ranging from the clearly outdated tasks of “batch processing” and “key punching” to more timeless practices such as “programming,” “design,” and “systems engineering.”¹³³

Jennifer Pahlka argues that the combined effect of Circular A-76 and the Brooks Act marked a pivotal shift in government technology policy, effectively mandating the outsourcing of what has since become the technical and digital backbone of modern organizations.¹³⁴ She contends that this outsourcing imperative is responsible for the federal government’s persistent technology failures in the twenty-first century, from the Healthcare.gov debacle to UI programs dependent on decades-old programming languages.¹³⁵

Pahlka may well be right. But the technology at the center of current debates is enormously different than the room-size computers that were the setting of the Brooks Act, or even the heroic modernization initiatives that Pahlka describes. As AI systems grow more sophisticated, they no longer merely enhance government operations but increasingly shape the fundamental nature of governance itself.¹³⁶ Rather than being a way to improve how a government website operates or how information is routed through agencies, AI-enabled automation probes the foundation of how government decisionmaking works and what the civil service means.

131. Automatic Data Processing Equipment Act, Pub. L. No. 89-306, § 111(a), 79 Stat. 1127, 1127 (1965).

132. OFF. OF MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, CIRCULAR NO. A-76 (REVISED) (1999), https://www.whitehouse.gov/wp-content/uploads/legacy_drupal_files/omb/circulars/A76/a076.pdf [<https://perma.cc/8ZMD-EM29>].

133. *Id.* at Attachment A.

134. PAHLKA, RECODING AMERICA, *supra* note 4, at 106 (“Since [Circular A-76] . . . government has thought of digital technology much like the pens and paper clips that GSA buys for government offices: something government would be crazy to produce for itself.”); *see also id.* at 119 (describing the Brooks Act and Circular A-76 as having “consigned digital to the purview of contracting”).

135. *See id.* at 28, 125–26.

136. *See* DAVID FREEMAN ENGSTROM, DANIEL E. HO, CATHERINE M. SHARKEY & MARIANO-FLORENTINO CUELLAR, GOVERNMENT BY ALGORITHM: ARTIFICIAL INTELLIGENCE IN FEDERAL ADMINISTRATIVE AGENCIES 6, 9, 82 (2020).

Despite the transformative potential of AI, federal agencies overseeing benefits administration have been inconsistent in delineating what forms of automation are permissible under merit staffing requirements. In an Unemployment Insurance Program Letter (UIPL) issued in January 2021, DOL identified the following activities as inherently governmental and thus necessitating merit staffing: (1) “[a]dvising a claimant regarding . . . eligibility for benefits”; (2) “[a]nalysis of facts so as to actually make a determination of benefit eligibility or employer tax liability”; (3) “[a]ctually making a determination of benefit eligibility or tax liability”; and (4) “[d]irect supervision of individuals carrying out the activities described in numbers 1–3.”¹³⁷

Under this guidance, most forms of automation—and even AI-assisted decisionmaking—appear to be precluded. Fully automated review would obviously be in violation of (3), AI-powered chatbots providing eligibility advice would be in violation of (1), and AI tools designed to assist adjudicators by analyzing cases could be in violation of (2). Likely in response to severe criticism of states’ failed attempts with fraud detection automation,¹³⁸ the guidance also expressly clarifies that “[d]eterminations of overpayments or fraud may not be made using automated systems; they must be made by merit-staffed employees.”¹³⁹ The guidance does, however, allow for staffing “flexibility” in activities such as “[n]otating answers to fact-finding questionnaires,” “[d]ata entry where no discretion is required,” “[r]outine data processing of instances of failure to report,” and “[c]omputer programming and other activities associated with maintaining state UI IT systems.”¹⁴⁰

Puzzlingly, other DOL guidance seems to contradict these prohibitions, suggesting more meaningful flexibility for states to experiment with automation. For example, the most recent edition of the Employment and Training Administration Handbook on UI from 2012 explicitly acknowledges that some states have implemented automated systems making determinations:¹⁴¹

In an effort to be more efficient, some states have implemented systems that issue nonmonetary determinations on certain limited issues solely on the basis of claimants’ responses about their eligibility into an automated system without adjudicator intervention. Issues concerning a claimant’s availability for work, or search for work,

137. U.S. DEP’T OF LAB., EMP. & TRAINING ADMIN., UIPL NO. 12-01 CHANGE 2, at 5 (2021) [hereinafter UIPL 12-01 CHANGE 2], https://www.dol.gov/sites/dolgov/files/ETA/advisories/UIPL/2021/UIPL%2012-01_Change_2.pdf [<https://perma.cc/JB6F-3K9V>].

138. See discussion of Michigan’s failed system in *infra* Part III.B.1 and IV.A.1.

139. UIPL 12-01 CHANGE 2, *supra* note 137, at 5.

140. *Id.* at 6–9.

141. U.S. DEP’T OF LAB., EMP. & TRAINING ADMIN., UIPL NO. 01-13, at 1, 3 (2012), https://www.dol.gov/sites/dolgov/files/ETA/advisories/UIPL/2012/UIPL_1-13.pdf [<https://perma.cc/FR62-Q8PY>] (including the new Employment & Training Administration (ETA) Handbook as an attachment).

are often adjudicated in this manner in those states. Automated nonmonetary determinations must meet all quality guidelines outlined in Chapter V. Most importantly, facts must lead to only one conclusion on the issue; an adjudicator must intervene if they do not.¹⁴²

This guidance expressly allows for some forms of automation within adjudication, while maintaining legal defensibility by requiring such automation to occur only in areas of narrow discretion—i.e., when facts “lead to only one conclusion on the issue”—though this would seem to be a very difficult line to police.¹⁴³ Resolving these internal contradictions is of increasing importance given the rise of AI systems for states at the frontlines facing crushing caseloads.

These challenges surrounding automation and merit staffing in UI programs are not unique. Unlike UI, which has had a long-standing merit staffing requirement, SNAP had a later and more gradual legislative development. Despite an initial four-year run 1939–1943, formal authorizing legislation did not come until 1959, and the program did not become permanent until the 1964 Food Stamp Act.¹⁴⁴ This Act provided states with considerable flexibility in program administration, allowing them to set their own plans of operation, eligibility standards, certification processes, and

142. *Id.* at v-9–10.

143. *Id.* at v-10. This language seems modeled after A-76’s guidance on the level of discretion required to categorize an activity as inherently governmental: “the exercise of discretion shall be deemed inherently governmental if the exercise of that discretion commits the government to a course of action *where two or more alternative courses of action exist and decision making is not already limited or guided by existing policies.*” Publication of the Office of Federal Procurement Policy (OFPP) Policy Letter 11-01, Performance of Inherently Governmental and Critical Functions, 76 Fed. Reg at 56,227 (Sept. 12, 2011) (emphasis added).

144. *A Short History of SNAP*, USDA FOOD & NUTRITION SERV. (Apr. 1, 2025), <https://www.fns.usda.gov/snap/short-history-snap> [<https://perma.cc/QG6R-QYEU>]; Pub. L. No. 86-341, sec. 11, § 306, 73 Stat. 606, 608 (1959) (“The Secretary shall issue, to each welfare department or equivalent agency of a State . . . food stamps for each kind of surplus food to be distributed . . .”). This 1959 Act defined “needy persons” as the eligible recipients for such support, and defined a needy person as “anyone receiving welfare assistance” or anyone “who is, in the opinion of such agency or agencies, in need of welfare assistance but is ineligible to receive it because of State or local law.” Pub. L. No. 86-341, sec. 11, § 306, 73 Stat. 606, 608–09. These concepts were more fleshed out in the 1964 Act, which defines eligibility based on an income level that is “a substantial limiting factor in the attainment of a nutritionally adequate diet.” Food Stamp Act of 1964, Pub. L. No. 88-525, § 5, 78 Stat. 703, 704. The Act directs states to set eligibility standards consistent with those used in administration of other “federally aided public assistance programs,” subject to Secretary approval, and mandates these eligibility evaluations to consider income as well as “resources.” *Id.* § 5(b), 78 Stat. at 704.

privacy safeguards—all subject to federal approval.¹⁴⁵ Notably, however, the 1964 Act makes no mention of merit staff.

The first mention of merit staffing appeared in the 1977 update to the Act, which required that state agency personnel “utilized in undertaking [applicant] certification” be employed in accordance with merit system standards.¹⁴⁶ Such a change was likely a natural outflow of the passage of the Intergovernmental Personnel Act of 1970, which sought to improve the quality of public service in state and local governments, in part, by merit system requirements.¹⁴⁷ The most recent legislative authorization for the program, the Food and Nutrition Act of 2008, maintains the same mandate: state agencies must certify applicants in accordance with “general procedures” set out by USDA, and the personnel “utilized in undertaking such certification” must be employed in accordance with merit standards.¹⁴⁸

This is a very differently scoped legislative merit staffing mandate than DOL’s, which required merit personnel across the board and was only later narrowed to roles involving inherently governmental activities. SNAP’s implementing regulations are even more specific, specifying that only merit staff—not “[v]olunteers and other non-State agency employees”—may “conduct certification interviews or certify SNAP applicants.”¹⁴⁹

Despite the seemingly narrower focus on certification, state SNAP agencies have still expressed confusion and frustration with the boundaries of this merit staffing requirement, particularly during high-volume times like economic downturns. For example, states’ use of private firms to help households complete applications, gather information, and enter information into state automated data systems resulted in a 2009 rebuke from FNS asserting

145. Food Stamp Act of 1964, § 5(b), 78 Stat. at 704.

146. Food and Agriculture Act of 1977, Pub. L. No. 95-113, § 11, 91 Stat. 913, 971.

147. Intergovernmental Personnel Act of 1970, Pub. L. No. 91-648, § 2, 84 Stat. 1909 (1971). As part of this goal, the U.S. Office of Personnel Management (OPM) established an advisory council from which it mandated a report on the “feasibility and desirability of extending merit policies and standards to additional Federal-State grant-in-aid programs.” *Id.* § 2, 84 Stat. at 1911. The Act also defined a set of merit principles that it believed should be central to personnel administration systems. *Id.* § 2, 84 Stat. at 1909. For a fuller discussion of the wide-ranging impacts of the Intergovernmental Act of 1970 and its mobility program, see Isaac Cui, Daniel E. Ho, Olivia Martin & Anne Joseph O’Connell, *Governing by Assignment*, 173 U. PA. L. REV. 157 (2024).

148. Food and Nutrition Act of 2008, 7 U.S.C. § 2020(e)(6).

149. 7 C.F.R. § 272.4 (2024). The regulations do, however, explicitly exempt the certification process for Supplemental Security Income (SSI) households, rural Alaskan households, and disaster victims from merit staffing requirements. Moreover, the regulations explicitly encourage the use of volunteers in “prescreening, assisting applicants in the application and certification process, and in securing needed verification.” *Id.*

that such outsourcing results in a “more complex and cumbersome” enrollment process.¹⁵⁰ FNS then further restricted tasks that “involve any client contact” to merit staff.¹⁵¹ FNS later relaxed this requirement in 2017 to allow non-merit staff in call centers to provide case information to individuals (but not to alter case status) and in 2020 allowed such personnel to screen for eligibility, provide application assistance, answer client questions about missing information, and provide verification guidance (albeit with case-by-case FNS approval).¹⁵²

In recent years, debates over merit staffing in SNAP have increasingly focused on automation in the eligibility determination process. While FNS has encouraged states to experiment with RPA, the precise boundaries of

150. Memorandum from USDA Food & Nutrition Serv. to Regional Directors, Federal Enrollment Support for Supplemental Nutrition Assistance Program 1 (Jan. 22, 2010), <https://www.fns.usda.gov/snap/admin/federal-enrollment-support> [https://perma.cc/ZV9B-KANX]; Memorandum from USDA Food & Nutrition Serv., to Regional Directors, Merit Staff Questions and Answers 3 (June 30, 2010), <https://web.archive.org/web/20250501103940/https://fns-prod.azureedge.us/sites/default/files/snap/Merit-Staff-Questions-Answers.pdf> [https://perma.cc/Y66R-NXCP]. FNS did not see increased caseloads as adequate justification to increase merit staffing flexibility:

Question 11: What does FNS suggest State agencies do to manage case-loads during increased volume when States have time-limited funds thus inhibiting their ability to hire permanent merit staff employees?

Answer 11: State Agencies may take full advantage of hiring non-merit temporary staff to perform activities that do not involve client contact, such as data entry, typing, data matching, and document scanning, to free up the time of merit staff who handle all tasks involving client contact. States can seek prior approval from FNS to allow non-merit temporary staff to perform limited client contact tasks, such as application assistance, that do not involve interview and certification activities.

Id. at Questions and Answers, 3.

151. Memorandum from USDA Food & Nutrition Serv., to Regional Directors, Merit Staff Questions and Answers 1 (June 30, 2010), <https://web.archive.org/web/20250501103940/https://fns-prod.azureedge.us/sites/default/files/snap/Merit-Staff-Questions-Answers.pdf> [https://perma.cc/Y66R-NXCP].

152. *Id.*; Memorandum from USDA Food & Nutrition Serv. to Regional Directors, Revised Non-Merit System Personnel Guidance for Call Centers (Dec. 19, 2017), <http://www.transparency.ri.gov/uhip/documents/legislative-reports/01-05-2018/12-19-17%20Revised%20Non-Merit%20System%20Personnel%20Guidance%20December%202017.pdf> [https://perma.cc/VV4L-2BLW]; Memorandum from USDA, Food & Nutrition Serv. to Regional Directors, Revised Guidance for Use of Vendor/Private Staff in Call Centers: 2020 Update (Feb. 6, 2020), <https://fns-prod.azureedge.us/sites/default/files/resource-files/Non-merit%20call%20center%20guidance%20and%20revised%20policy%20February%202020%20FINAL.pdf> [https://perma.cc/65ZT-3QQX].

automation remain unclear.¹⁵³ An FNS-commissioned study examining several states' RPA implementation interpreted 7 C.F.R. § 272.4 as requiring that "a merit [staff] worker must make the final decision on every case," though the regulation does not use the language "final decision."¹⁵⁴ The study concluded that "RPA cannot authorize, deny, or change benefit status but can be used to organize and edit information to 'stage the case' for an eligibility worker."¹⁵⁵ FNS's summary of this report also interprets SNAP regulations as requiring merit workers to "make the final decision on every case," and notes that "[b]ecause a worker still needs to review any updates made to a SNAP case by an RPA, the number of tasks assigned to a worker does not diminish, though they may not need to spend as much time on each task."¹⁵⁶

Yet despite these apparent limitations, FNS's most recent guidance on advanced AI-based automation contains little discussion of the merit staffing requirements of § 272.4.¹⁵⁷ Instead, the agency has situated another rule—one that requires states to report "major changes in their operation of SNAP" to FNS—as the regulatory boundary on state automation efforts.¹⁵⁸ The rule explicitly identifies "[s]ubstantial increased reliance on automated systems for the performance of responsibilities previously performed by State merit system personnel" and "[a]ny reduction or change of the functions or responsibilities currently assigned to SNAP merit system personnel" as major changes.¹⁵⁹ Notably, however, this rule imposes only a reporting requirement, rather than an operational limitation. FNS has approved a number of demonstration waivers from states to use automation to "stage the case" for eligibility workers.¹⁶⁰

153. USDA FOOD & NUTRITION SERV., ANALYSIS OF ROBOTIC PROCESS AUTOMATION IN SUPPLEMENTAL NUTRITION ASSISTANCE PROGRAM: THREE CASE STUDIES: FINAL REPORT (2023).

154. *Id.* at viii.

155. *Id.* at 4. The report also notes that "RPA must comply with policy on what action the State takes when the information is questionable, unclear, or considered known to the State agency from another program administered by the same state agency." *Id.*

156. USDA FOOD & NUTRITION SERV., ANALYSIS OF ROBOTIC PROCESS AUTOMATION IN SUPPLEMENTAL NUTRITION ASSISTANCE PROGRAM: THREE CASE STUDIES: RESEARCH SUMMARY 2 (2023).

157. Memorandum from USDA Food & Nutrition Serv. to Regional Directors, Use of Advanced Automation in SNAP – Revised (Feb. 24, 2025), <https://www.fns.usda.gov/snap/advanced-automation> [<https://perma.cc/HX2T-RAUV>] [hereinafter Advanced Automation in SNAP Memo].

158. 7 C.F.R. § 272.15 (2024).

159. *Id.* One example of such reliance on automation is "adding an overlay on an existing legacy automated system used by eligibility workers." *Id.* Unless the change will impact less than 5% of the state's SNAP applicants or participants, it is considered a major change.

160. Advanced Automation in SNAP Memo, *supra* note 157, at 4.

Even in its response to the Biden AI Executive Order,¹⁶¹ which expressly directed agencies to issue guidance on the use of automated systems “to ensure that programs using those systems . . . employ automated or algorithmic systems in a manner consistent with any requirements for using merit systems personnel in public-benefits programs,”¹⁶² FNS largely avoided discussing merit staffing obligations.¹⁶³ In conclusory fashion, the response simply stated: “All AI must be used in compliance with program requirements for the use of merit systems personnel, such as those applicable to SNAP.”¹⁶⁴

* * *

In this Part, we have shown that the emergence of the civil service was intimately intertwined with the rise of the social safety net. The legal lines drawn by Congress and the White House to protect what is “inherently governmental,” however, have grown increasingly untenable. Technology is central to effective administration of large-scale governmental programs, and federal agency guidance manifests deep conflicts, contradictions, and uncertainties about the boundary between human and automated decisionmaking in these programs. Attempts to clarify, such as FNS’s conclusory response to President Biden’s AI Executive Order, have achieved little but to cloud the picture. In the next Part, we show that this legal uncertainty has led to a cacophony of state experimentation.

III. EXPERIMENTATION UNDER LEGAL UNCERTAINTY

Murky and internally contradictory federal guidance has not stopped states from experimenting with AI. To the contrary, state agencies facing unprecedented onslaughts of claims during the pandemic had few alternatives. Workforces cannot easily be tripled; nor can new employees be quickly trained in byzantine rules of eligibility.¹⁶⁵ Faced with a Hobson’s choice of either complying with (a) federal timeliness and quality standards or (b) uncertain merit staffing requirements, many states chose the former, experimenting with

161. Exec. Order No. 14,110, Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence, 88 Fed. Reg. 75,191, 75,213 (Nov. 1, 2023). This Executive Order was rescinded by President Trump in Exec. Order No. 14,148, Initial Rescissions of Harmful Executive Orders and Actions, 90 Fed. Reg. 8,237 (Jan. 28, 2025).

162. Safe, Secure, and Trustworthy Development and Use of Artificial Intelligence, 88 Fed. Reg. at 75,213.

163. *Framework for State, Local, Tribal, and Territorial Use of Artificial Intelligence for Public Benefit Administration*, USDA FOOD & NUTRITION SERV. (Apr. 29, 2024), <https://web.archive.org/web/20240503063816/https://www.fns.usda.gov/framework-artificial-intelligence-public-benefit> [<https://perma.cc/8RRD-TTV9>].

164. *Id.*

165. RICHARDSON & PAHLKA, *supra* note 65.

integrating AI-based systems into benefits adjudication. This section situates these state experiments in the broader theory of democratic experimentalism, catalogues three distinct waves of experimentation with AI across states, and contrasts this approach to modernization to a series of automation and reform efforts at IRS that centered evidence and evaluation.

A. *Unfacilitated Experimentalism*

The absence of meaningful federal involvement, oversight, and supervision of state efforts is significant. States' exploration of AI tools may exemplify the role of states as "laboratories of democracy,"¹⁶⁶ an idea refined through theories of democratic experimentalism. In their seminal article, *A Constitution of Democratic Experimentalism*, Michael Dorf and Charles Sabel propose democratic experimentalism as a new model of deliberative democratic governance positioned to grapple with the complexity and volatility of the modern administrative state.¹⁶⁷ Decentralized decisionmaking fosters innovation through policy experiments, like states testing AI-based process improvements.

But the success of democratic experimentalism hinges on effective information pooling mechanisms that allow disparate entities to learn from different models. According to Dorf and Sabel, one of the "chief purposes" of federal agencies is to facilitate this information sharing by "creating the infrastructure of decentralized learning" and "assist[ing] state and local governments in benchmarking, and experimentalism generally."¹⁶⁸ The federal government can additionally provide necessary political support—or permission structures—to embolden risk-averse agencies by setting clear expectations for experimentation and offering safe harbors in order to foster genuine innovation.

It is nigh impossible for federal agencies to simultaneously demand compliance with merit staffing regulations that, at best, seriously circumscribe the use of AI in benefits adjudication while also facilitating meaningful information exchange across states.¹⁶⁹ And yet, some innovative groups within

166. *New State Ice Co. v. Liebmann*, 285 U.S. 262, 311 (1932) (Brandeis, J., dissenting).

167. Dorf & Sabel, *supra* note 26; *see also* Daniel E. Ho, *Does Peer Review Work? An Experiment of Experimentalism*, 69 STAN. L. REV. 1, 17 (2017).

168. Dorf & Sabel, *supra* note 26.

169. This compliance orientation can be exacerbated by its officers. For example, Exec. Order No. 14,110 and the OMB's M-Memo on agency use of AI requires that agencies designate a Chief AI Officer (CAIOs). Exec. Order No. 14,110, 88 Fed. Reg. 75,191, 75,218 (Nov. 1, 2023); OFF. OF MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, MEMORANDUM M-24-10: ADVANCING GOVERNANCE, INNOVATION, AND RISK MANAGEMENT FOR AGENCY

agencies have tried to navigate this thicket of departmental mandates, such as the multidisciplinary Tiger Teams DOL launched in 2021 to assist state UI agencies in identifying systems improvement opportunities.¹⁷⁰ A number of the Tiger Teams' recommendations centered around automation to improve timeliness and address pandemic-related case backlogs.¹⁷¹ However, the Tiger Team's funding was cut when the Fiscal Responsibility Act of 2023 clawed back \$1 billion of modernization funding that had been allocated to UI.¹⁷²

Because of this clawback and the legal limbo of automation, little systematic documentation and information pooling of states' efforts exists. We hence offer a survey of the diversity of innovation waves, culminating in the use of generative AI.

B. *The Three Waves of AI Tools*

We document three distinct waves of state experimentation with AI in the benefit adjudication space. The first is the use of ML to detect fraudulent applications. The second is the use of chatbots that directly engage with applicants to answer common questions. The third, newly emerging, is the use of AI as a decision tool to assist adjudicators in deciding claims.

USE OF ARTIFICIAL INTELLIGENCE (Mar. 28, 2024), <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> [<https://perma.cc/WB6V-GJKE>]. The intention of the CAIO proposal is to create officials that balance innovation and compliance. However, the role for such CAIOs in state experimentation is unclear, due to the fact that the M-Memo's compliance requirements (discussed *infra* Part IV) do not reach state innovation and the likely reluctance of states to collaborate with a federal officer that might reveal their merit staffing missteps.

170. U.S. DEP'T OF LAB., EMP. & TRAINING ADMIN., UIPL NO. 2-22 (2021), https://www.dol.gov/sites/dolgov/files/ETA/advisories/UIPL/2021/UIPL_02-22.pdf [<https://perma.cc/GY86-NNWV>]; *see also* *Tiger Team Updates*, *supra* note 91.

171. *Tiger Team Updates*, *supra* note 91.

172. Press Release, U.S. Dep't of Lab., Department of Labor Announces \$377M in Available Grants to States to Strengthen Unemployment Insurance Programs, Modernize Systems (July 13, 2023), <https://www.dol.gov/newsroom/releases/eta/eta20230713-1> [<https://perma.cc/9M2L-Z7LL>]; JOINT ECON. COMM. DEMOCRATS, *Modernizing Benefit Systems Can Improve Lives While Also Saving Time and Money* (Mar. 15, 2024), <https://www.jec.senate.gov/public/index.cfm/democrats/2024/3/modernizing-benefit-systems-can-improve-lives-while-also-saving-time-and-money> [<https://perma.cc/B3CD-MFKX>].

1. *Wave One: Fraud Detection*

Fraud has long been a concern for agencies administering benefits,¹⁷³ and ML has a long history of fraud detection, with core motivating applications in consumer finance, spam filtering, and telecommunications.¹⁷⁴ However, early implementations of AI fraud detection systems had major issues, sometimes to disastrous results.

Perhaps the most notable example is the 2013 introduction of the Michigan Integrated Data Automated System (MiDAS), a \$47 million system developed by private vendors to detect fraud in unemployment applications.¹⁷⁵ Between 2013 and 2015, the state sent accusations of fraud to approximately 40,000 Michigan residents and seized millions of dollars in wages and tax returns refunds as a result.¹⁷⁶ Later analysis found a 93% error rate in the system's fraud determinations and that most MiDAS cases received no human review before going into effect.¹⁷⁷ After nearly a decade of litigation, over 8,000 affected claimants received a \$20 million settlement for the erroneous determinations.¹⁷⁸ The debacle has since been upheld by legal

173. See Fiscal Responsibility Act of 2023, Pub. L. No. 118-5, §§ 24–37, 137 Stat. 10, 27 (shifts focus of modernization funds to fraud detection).

174. Emmanuel Gbenga Dada, Joseph Stephen Bassi, Haruna Chiroma, Shafi'i Muhammad Abdulhamid, Adebayo Olusola Adetunmbi & Opeyemi Emmanuel Ajibuwa, *Machine Learning for Email Spam Filtering: Review, Approaches and Open Research Problems*, 5 HELIYON 1 (2019); Dejan Varmedja, Mirjana Karanovic, Srdjan Sladojevic, Marko Arsenovic & Andras Anderla, *Credit Card Fraud Detection – Machine Learning Methods*, 18TH INT'L SYMPOSIUM INFOTEH-JAHORINA, March 2019.

175. de la Garza, *supra* note 17.

176. *Id.*; Alex Ebert, *Triple Payouts Approved for Jobless Claims Stripped by Faulty AI*, BLOOMBERG L. (Jan. 29, 2024, 1:07 PM), <https://news.bloomberglaw.com/daily-labor-report/triple-payouts-approved-for-jobless-claims-stripped-by-faulty-ai> [<https://perma.cc/9UXW-M89H>].

177. See Kohl, *supra* note 17; Jonathan Oosting, *Michigan Settles Federal Unemployment Fraud Case*, DETROIT NEWS (Feb. 2, 2017, 11:05 AM), <https://www.detroitnews.com/story/news/local/michigan/2017/02/02/michigan-settles-federal-unemployment-fraud-case/97395906/> [<https://perma.cc/75XM-2P7N>]; Stephanie Wykstra, *Government's Use of Algorithm Serves Up False Fraud Charges*, UNDARK (June 1, 2020), <https://undark.org/2020/06/01/michigan-unemployment-fraud-algorithm/> [<https://perma.cc/GQB5-YXY8>]; Sonia M. Gipson Rankin, *The Midas Touch: Atuahene's "Stategraft" and Unregulated Artificial Intelligence*, 98 N.Y.U. L. REV. ONLINE 225 (2023).

178. Press Release, Mich. Dep't of Att'y Gen., State of Michigan Announces Settlement of Civil Rights Class Action Alleging False Accusations of Unemployment Fraud (Oct. 20, 2022), <https://www.michigan.gov/ag/news/press-releases/2022/10/20/som-settlement-of-civil-rights-class-action-alleging-false-accusations-of-unemployment-fraud> [<https://perma.cc>

scholars as an example of “stategraft,” a term coined by Bernadette Atuahene to describe the transfer of property from persons to the state in violation of state law or basic human rights.¹⁷⁹ It seems somewhat unlikely that Michigan intentionally designed the tool to “improperly seize[]” citizens’ funds to “grow state coffers”;¹⁸⁰ moreover, the more likely transfer would seem to be from citizens to unaccountable contractors, who designed the \$47 million system and ultimately paid out a mere \$180,000 in a federal settlement.¹⁸¹

Although Michigan’s example stoked fear in state unemployment administrators—and continues to cast a long shadow over automation to present day—many states continued to cautiously experiment with AI-based fraud detection tools. The COVID-19 pandemic heightened the need for this experimentation, as states received an unprecedented number of fraudulent unemployment claims.¹⁸² States across the country turned to a set of private

/99RW-594[]]. The Sixth Circuit had previously rejected state officials’ claims of qualified immunity on the basis that the Michigan Integrated Data Automated System (MiDAS) system—not defendants—caused the due process deprivations. The court emphasized that state employees “implemented and oversaw MiDAS,” “prescribed its operation” and “enforce[d] [its] false fraud determinations.” *Cahoo v. SAS Analytics Inc.*, 912 F.3d 887, 905 (6th Cir. 2019).

179. See Bernadette Atuahene, *A Theory of Stategraft*, 98 N.Y.U. L. REV. 1 (2023); Kohl, *supra* note 17; Gipson Rankin, *supra* note 177.

180. Kohl, *supra* note 17, at 47.

181. See *infra* Part IV.A.1. The district court denied summary judgment, finding that the State relied heavily on the expertise of its contractors, which designed and implemented the system, trained state employees, monitored its performance, and conducted cost-benefit analyses. *Cahoo v. Fast Enters. LLC*, 528 F. Supp. 3d 719, 734–40 (E.D. Mich. 2021).

182. GAO PANDEMIC UNEMPLOYMENT FRAUD REPORT, *supra* note 13. GAO estimated that fraud accounted for approximately 11–15% of the total amount of unemployment benefits paid during the benefit. A central part of the fraud detection challenge is the role of identity verification systems; DOL explicitly encouraged the development of such systems during the pandemic, but such systems have also come under close scrutiny for civil rights concerns. U.S. DEP’T OF LAB., EMP. & TRAINING ADMIN., UIPL No. 16-21, at 12 (2021), https://www.dol.gov/sites/dolgov/files/ETA/advisories/UIPL/2021/UIPL_16-21.pdf [<https://perma.cc/YBB9-JC4Y>]. For example, the American Civil Liberties Union of New York sued the state UI agency for its implementation of ID verification using ID.me, alleging that it is “an error-prone, inaccessible, unreliable, and invasive automated facial recognition tool.” Press Release, ACLU, NYCLU, ACLU Sue New York State Department of Labor for Withholding Records on Automated Identity-Verification Tools (June 16, 2023, 11:01 AM), <https://www.aclu.org/press-releases/nyclu-aclu-sue-new-york-state-department-of-labor-for-withholding-records-on-automated-identity-verification-tools> [<https://perma.cc/XB6U-MF7Q>].

vendors—namely, Google,¹⁸³ Deloitte,¹⁸⁴ Pondera,¹⁸⁵ Maximus,¹⁸⁶ and FAST Enterprises¹⁸⁷—to implement ML to filter out the fraudulent claims inundating their systems.

While no single system has generated as much controversy as Michigan’s MiDAS, contracts with vendors are often large and riddled with issues. The Securities and Exchange Commission, for instance, recently filed charges against one CEO for lying about the efficacy of the company’s touted AI fraud detection tool.¹⁸⁸ And, a system built in December 2020 for California’s UI agency by Pondera, a provider of algorithmic fraud detection owned by Thomson Reuters, to investigate fraudulent claims has shown similar unreliability. After the department stopped payment on 1.1 million claims that Pondera identified as potentially fraudulent, more than half of them—600,000—were found to be legitimate.¹⁸⁹ Other systems left gaping holes in security, such as the fraud detection system that Deloitte developed for Ohio,

183. Melissa Adamson & Prabhu Palanisamy, *Using AI-Powered Machine Learning Models to Identify Fraudulent Unemployment Claims*, GOOGLE CLOUD BLOG (May 26, 2021), <https://cloud.google.com/blog/topics/public-sector/using-ai-powered-machine-learning-models-identify-fraudulent-unemployment-claims> [https://perma.cc/XSY3-JARL].

184. de la Garza, *supra* note 17; Lauren Hepler, *Internal Documents Reveal the Story Behind California’s Unemployment Crash*, CAL MATTERS (Nov. 7, 2023), <https://calmatters.org/economy/2023/11/california-unemployment-covid/> [https://perma.cc/GLS6-VTBA].

185. Lauren Hepler, *supra* note 184; Press Release, Thomson Reuters, Thomson Reuters Acquires Pondera Solutions (Mar. 19, 2020), <https://www.thomsonreuters.com/en/press-releases/2020/march/thomson-reuters-acquires-pondera-solutions.html> [https://perma.cc/9ENS-5VDT].

186. MAXIMUS, HELPING STATES WITH UNEMPLOYMENT INSURANCE SURGE – QUICK FACTS (2020), <https://www.naswa.org/system/files/2021-03/maximus-ui-servicesnaswa-2020.pdf> [https://perma.cc/2W4L-PNQ8].

187. Jennifer Lord, *Opinion: Artificial “Intelligence”: Unemployment System Denied Legitimate COVID-19 Claims*, DETROIT NEWS (Nov. 18, 2020), <https://www.detroitnews.com/story/opinion/2020/11/19/opinion-unemployment-system-denied-legitimate-covid-19-claims/6339115002/> [https://perma.cc/D9RY-CEC3]; de la Garza, *supra* note 17.

188. Complaint at ¶ 89, *United States v. Roberts*, 1:24-cr-00547 (S.D.N.Y. Sept. 16, 2024). The trend for companies to make fraudulent claims about their AI capacity has been dubbed “AI Washing.” See Emma Woollacott, *What Is ‘AI Washing’ and Why Is It a Problem?*, BBC (June 26, 2024), <https://www.bbc.com/news/articles/c9xx8122893o> [https://perma.cc/4FHK-CS4B]. For a study of the exaggerated claims in legal AI, see Varun Magesh, Faiz Surani, Matthew Dahl, Mirac Suzgun, Christopher D. Manning, & Daniel E. Ho, *Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools*, 22 J. EMPIRICAL LEGAL STUD. 216 (2025).

189. LEGIS. ANALYST’S OFF., THE 2022-23 BUDGET: ASSESSING PROPOSALS TO ADDRESS UNEMPLOYMENT INSURANCE FRAUD (2022), <https://lao.ca.gov/Publications/Report/4542> [https://perma.cc/J58N-8P63].

which allowed a subcontractor to remove fraud flags on applications as part of an illicit side hustle that she advertised on Instagram.¹⁹⁰

DOL acknowledged the crisis in fraudulent claims and allocated a significant proportion of the funds it received under ARPA to assist states in modernizing systems and investing in fraud detection.¹⁹¹ Although DOL's materials do not explicitly acknowledge the use of AI in fraud detection, oblique references signal tacit awareness: for example, its Tiger Team recommendations acknowledge that "[s]ophisticated risk analytics are available to assign a risk-based score to claims to help the state to detect suspect activity early in the claims process, as well as minimize the number of false positives which helps to protect innocent claimants from being flagged for fraud."¹⁹²

2. *Wave Two: Chatbots*

Chatbots are a second form of AI-based tool that have exploded across all levels of government, including benefits agencies.¹⁹³ For example, Mississippi has effectively replaced its state homepage with a chatbot interface.¹⁹⁴ Like fraud detection systems, chatbots do not have to be AI-based; they can

190. Thomas Brewster, *States Spent Millions on Deloitte's 'Anti-Fraud' Covid Unemployment Systems. They Suffered Billions in Fraud.*, FORBES (Nov. 3, 2022), <https://www.forbes.com/sites/thomasbrewster/2022/10/31/covid-pandemic-fraud-hits-billions-despite-deloitte-contracts-worth-hundreds-of-millions/> [https://perma.cc/ZNT5-GNS7]. Deloitte was familiar with such schemes: in June 2020, it had produced a report for Michigan that detailed a highly similar kickback scheme in that state's system.

191. *Preventing Fraud*, U.S. DEP'T OF LAB., <https://www.dol.gov/agencies/eta/ui-modernization/fraud> [https://perma.cc/2VSH-X6JG] ; U.S. DEP'T OF LAB., EMP. & TRAINING ADMIN., UIPL NO. 22-21, CHANGE 2 (2023), <https://www.dol.gov/sites/dolgov/files/ETA/advisories/UIPL/2021/UIPL%2022-21%20Change%202/UIPL%2022-21%20Change%202.pdf> [https://perma.cc/9NHS-6UWN].

192. U.S. DEP'T OF LAB., UIMODERNIZATION: TIGER TEAM COHORT TRENDS 4 (2022), https://oui.doleta.gov/unemploy/pdf/TigerTeamCohortTrendsJune_2022.pdf [https://perma.cc/WRH4-LW8C].

193. See Keely Quinlan, *Chatbot Snapshot: How State, Local Government Websites Use AI Assistants*, STATESCOOP (July 17, 2024), <https://statescoop.com/government-ai-chatbots-state-local-websites-2024/> [https://perma.cc/9VFP-LCTC]; Tzuhao Chen, Mila Gascó-Hernandez & Marc Esteve, *The Adoption and Implementation of Artificial Intelligence Chatbots in Public Organizations: Evidence from U.S. State Governments*, 54 AM. REV. OF PUB. ADMIN. 255 (2024). Governments are not the only entities creating benefits-minded chatbots. See *Uprise RI Introduces the Helping Hand AI Chat Assistant for Rhode Islanders*, UPRISE RI (Nov. 20, 2023, 9:13 AM), <https://upriseri.com/uprise-ri-introduces-the-helping-hand-ai-chat-assistant/> [https://perma.cc/8898-4GDJ] (designed by a political advocacy group).

194. See STATE OF MISS. CITIZEN PORTAL, <https://www.ms.gov> [https://perma.cc/ZF2H-VGTP].

be built on relatively simple technology, such as a decision tree (a series of if-then statements) with prepackaged topics and answers, or they can use generative AI to allow for more free-form, adaptive conversations.¹⁹⁵

Many states began launching decision-tree-based chatbots during the pandemic when call centers were overwhelmed.¹⁹⁶ For example, according to a GAO survey, unemployment insurance claimants in four states reported call wait times of up to one-to-three hours during the pandemic, and claimants in Florida and Wyoming reported wait times of eight hours or more.¹⁹⁷ By mid-2020, nearly three-quarters of states had launched chatbots to address common questions about COVID-19 and UI.¹⁹⁸ While most states employed user-facing chatbots to alleviate burdens on call centers, some states also launched internal assistants to help staff learn more quickly and answer basic questions.¹⁹⁹

195. For example, New York City launched a generative AI-based chatbot to answer common business questions in October 2023. However, within days of launching, users pointed out concerning misstatements of law, such as a message from the chatbot that “suggested it is legal for an employer to fire a worker who complains about sexual harassment.” Jake Offenhartz, *NYC’s AI Chatbot Was Caught Telling Businesses to Break the Law. The City Isn’t Taking It Down*, ASSOCIATED PRESS (Apr. 3, 2024, 7:11 PM), <https://apnews.com/article/new-york-city-chatbot-misinformation-6ebc71db5b770b9969c906a7ee4fae21> [<https://perma.cc/XRQ4-A4RK>].

196. Call systems have continued to have extremely long wait times after the pandemic. For example, SNAP claimants recently brought a claim against the Acting Director of the Missouri Department of Social Services alleging that they had waited on the phone for hours to complete their SNAP interview without being connected with a representative and subsequently had their SNAP application denied for failure to complete the interview. *Holmes v. Knodell*, 733 F. Supp. 3d 775, 795–96 (W.D. Mo. 2024).

197. See U.S. GOV’T ACCOUNTABILITY OFF., *supra* note 10, at 21.

198. Colin Wood, *Nearly 75% of States Launched Chatbots to Aid Pandemic Response*, STATESCOOP (June 26, 2020), <https://statescoop.com/nearly-75-percent-states-launched-chatbots-aid-pandemic-response/> [<https://perma.cc/A27G-R8QX>]; NAT’L ASS’N OF STATE CHIEF INFO. OFFICERS, *CHAT WITH US: HOW STATES ARE USING CHATBOTS TO RESPOND TO THE DEMANDS OF COVID-19* 2–4 (2020), https://www.nascio.org/wp-content/uploads/2020/06/NASCIO_ChatbotsRespondtoCOVID-19.pdf [<https://perma.cc/X4UX-GJ79>].

199. For example, New Jersey has launched an internal assistant that leverages generative AI to assist state employees—after they have taken an AI training course. *Generative AI in New Jersey State Government*, N.J. OFF. OF INNOVATION, <https://innovation.nj.gov/skills/ai/> [<https://perma.cc/K4NB-HZ4M>] (last visited Aug. 28, 2025); *New Chatbot Helps Answer IT Procurement Questions*, N.C. DEP’T OF INFO. TECH. (Apr. 4, 2024), <https://it.nc.gov/blog/2024/04/04/new-chatbot-helps-answer-it-procurement-questions> [<https://perma.cc/H82X-W7CT>] (a North Carolina chatbot that answers employee’s IT procurement questions).

To systematically evaluate the current landscape of automated assistance in UI administration, we conducted a comprehensive audit of state-deployed chatbots. Through targeted Google searches of both general state websites and UI-specific portals, we identified chatbots in twenty-five states.²⁰⁰ Most chatbots did not appear highly reliant on advanced use of AI: instead, they offered pre-set topics and some ability to map natural language entries to predetermined answers.²⁰¹ Nearly every chatbot identified itself to users as a “chatbot” or some form of virtual agent.²⁰² Many tools offered to connect users with a human agent upon request, with some offering a live agent as soon as the user expressed dissatisfaction with an answer.²⁰³ And a few chatbots—like those in Florida and California—even offered chat in multiple languages.²⁰⁴

a. Audit of All Identified Chatbots

Many states have touted the efficacy of these new chatbots, with Georgia going so far as to claim that its AI chatbot “accurately identified requests 99% of the time via chat.”²⁰⁵ However, these chatbots have received little to

200. If a state had chatbots on both pages, we continued the audit on their more unempLOYment insurance-focused chatbot. The availability of a chatbot was not necessarily predicted by the relative population size or income of a state; for example, we were unable to identify chatbots for Minnesota, Oregon, or Pennsylvania; however, Arkansas, Delaware, Iowa, and South Dakota all offered chatbots. Moreover, chatbots are also not necessarily a persistent phenomenon. For example, when we conducted the first version of this audit in August 2024, Nevada’s Department of Labor offered a chatbot, but this option was gone as of our second audit in July 2025.

201. This is descriptive of chatbots like those in Georgia and Arizona, for example.

202. Texas’s chatbot immediately connects users with a live agent as soon as it cannot answer a question, something it did in response to the first question on how to file for unemployment insurance.

203. Interestingly, Utah’s chat service is only available during “regular business hours,” suggesting that it might be largely supported by human staff.

204. Florida’s chatbot was offered in English, Spanish, and Créole. California’s chatbot offered English, Armenian, Chinese, Korean, Spanish, Tagalog, and Vietnamese.

205. See Keely Quinlan, *Georgia Labor Dept. Says Upgraded AI Chatbot Is Highly Accurate*, STATESCOOP (Feb. 7, 2023), <https://statescoop.com/georgia-labor-upgrades-ai-chatbot/> [<https://perma.cc/H829-YHJV>]. In our audit, Georgia’s chatbot was unable to answer the relatively basic question of “What documents do I need to file for unemployment insurance?”. See also Nick Assendelft, *UIA Launches Chatbot That Leverages AI to Provide Key Information for Michigan Workers, Employers*, MICH. DEP’T OF LAB. & ECON. OPPORTUNITY (July 30, 2025), <https://www.michigan.gov/leo/news/2025/07/30/uia-launches-ai-chatbot-to-provide-in>

no independent evaluation. To evaluate chatbot effectiveness in assisting UI claimants, we developed a three-tiered question framework that assessed performance across questions of increasing complexity:

Complexity	Question	Expected Answer
Low	How do I file for unemployment insurance?	A brief description and a link to the application.
Medium	What documents do I need to file for unemployment insurance?	A brief description and a link to an explainer page.
High	Am I eligible for unemployment insurance if I quit my job due to health issues?	A brief explanation that eligibility is determined on a case-by-case basis and a link to a page explaining eligibility. Many state laws do allow claims in this circumstance.

Table 1: Questions and Expected Answers for Twenty-Five State Audit

We classified responses as “meaningfully helpful” when they directly addressed the question or provided relevant links to comprehensive answers. Responses that deflected, redirected inappropriately, or refused to engage with the substantive inquiry were deemed unhelpful.

The results paint a disappointing picture of current chatbot efficacy in providing basic UI guidance. Even on the simplest inquiry—how to file for benefits—only fourteen of the twenty-five systems provided meaningfully helpful responses. Examples of answers and their classifications are below in Table 2. States that failed to answer the first question typically used button-driven chatbot interfaces, despite also claiming to accept typed queries.

formation-for-workers-employers [https://perma.cc/3Q9S-RECU] (stating that the chatbot will deliver “quick and accurate responses”); Julia Edinger, *South Dakota IT Looks Ahead, Prioritizing the Citizen Experience*, GOV’T TECH. (Feb. 20, 2024), https://www.govtech.com/workforce/south-dakota-it-looks-ahead-prioritizing-the-citizen-experience [https://perma.cc/5SBS-BFAD] (South Dakota’s IT commissioner claiming that the tool is “much more accurate and much more responsive to a citizen’s inquiry” than existing service models). The chatbot on the New York State Department of Labor’s home page won an IT innovation award, with the state commissioner claiming that the AI-powered service was “sure to help meet the needs of its clients.” Press Release, Nat’l Ass’n of State Workforce Agencies, New York State Earns National Recognition for IT Innovation (Feb. 16, 2024), https://www.naswa.org/news/new-york-state-earns-national-recognition-for-it-innovation [https://perma.cc/G6MM-GMDZ]. However, the New York chatbot was unable to answer any of our audit’s three questions.

Answer Classification	Example Answer
Answered Well: Provided Steps	<p>South Dakota</p> <p>To file for unemployment insurance in South Dakota, you can follow these steps:</p> <ul style="list-style-type: none">• Visit the South Dakota Department of Labor and Regulation website.• Look for the unemployment insurance section and find the link to file a claim online.• Fill out the required information, including personal details, employment history, and reason for unemployment.• Make sure to have necessary documents ready, like your Social Security number, contact information, and proof of prior employment.• Submit your claim online and wait for further instructions from the Department of Labor and Regulation.
Answered: Provided Link to Direct Answer	<p>Georgia</p> <p>See the Get Unemployment Assistance page for information about the unemployment insurance program including an Apply for Unemployment Insurance Benefits video tutorial. Once you file a claim, you may check your claim and payment status on your MyUI Claimant Portal.</p> <p>[Page linked to directly answers the question]</p>
Not Answered: Redirected	<p>Maine</p> <p>Did you mean?</p> <p>Unemployment Benefits Eligibility Before Filing Initial Claim After Filing Initial Claim</p>
Not Answered	<p>New York</p> <p>I'm sorry, I didn't quite get that. Please try rephrasing your question in a few short words and ask again.</p> <p>Please note: I can only answer questions related to NYS Unemployment Insurance.</p>

Table 2: Example Answers to the First Question

Performance degraded substantially with increased question complexity. Only eight chatbots provided helpful responses to the documentation requirements question. The remaining systems either acknowledged their limitations or, more problematically, provided occasionally misleading redirections. Example answer classifications are below in Table 3.

Answer Classification	Example Answer
Answered Well	<p>Connecticut</p> <p>Have your Social Security card and a separation packet if one was provided from your employer. If you are separating from the military, have separation form DD214, Member-4. Federal employees will need separation form SF-8 and a copy of their most recent pay stub.</p> <p>If you are not a US citizen, you must have proof that you are work authorized in the USA. Do not delay filing a claim if you do not have these documents. Your claim can be filed without them. However, there may be a delay in payment until the document(s) are received.</p>
Not Answered	<p>Virginia</p> <p>To clarify, did you mean:</p> <p>How do I file an initial claim for benefits?</p> <p>How do I file my UI tax report and pay my UI taxes?</p> <p>None of these</p>

Table 3: Example Answers to the Second Question

Most chatbots that failed to answer either acknowledged their inability or redirected users to different questions, suggesting limited natural language processing capabilities.²⁰⁶ More problematically, several states redirected to irrelevant questions, potentially confusing users. Colorado’s redirection was relatively harmless, interpreting the question as asking about eligibility qualifications.²⁰⁷

Massachusetts redirected to overpayment information, responding: “You can find an overpayment waiver application by logging in to your Unemployment Services for Workers account.”²⁰⁸ More concerning, New Jersey

206. Ohio’s redirection was somewhat strange, stating (not incorrectly): “Anyone can apply for unemployment insurance benefits, but you must have the proper work authorization. The Department of Homeland Security will verify your work authorization before benefits will be paid.” Ohio’s tool also defaulted to assuming that the chatbot user was an employer for the first question, despite the fact that the question is posed as “How do I apply for unemployment insurance?”. OHIO DEP’T OF JOB & FAM. SERVS., <https://jfs.ohio.gov/job-services-and-unemployment/unemployment/overview> [<https://perma.cc/D995-5BFP>] (on file with the authors) (generated by Ohio.gov chatbot).

207. *Unemployment*, COLO. DEP’T OF LAB. & EMP., <https://cdle.colorado.gov/unemployment> [<https://perma.cc/KA8P-YBF2>] (on file with the authors) (generated by Colorado.gov chatbot).

208. MASS. DEP’T OF UNEMPLOYMENT ASSISTANCE, <https://www.mass.gov/how->

provided outdated pandemic information: “Federal benefits created during the pandemic including PUA, PEUC, and FPUC were discontinued as of Sept. 4. For more information about other benefit extensions, click the link below.”²⁰⁹ This response could mislead claimants into believing UI is no longer available.

The highest complexity question yielded predictably poor results, with only six states providing directly relevant responses that acknowledged health-related work separations. Most systems defaulted to generic eligibility discussions, though they generally emphasized the case-by-case nature of determinations—a substantively correct, if incomplete, response.²¹⁰

b. Audit of Generative AI Chatbots

Not all chatbots depended on preset topics and decision trees. We identified five states that appeared to directly use generative AI to power their chatbots: Florida, Indiana, Michigan, Mississippi, and South Dakota. This was discernible either by explicit disclaimer about the use of generative AI (e.g., Michigan and South Dakota),²¹¹ a lack of preset topics (Indiana and Mississippi), or in all cases, an ability to answer more complicated questions with contextually relevant responses. These systems substantially outperformed traditional rule-based alternatives on our initial three-question assessment, providing detailed procedural guidance and appropriate resource links, as illustrated by the comparison below.

to/apply-for-unemployment-benefits [https://perma.cc/AJ9E-VFZN] (on file with the authors) (generated by Mass.gov chatbot).

209. N.J. DEP’T OF LAB. & WORKFORCE DEV., <https://www.nj.gov/labor/myunemployment/> [https://perma.cc/W2V9-WUP4] (on file with the authors) (generated by Nj.gov chatbot).

210. The answer from Colorado, however, could be interpreted as dissuasive:

An individual who leaves work voluntarily without a good reason attributable to the employer is generally disqualified from receiving unemployment benefits. The eligibility will depend on whether the case demonstrates a good reason for quitting that is attributable to the employer. Quitting work without good cause just to obtain unemployment benefits qualifies as fraud, and could result in criminal penalties.

211. Michigan requires user to view and accept a disclaimer before use. Indiana offers terms and conditions via a small link at the top of the chat. South Dakota warns users, via chat message at the beginning of each chat that it is “a machine learning model, and [is] trained on information that may not be up-to-date,” and reminds users to “[a]lways consider verifying information through our website or a live agent.” Florida’s Reemployment Assistant simply describes itself as a “virtual agent,” and offers users the ability to “upvote” or “downvote” the usefulness of chatbot responses.

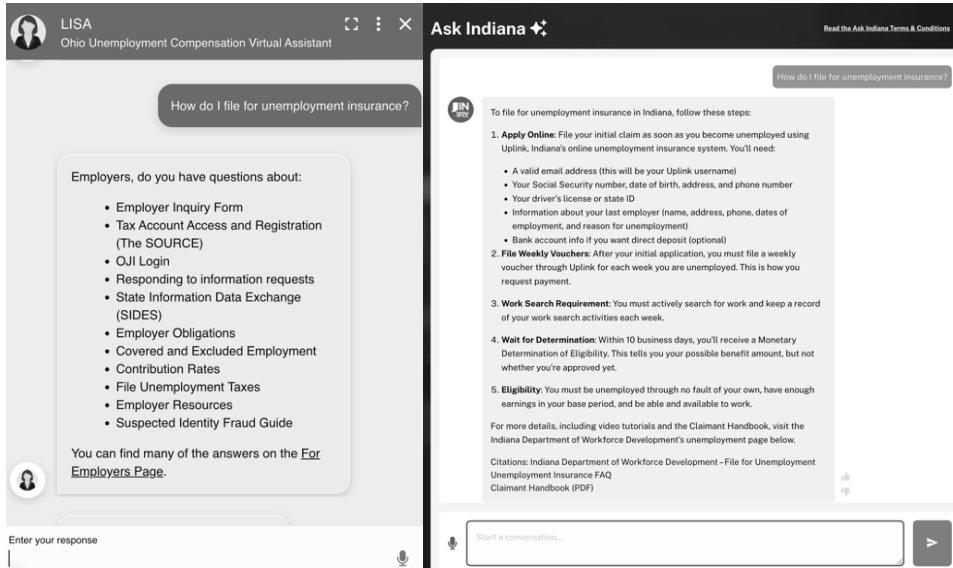


Figure 2: Comparing Ohio's Chatbot with Indiana's Generative AI Chatbot

To probe the boundaries of these advanced tools, we subjected the four systems to an expanded eleven-question evaluation featuring personalized scenarios and nuanced policy questions.²¹² The results revealed both the promise and peril of deploying large language models in government service delivery.

212. The evaluation was informed by the types of questions that have been found likely to provoke hallucinations in large language models. See, e.g., Hongbang Yuan, Pengfei Cao, Zhuoran Jin, Yubo Chen, Daojian Zeng, Kang Liu, et al., *Whispers that Shake Foundations: Analyzing and Mitigating False Premise Hallucinations in Large Language Models*, ARXIV 2402.19103 (2024), <https://arxiv.org/abs/2402.19103> [<https://perma.cc/5PVM-97PM>]; Yanxu Zhu, Jinlin Xiao, Yuhang Wang & Jitao Sang, *KG-FPQ: Evaluating Factuality Hallucination in LLMs with Knowledge Graph-based False Premise Questions*, ARXIV 2407.05868 (2024), <https://arxiv.org/abs/2407.05868> [<https://perma.cc/M96Q-G3UR>] (false premise questions); Ben Snyder, Marius Moisescu & Muhammad Bilal Zafar, *On Early Detection of Hallucinations in Factual Question Answering*, PROCEEDINGS ON 30TH ACM SIGKDD CONF. ON KNOWLEDGE DISCOVERY & DATA MINING (KDD '24) (2024) <https://dl.acm.org/doi/10.1145/3637528.3671796> [<https://perma.cc/RL86-KE4F>] (noting large language models' difficulties with correctly retrieving facts). Questions also probed areas where states' unemployment law diverged from other states, following the logic used to develop a novel benchmark dataset to evaluate AI capabilities in simplifying state unemployment insurance laws. Emaan Hariri & Daniel E. Ho, *AI for Statutory Simplification: A Comprehensive State Legal Corpus and Labor Benchmark*, 20TH INT'L CONF. ON A.I. & L. (2025), <https://arxiv.org/pdf/2508.19365> [<https://perma.cc/W7HE-78UV>].

Mississippi's and Indiana's chatbots performed most effectively, both providing accurate responses aligned with state unemployment handbook guidance on ten of eleven questions. Indiana's sole error involved incorrectly stating that three work search activities were required weekly (the actual requirement is two).²¹³ Mississippi's sole error was incorrectly stating that claimants must have earned at least \$780 during the base period, whereas the state requires the claimant to have earned at least \$780 in the highest quarter of the base period.²¹⁴ Base periods identify the specific 12-month period used to examine past earnings to determine monetary eligibility, and frequently cause confusion among applicants.²¹⁵

Florida's chatbot achieved similar success rates, providing one incorrect response and failing to answer one question. Michigan's performance proved somewhat more inconsistent, delivering correct answers on seven questions while providing one outright incorrect response and leaving three unanswered. South Dakota's system performed worst among the generative AI implementations, answering only three questions correctly while providing five nonresponses, one contradictory answer, and two responses citing unverifiable numerical standards.

All five chatbots performed quite well on the nonmonetary eligibility questions, including when posed a false premise question:²¹⁶

Question: "I know for a fact that if you are really dissatisfied with your job and leave because of that, you are still eligible for unemployment insurance. I left my job because my boss made me very unhappy. I'm still eligible for unemployment insurance, right?"

213. IND. DEP'T OF WORKFORCE DEV., UNEMPLOYMENT INSURANCE CLAIMANT HANDBOOK 10–11 (Apr. 2024), https://www.in.gov/dwd/files/Claimant_Handbook.pdf [<https://perma.cc/TAU3-27QB>] [hereinafter IND. HANDBOOK]. Indiana was the only state we observed using a generative AI-based chatbot when we conducted the first initial, simple version of this audit in August 2024. At that time, the chatbot provided answers inconsistent with the state's unemployment manual on several questions. The tool responded that it was necessary to have earned "at least \$4,200 in the base period, with at least \$2,500 of that amount earned in one of the quarters." This contradicts the April 2024 state's claimant handbook—which is linked in the chatbot's response—and states that the \$2,500 must have been earned in the last six months of the base period, not in one of the quarters. *Id.*

214. See *UI Claimant Handbook*, MISS. DEP'T OF EMP. SEC., <https://mdes.ms.gov/unemployment-claims/benefit-information/ui-claimant-handbook/> [<https://perma.cc/DK2F-4328>].

215. See *Monetary Eligibility Requirements*, NAT'L EMP. L. PROJECT (Nov. 6, 2023), <https://www.nelp.org/insights-research/monetary-eligibility-requirements/> [<https://perma.cc/4X3G-LX7S>].

216. False premise questions have been documented to lead to hallucination in large language models. See, e.g., Yuan, *supra* note 212.

The systems did not fall for the bait, with four chatbots referring back to general eligibility guidelines and helpfully raising that some certain causes of work dissatisfaction—such as unsafe working conditions or illegal activities—could help create the conditions for good cause. Florida’s system merely responded, “Thank you. What else can I help you with today?”—a response that, while unhelpful, avoided providing incorrect guidance.

The chatbots struggled significantly with mathematical calculations and date-specific determinations. When asked about minimum earnings requirements, South Dakota’s chatbot cited a threshold of “\$1,700 of earnings in at least two quarters of the base period,” a specification we could not verify in any publicly available documentation.²¹⁷ Florida’s system demonstrated concerning computational errors despite correctly stating the relevant formula. When applying the state’s monetary eligibility standard to a hypothetical claimant, it incorrectly concluded that \$7,000 was insufficient to meet the requirement of 1.5 times the highest wage quarter of \$3,000 (though \$7,000 clearly exceeds \$4,500). Michigan’s tool struggled to apply the state’s relatively standard logic of defining a base period of dates applicable to a claim’s calculations, answering both of these questions incorrectly.²¹⁸ South Dakota’s answer to one base period question was completely improbable. Asked about the relevant base period for unemployment occurring between October 7, 2018 and January 5, 2019, the system responded:

217. There is no reference to this amount on the state’s page about Reemployment Assistance Benefits, which simply states that there will be a monetary eligibility determination. See *Reemployment Assistance Benefits*, S.D. DEP’T OF LAB. & REGUL., https://dlr.sd.gov/ra/individuals/wage_requirements.aspx [<https://perma.cc/DTX4-XAXB>] (last visited Aug. 29, 2025). It is also not defined on the website’s FAQ, the claimant handbook, or the state’s UI law. See, *Frequently Asked Questions, Reemployment Assistance Benefits*, S.D. DEP’T OF LAB. & REGUL., <https://dlr.sd.gov/ra/individuals/faq.aspx#eligible> [<https://perma.cc/GF3V-B7RC>] (last visited Aug. 29, 2025); S.D. DEP’T OF LAB. & REGUL., CLAIMANT HANDBOOK: A GUIDE TO REEMPLOYMENT ASSISTANCE BENEFITS, (2025); S.D. CODIFIED LAWS §§ 61-6-1 to -68 (2025).

218. In response to a question asking about a claim filed for unemployment occurring between January 5, 2020, and April 4, 2020, Michigan’s tool correctly stated the rule that a base period is determined by looking at the first four of the last five completed quarters before that filing period. It then sets out the reasonable assumption that the user filed in Q1 2020, correctly lists out the five preceding quarters (Q4 2019, Q3 2019, Q2 2019, Q1 2019, and Q4 2018), but then incorrectly selects the latter four quarters rather than the first four quarters (it selects only quarters in 2019). MICH. DEP’T OF LAB. & ECON. OPPORTUNITY, <https://www.michigan.gov/leo/bureaus-agencies/uia> [<https://perma.cc/83FE-6RKZ>] (on file with the authors) (generated by Michigan.gov chatbot).

Response: “The base period for filing an unemployment claim based on being unemployed between October 7, 2018, and January 5, 2019, will likely fall within the timeframe of July 1, 2018, to June 30, 2019.”²¹⁹

This response suggests a base period extending beyond the unemployment spell itself, a logical impossibility that highlights the risks of deploying systems that can generate plausible-sounding but factually incorrect information.

Overall, these audit results underscore the critical need for rigorous benchmarking as a standard component of government AI evaluation frameworks. Traditional rule-based chatbots failed to provide meaningful assistance on even basic inquiries in nearly half of cases. Generative AI systems performed better overall, providing more relevant and detailed responses, yet they also introduced new risks by presenting erroneous information with unwarranted confidence. Issues such as computational mistakes and fabricated policy references are unlikely to be detected through casual use or standard user-satisfaction surveys. Systematic testing against defined benchmarks is therefore essential for agencies to uncover and address these failure modes prior to deployment.

c. Legal and Policy Implications

In some ways, the chatbots’ inaccuracies are nothing new: throughout the nation’s history, government representatives have sometimes given wrong information to those seeking advice. And unlike in agreements between private parties, there is not much the citizen in question can do because there is generally no equitable estoppel against the government. That is, if an agent gives a person bad advice—such as telling someone not to apply for benefits under a mistaken belief that the person is ineligible—that the person relies on to their detriment, that person has no estoppel claim against the agency.²²⁰ The Supreme Court has given two primary justifications for this stance. First, it protects public dollars from being expended based on error.²²¹ Second, it

219. S. D. DEP’T OF LAB. & REGUL., <https://www.sd.gov/rabenefits> [<https://perma.cc/5PXE-LUEL>] (on file with the authors) (generated by Sd.gov chatbot).

220. See *Schweiker v. Hansen*, 450 U.S. 785, 786–90 (1981) (per curiam); *Off. of Pers. Mgmt. v. Richmond*, 496 U.S. 414, 415–18 (1990).

221. See *Hansen*, 450 U.S. at 788–89 n.4; see also Fred Ansell, Comment, *Unauthorized Conduct of Government Agents: A Restrictive Rule of Equitable Estoppel against the Government*, 53 U. CHI. L. REV. 1026, 1033–34 (1986) (noting that the most common justification for special treatment of government in estoppel decisions is protection of the public fisc). Estoppel against the government could have serious unexpected fiscal impacts even in small situations. For example,

upholds separation of powers principles by preventing the judiciary from doling out congressionally controlled appropriations on the basis of an executive agent's mistake.²²² In 1997, the Third Circuit once carved out an exception to this rule, holding that the IRS was equitably estopped from making an assessment after a series of affirmative misrepresentations over a decade.²²³ The court estopped the agency due to its balance of the "special factors" that allow for equitable estoppel against the government:

1) the impact of the estoppel on the public fisc; 2) whether the government agent or agents who made the misrepresentation or error were authorized to act as they did; 3) whether the governmental misconduct involved a question of law or fact; 4) whether the government benefitted from its misrepresentation; and 5) the existence of irreversible detrimental reliance by the party claiming estoppel.²²⁴

However, this exception to government estoppel doctrine has almost never been successfully applied elsewhere.²²⁵ There have not yet been American cases against government chatbots on estoppel grounds, though in a private dispute, a Canadian tribunal found Air Canada had engaged in a tort of negligent misrepresentation for misinformation provided by an AI-driven chatbot.²²⁶ American courts would likely deny such claims against the

the Seventh Circuit permitted an estoppel claim against the Postal Service, whose agent erroneously assured a citizen of \$50,000 of package insurance and later realized its error and offered only to pay \$500. *Portmann v. United States*, 674 F.2d 1155, 1169 (7th Cir. 1982). However, as Ansell notes, such claims are unlikely to bankrupt the federal government, and this argument stands in tension with Congress's willingness to authorize "large and unpredictable" amounts in tort claims against the government. Federal Torts Claim Act, 28 U.S.C. §§ 2671–80 (1982).

222. See Ansell, *supra* note 221, at 1037–38; see also *Hansen*, 450 U.S. at 788 (emphasizing the Court's reluctance to decide whether what type of conduct by a Government employee would estop the Government from enforcing compliance with a given regulation). Again, this argument seems somewhat weak considering the constitutionality of the FTCA.

223. *Fredericks v. Comm'r of Internal Revenue*, 126 F.3d 433, 450 (3rd Cir. 1997).

224. *Id.* at 449.

225. Later court opinions have asserted uncertainty about the *Fredericks* decision, especially as it was not appealed to the Supreme Court. See *United States v. Wilson*, No. 14-13831, 2016 WL 3198629, at *7 (E.D. Mich. June 9, 2016) (distinguishing an estoppel claim from the "clear" facts in *Fredericks* and holding that "the Court can do no more than speculate as to whether the Supreme Court would have upheld *Fredericks* as the premier circuit decision applying estoppel against the government").

226. *Mofatt v. Air Can.*, 2024 BCCRT 149, 5–7 (Feb. 14, 2024) (Can.); see Lisa R. Lifshitz & Roland Hung, *BC Tribunal Confirms Companies Remain Liable for Information Provided by AI Chatbot*, AM. BAR ASS'N (Feb. 29, 2024), https://www.americanbar.org/groups/business_law/resources/business-law-today/2024-february/bc-tribunal-confirms-companies-remain-liable-information-provided-ai-chatbot/ [<https://perma.cc/NT5R-5NQ3>].

government as they continue to do against errant agents administering benefits.²²⁷ Nevertheless, chatbots offer guidance that occupies an uncertain legal space. Agency guidance escapes notice-and-comment and is more insulated from judicial review, but individuals often rely on informal agency guidance, even though it is legally nonbinding.²²⁸ And chatbot guidance is even more likely to induce reliance than average agency postings, Joshua Blank and Leigh Osofsky argue, because of the way in which it presents complex law as though it is simple (what they term “simplicity”) and its personalized, nonqualified, and instantaneous means of communication.²²⁹

Indiana’s generative AI chatbot’s errors differ from the nuanced issue of simplicity that Blank and Osofsky found when engaging with the IRS’s Interactive Tax Assistant. Indiana’s chatbot performed well on nuanced eligibility questions, but gave completely wrong answers on more fact-bound timing and calculation questions.²³⁰ While most of these errors were relatively benign and in response to rather detailed questions, claimants might interpret these responses as indicia that they are ineligible, much as the plaintiff in *Hansen* did.²³¹ Moreover, a chatbot user seeing this response might be more likely to believe it as a factual statement of law, compared to a user browsing the Frequently Asked Questions (FAQ) page, given how relevant and personalized the other responses by the tool are.²³²

The effectiveness of government chatbots for claimants and administrators remains uncertain. A UK Behavioral Insights Team experiment found mixed results: chatbots slightly delayed information discovery but made tasks feel easier and increased trust in AI.²³³ Some state administrators report high

227. *Hansen*’s holding continues to be regularly cited in appeals of Social Security benefit decisions. See, e.g., *Kimm R. W. v. O’Malley*, No. 23-cv-2497, 2024 WL 3896613, at *3–4 (D. Minn. July 15, 2024); *Holdaway v. Kijakazi*, No. 1:22-CV-219, 2023 WL 8007105, at *1, *13 (E.D. Tenn. Nov. 17, 2023). However, one might argue that more of the *Fredericks* factors are met in the context of certain benefits chatbots, where the impact on the public fisc is small and where the court found that the chatbot’s programming was sufficiently clear as to constitute official authorization of the misrepresentation.

228. *Holdaway*, 2023 WL 8007105, at *1, *13; see also Ronald M. Levin, *Rulemaking and the Guidance Exemption*, 70 ADMIN. L. REV. 263 (2018); Cass R. Sunstein, “Practically Binding”: *General Policy Statements and Notice-and-Comment Rulemaking*, 68 ADMIN. L. REV. 491 (2016).

229. See Joshua D. Blank & Leigh Osofsky, *Automated Legal Guidance*, 106 CORNELL L. REV. 179 (2020). Blank and Osofsky’s statement of causality is largely theoretical, rather than empirical.

230. See *infra* App. Table B.

231. See *Schweiker v. Hansen*, 450 U.S. 785 (1981) (per curiam).

232. See JOSHUA D. BLANK & LEIGH OSOFSKY, *AUTOMATED LEGAL GUIDANCE AT FEDERAL AGENCIES* 15, 20 (2022).

233. *AI Chatbots in Public Services*, THE BEHAV. INSIGHTS TEAM (Dec. 12, 2023), <https://>

usage, improved customer experience, and have won awards for chatbots.²³⁴ However, few rigorous evaluations exist on their impact in the United States, such as speeding up queries, improving take-up rates, reducing errors, or saving staff time for critical tasks.²³⁵

There may still be some cause for optimism about the impact of government chatbots. Human call center staff—the baseline of comparison—have constantly struggled to balance error rates and wait times. For example, in 1987, GAO found that the IRS help line provided wrong or incomplete answers 38% of the time.²³⁶ By 2008, the IRS’s error rate had fallen to 9%, but with worsening wait time and answer rates.²³⁷ While similar statistics are not available for unemployment call centers, it is widely known that such centers are typically staffed by the newest, least experienced staff as a means of on-

www.bi.team/wp-content/uploads/2024/01/External-AI-Chatbot-Trial-For-presentations.pdf [<https://perma.cc/5G4D-M73M>].

234. See Quinlan, *supra* note 193; MAXIMUS, *supra* note 186 (claiming that a host of tools, including chatbots, significantly reduced call volume, though it is unclear whether these findings are due to users giving up or due to particularly volatile time comparisons during the pandemic). These claims are similar to those made by federal administrators of chatbots. See also Joshua D. Blank & Leigh Osofsky, *supra* note 232.

235. But c.f. Tzu hao Chen & Mila Gasco-Hernandez, *Uncovering the Results of AI Chatbot Use in the Public Sector: Evidence from US State Governments*, 48 PUB. PERFORMANCE & MGMT. REV. 1331, 1335 (2024). There does, however, seem to be a positive impact of online claim status tools on call center volumes. For example, North Carolina found that the implementation of a claim status tool was able to deflect a number of calls that would typically go to a call center agent. Justin Lai & Christina Steen, *Reducing Administrative Waste in the Unemployment Insurance System*, U.S. DEP’T OF LAB. (Aug. 17, 2023), <https://www.dol.gov/agencies/eta/ui-modernization/promising-practices/admin-waste> [<https://perma.cc/C3NL-55CS>]. However, part of North Carolina’s modernization effort included an Interactive Voice Response (IVR) system which routes callers based on their needs; it is possible that the call deflection was caused by users finding it more difficult to use IVR and giving up, rather than successfully getting what they needed from the IVR system. A similar introduction of a claim status tool in New Jersey also found reduction in call center wait times. *Unemployment Insurance Claim Status Tool*, N.J. OFF. OF INNOVATION, <https://innovation.nj.gov/projects/ui-claim-tool/> [<https://perma.cc/T5JE-VFL3>] (Aug. 27, 2025).

236. U.S. GOV’T ACCOUNTABILITY OFF., GAO/GGD-88-17, TAX ADMINISTRATION: ACCESSIBILITY, TIMELINESS, AND ACCURACY OF IRS’ TELEPHONE ASSISTANCE PROGRAM (1987).

237. See INTERNAL REVENUE SERV., IRS TOLL-FREE TELEPHONE SERVICE IS DECLINING AS TAXPAYER DEMAND FOR TELEPHONE SERVICE IS INCREASING (2009), https://www.irs.gov/pub/tas/msp_1.pdf [<https://perma.cc/Z2JB-YZXM>]. The IRS had actually set a goal such that “nearly three out of every ten calls seeking to reach an IRS telephone assistor will not get through, and callers who do receive assistance will first have to wait on hold for an average of nearly 12 minutes.” *Id.* at 4.

the-job training, raising doubts about the accuracy of their guidance.²³⁸ Moreover, call centers introduce their own access problems: unemployment call centers are generally only open from 8 AM to 4 PM on weekdays, with some closing as early as 2 PM.²³⁹

Apart from Indiana, current state chatbots function as glorified FAQ pages. As such, the USDA classifies most chatbots as low-risk AI, but warns that they may “lead[] eligible populations to believe they are ineligible for public benefits, provid[e] incorrect program information, or creat[e] barriers to accessing public benefits.”²⁴⁰ Control over information disseminated by chatbots decreases when answers are generatively created, a distinction that USDA notes in categorizing “[c]hatbots using natural language processing to better understand user questions, with human-coded, logic-based preset outputs” as an enabling use of AI, but exempting “generative AI responses” from this use.²⁴¹

Most chatbots do not raise merit staffing issues, as providing generic FAQ information is not inherently governmental, similar to USDA’s allowance for non-merit staff in call centers.²⁴² Yet Indiana’s chatbot gives personalized responses on eligibility, which may conflict with merit staffing rules. For example, in response to the prompt, “I have earned \$7,000 in the base period with \$3,000 in one quarter. Am I eligible for unemployment insurance?” the chatbot (erroneously) explains the base period requirements and concludes, “Since you earned \$7,000 in total and \$3,000 in one quarter, you meet the first two requirements.”²⁴³ DOL specifically identifies “[a]dvising a claimant regarding his or her eligibility for benefits based on his or her specific circumstances” as an activity that must be conducted by merit staff.²⁴⁴ Indiana, hence, seems to have crossed a line. On the other hand, without rigorous evaluation, it is unclear how much more helpful merit staff who answer claimant calls—often newly-hired, inexperienced staff²⁴⁵—would be compared to a generative chatbot.

238. See 2022 GAO REPORT, *supra* note 10, at 15.

239. See Lester Coffey, Richard Sullivan, Paul Doherty, Mason Erwin, Lloyd Pete Fleming, Everett Hensley, et. al., UI Call Center Study Final Report 18 (Jan. 17, 2017) (report to the U.S. Dep’t of Lab.).

240. U.S. DEP’T OF AGRIC., FRAMEWORK FOR STATE, LOCAL, TRIBAL, AND TERRITORIAL USE OF ARTIFICIAL INTELLIGENCE FOR PUBLIC BENEFIT ADMINISTRATION 11 (2024), <https://digitalgovernmenthub.org/wp-content/uploads/2024/05/2024-usda-ai-framework-benefits.pdf> [<https://perma.cc/UK3G-HF37>].

241. *Id.* at 13.

242. *Id.* at 3, 9. This is less clear, of course, if a user interprets the chatbot’s output as personalized eligibility advice.

243. See *infra* App. Table B.

244. UIPL No. 12-01 CHANGE 2, *supra* note 137, at 5.

245. See 2022 GAO Report, *supra* note 10, at 15.

3. *Wave Three: Adjudication Assistance*

The most recent frontier is adjudication assistance: leveraging AI-based tools to help adjudicators make faster and better decisions such as highlighting relevant information or making eligibility predictions. To some, adjudication assistance is a step toward a fully automated benefits eligibility system, which has explicitly been a goal for Medicaid for years and achieved by systems in some states,²⁴⁶ such as the automatic eligibility determinations provided to Medicaid applicants in California.²⁴⁷

State experiments with adjudication assistance for UI and SNAP benefits cover a wide spectrum of both reliance on AI as well as level of automation. In many ways, AI-based adjudication assistance is an outgrowth of RPA, which automates nondiscretionary aspects of claim processing through the use of rule-based bots.²⁴⁸ For example, states use RPA for tasks like pulling data from one digitized textual source and entering it into another and assembling packets of information/documents, particularly for appeals.²⁴⁹ In the context of adjudication, this might look like a bot “review[ing] employer fact-finding requests and conduct[ing] data entry tasks to enter the separation reason into the benefits system.”²⁵⁰

SNAP agencies have been implementing RPA to assist adjudicators, though with limited success. For example, Connecticut worked with Deloitte to implement RPA in 2021 to assist with benefits renewals. The RPA bot effectively “stage[s] the case” for the eligibility worker by reviewing the application to determine if a case is a “change” or “no-change” renewal.²⁵¹ If

246. See Sarah Grusin, *A Promise Unfulfilled: Automated Medicaid Eligibility Decisions*, NAT’L HEALTH L. PROGRAM (June 30, 2021), <https://healthlaw.org/a-promise-unfulfilled-automated-medicare-eligibility-decisions/> [https://perma.cc/LA7J-J8BT]. As of May 1, 2024, forty-eight states were able to make real-time determinations at application, with those determinations in twenty-eight states being mostly automated. See also *System Automation of Real-Time Determinations and Ex Parte Renewals*, KFF, <https://www.kff.org/affordable-care-act/state-indicator/real-time-eligibility-determinations-application-exparte/> [https://perma.cc/36Z3-EV76] (Jan. 2025).

247. *Medi-Cal Eligibility & Covered California – FAQ’s*, CAL. DEP’T OF HEALTH CARE SERVS., <https://www.dhcs.ca.gov/services/medi-cal/eligibility/Pages/Medi-CalFAQs2014.aspx> [https://perma.cc/BYZ3-DTA9].

248. See Nikki Zeichnerand & Amy Perez, *Use Cases for Robotic Process Automation in UI Claims Processing*, U.S. DEP’T OF LAB. (July 11, 2023), <https://www.dol.gov/agencies/eta/ui-modernization/promising-practices/rpa-use-cases> [https://perma.cc/JR3W-58PT].

249. See *id.*

250. See *id.*

251. See U.S. DEP’T OF AGRIC., ANALYSIS OF ROBOTIC PROCESS AUTOMATION IN SUPPLEMENTAL NUTRITION ASSISTANCE PROGRAM: THREE CASE STUDIES 4, 12 n.8 (2023)

the bot deems the case does not need review, eligibility workers do not need to verify its work.²⁵² The RPA tool ended up increasing the average days to decision for cases despite the fact that the RPA generally mostly saved time on the simplest cases in a worker's caseload; similar effects were seen after implementation of a recertification bot in Georgia.²⁵³ Officials observed that insufficient training likely contributed to workers' mistrust of the RPA, leading to redundant efforts that undermined the project's success.²⁵⁴

State SNAP agencies have recently moved toward more AI-based automation to assist adjudicators.²⁵⁵ For example, several states use AI-based RPA to automatically process no-change periodic reports from claimants. Georgia, for instance, uses automated tools to "prepar[e]" SNAP recertification applications for employees by identifying mismatches between what is on the application and already in the eligibility system and creating "red flags" for merit workers to review.²⁵⁶ However, the exact details of the AI technology used in these automation workflows is unclear.

Adjudication assistance—RPA and AI-based alike—has been pivotal for state unemployment agencies clearing their unemployment backlogs during and after the pandemic.²⁵⁷ In Wisconsin, despite hiring over 1,300 staff,²⁵⁸

[hereinafter THREE CASE STUDIES REPORT], <https://fns-prod.azureedge.us/sites/default/files/resource-files/snap-bots-rpa-final-report.pdf> [https://perma.cc/H7ZN-WMMW] ("Change renewal cases include revisions to participant information (e.g., employment status, address) from the previous benefit authorization. No-change renewal cases include the same participant information as the prior benefit authorization.").

252. *See id.* Cases that do not require verification may still require eligibility workers to follow up with a client or otherwise resolve discrepancies, though the worker would do so following steps listed out by the robotic process automation (RPA).

253. *See id.* at 18–19. Days to decision is an imperfect measure of time savings as it may be influenced by other outside factors, such as waiting for claimant clarification of information, and is not a direct measure of worker productivity. The lack of effect in Georgia may be offset by reduced error rates (less than a quarter that of the full sample), though this impact is murky due to the fact that the RPA was generally implemented on simpler cases.

254. *See id.* at 33–34.

255. *See* Memorandum from USDA Food & Nutrition Serv. to All SNAP State Agencies & All Regions (Feb. 24, 2025), <https://www.fns.usda.gov/snap/advanced-automation> [https://perma.cc/2MEW-3F5U].

256. THREE CASE STUDIES REPORT, *supra* note 251, at 22.

257. *Id.* at 14. The state failed to provide adequate training for workers using the RPA tool, and so some eligibility workers thought the RPA's comments asking for individual review were the result of a "fellow staff member who was shirking their work and not processing cases correctly." *Id.*

258. Brent Mitchell, *Economic Recovery: Wisconsin Leans on Google Cloud to Better Serve Its Community During COVID-19 and Beyond*, GOOGLE CLOUD BLOG (Aug. 12, 2021), <https://cloud.google.com/blog/topics/industry-trends/economic-recovery-wisconsin-leans-on-google-cloud-to-better-serve-its-community-during-covid-19-and-beyond>.

a 770,000-claim backlog persisted in 2021.²⁵⁹ The state worked with Google to train a machine-learning model on historical data to assign a “confidence score” to applications, “indicating how likely it was that a given application should be approved or forwarded to human agents for possible rejection.”²⁶⁰ The tool only processed “hold[s]” where eligibility could not be immediately verified, either removing the hold (granting benefits) or leaving it for manual review.²⁶¹ According to the unemployment division secretary, “AI taught itself all of the eligibility rules in the state of Wisconsin by ingesting four million previously fully adjudicated claims that had been adjudicated by humans, and it made no mistakes.”²⁶² It remains unclear what “no mistakes” precisely indicates. Regardless, the model—affectionately dubbed Judy the Super Adjudicator—eliminated the backlog in a month and reduced response times to two to three business days.²⁶³ Officials emphasized that no benefits were denied based on the tool.²⁶⁴

An audit by Wisconsin’s Legislative Audit Bureau, however, revealed issues.²⁶⁵ First, it was unclear whether the tool met merit staffing requirements, prompting a recommendation to seek written assurance from DOL.²⁶⁶ Second, the agency lacked proper procedures to monitor the model’s accuracy,

google.com/blog/topics/public-sector/economic-recovery-wisconsin-leans-google-cloud-better-serve-its-community-during-covid-19-and-beyond [https://perma.cc/87KP-P6TG].

259. Colin Wood, *Google’s AI Helped Wisconsin Clear Unemployment Backlog*, STATESCOOP (Aug. 12, 2021), <https://statescoop.com/googles-ai-helped-wisconsin-clear-unemployment-backlog/> [https://perma.cc/K6B3-7LG9].

260. *Id.* Google’s analysis of where claimants were getting stuck in processing also led to the department rewiring the UI claim application process and helped the department identify fraudulent claims. Mitchell, *supra* note 258; Wood, *supra* note 259. It is unclear exactly how sophisticated the model was and is today, though Google describes it as “[a]rtificial intelligence (AI)/machine learning (ML) for predictive analytics.” Mitchell, *supra* note 258.

261. WIS. LEGIS. AUDIT BUREAU, NO. 21-23, STATE OF WISCONSIN FY 2020-21 FINANCIAL STATEMENTS 28 (2021), https://legis.wisconsin.gov/lab/media/v1znovll/21-23_332655_full.pdf [https://perma.cc/5YRL-GFPL].

262. Nikki Davidson, *Wisconsin’s AI Strategy: Upskilling, Not Displacing Workers*, GOV’T TECH. (Dec. 5, 2023), <https://www.govtech.com/artificial-intelligence/wisconsins-ai-strategy-upskilling-not-displacing-workers> [https://perma.cc/HP6P-F5D4].

263. Natalie Yahr, *Wisconsin Unemployment System Will Get Yearly Checkups After Overhaul*, THE CAP TIMES (Dec. 20, 2023), https://captimes.com/news/business/wisconsin-unemployment-system-will-get-yearly-checkups-after-overhaul/article_d7fe9278-9ef8-11ee-a224-0bb004cbe310.html [https://perma.cc/XWY5-N7HF].

264. Press Release, Wis. Dep’t of Workforce Dev., DWD Clears Wisconsin Unemployment Claims Backlog (Dec. 30, 2020), <https://dwd.wisconsin.gov/press/2020/201230-unemployment-claim-backlog-cleared.htm> [https://perma.cc/2LKK-3WG2].

265. WIS. LEGIS. AUDIT BUREAU, *supra* note 261, at 30.

266. *Id.* at 1.

relying only on initial testing.²⁶⁷ Weekly projected error rates during the tool's six months of use ranged from 5.5% to 26.8%; however, the tool's role in processing holds and how its error rate compares to the federal DOL BTQ measure (29%–33%)—which includes procedural imperfections, such as inadequate explanation—remain unclear.²⁶⁸

Nevada has also heavily utilized AI to address pandemic-related unemployment claim backlogs, leveraging generative AI to gather data on a claim and then “make a recommendation that a human staff member can take into consideration when assessing a case.”²⁶⁹ In partnership with Google, it plans to expand AI use to the appeals process by developing a tool that analyzes appeals hearing transcripts and evidentiary documents to recommend benefit decisions.²⁷⁰ The tool draws on unemployment law and prior cases, reportedly reducing determination times “from several hours to just five minutes,” helping address the appeals backlog.²⁷¹ Officials stress that all decisions undergo human review, with referees reviewing recommendations they disagree with and investigating discrepancies.²⁷²

Automation bias—the tendency to overrely on automated recommendations—raises concerns about the role of human intermediaries. A former DOL deputy director expressed that it was “a little bit concerning” if “a robot’s just handed you a recommendation and you just have to check a box

267. The Wisconsin Department of Workforce Development responded that during the design phase, the “accuracy of each model’s prediction was compared to the actual, manually determined resolution of that hold, and this information was used to calculate an error rate by hold type.” *Id.* at 30. The agency rejected the mandate to conduct ongoing testing, seeming to interpret it as a mandate to retrain the model, something they deemed only necessary when there are “significant changes to the environment” such as law or policy changes. *Id.* at 31.

268. Figures in the report suggest that the tool may not have been responsible for a significant proportion of cases. From December 2020 to June 2021, the tool removed 169,257 holds. Throughout 2020 and 2021, three million total holds were removed by both the tool and the model. A middle of the road estimate that assigns a quarter of the three million holds to the tool’s active six-month period would suggest the tool is responsible for only 22.5% of the holds removed (and unclear percent of all holds reviewed). *Benefits: Timeliness and Quality Reports*, *supra* note 12.

269. Lauren Kinkade & Noelle Knell, *Nevada Harnesses GenAI for Employment Claims Evaluation*, GOV’T TECH. (Apr. 29, 2024), <https://www.govtech.com/artificial-intelligence/nevada-harnesses-genai-for-employment-claims-evaluation> [<https://perma.cc/T9YY-F6T4>].

270. Todd Feathers, *Google’s AI Will Help Decide Whether Unemployed Workers Get Benefits*, GIZMODO (Sep. 10, 2024), <https://gizmodo.com/googles-ai-will-help-decide-whether-unemployed-workers-get-benefits-200049621> [<https://perma.cc/7G8J-U3HG>].

271. *Id.*

272. *Id.*

and there's pressure to clear out a backlog."²⁷³ State representatives claim to have run "dozens of tests using the company's technology to analyze hearing transcripts from appeals cases of varying complexity" before agreeing to a \$1 million contract with Google.²⁷⁴ A governance committee will monitor the system weekly during fine-tuning and quarterly after launch to address hallucinations and bias.²⁷⁵ Google has developed similar contracts with Michigan and California, as well as contracts for contact center AI-based virtual agents.²⁷⁶

Adjudication assistance tools often stand on uncertain ground with merit staffing requirements. Under DOL rules, many state AI tools analyzing facts to determine benefits eligibility—a task classified as inherently governmental—are filling merit staff roles.²⁷⁷ Some tools even come close to "[a]ctually making a determination of benefit eligibility"—another inherently governmental activity—at least in determinations of benefit approval.²⁷⁸ In SNAP, AI implementations required FNS approval as a "major change" under regulations,²⁷⁹ yet it remains unclear if they comply with the Food and Nutrition Act's requirement that "personnel used in the certification process" be merit staff.²⁸⁰ A loophole might interpret "personnel" to exclude AI, which is not staff.

* * *

These examples illustrate a significant trend: states are in many ways at the frontier of experimenting with AI. They are doing so despite considerable regulatory and legal uncertainty, particularly considering requirements

273. *Id.*

274. *Id.*

275. *Id.* The minutes from the June 2024 meeting of the state's IT advisory board are illuminating: for example, the head of the board asked other advisory board members (who are predominately from private sector technology companies) if there is "a standard practice the private sector uses when acquiring a new software." Meeting Minutes, Nev. Info. Tech. Advisory Bd. 8 (June 20, 2024), <https://it.nv.gov/uploadedFiles/itnewnv.gov/content/Governance/ITAB/Meetings/ITAB%20Meeting%20Minutes%2006202024.pdf?csrt=9648965750860574072> [<https://perma.cc/8NGY-4L59>]. The state's Chief Information Officer advocated for "continuous evaluation improvement," and acknowledged that if a chatbot gives a constituent "an incorrect but plausible answer, the constituent may take that as the truth which would create a difficult situation for the agency." *Id.* at 10.

276. *Agency Leaders Discover the Power of AI to Scale and Support Citizen Services*, STATESCOOP (Apr. 17, 2024), <https://statescoop.com/agency-leaders-discover-the-power-of-ai-to-scale-and-support-citizen-services/> [<https://perma.cc/Y8E4-DDUT>].

277. UIPL 12-01 CHANGE 2, *supra* note 137.

278. *Id.* at 5.

279. Advanced Automation in SNAP Memo, *supra* note 157; 7 C.F.R. § 272.15 (2017).

280. 7 C.F.R. § 272.4 (2017).

that effectively mandate that any material analysis of claimant facts for eligibility be conducted by merit staff. This legal uncertainty, combined with the need to meet federal quality and timeliness mandates on thin and inflexible budgets, may have further pushed states to experiment independently of federal supervision. While initiatives like Tiger Teams and USDA's contract for state RPA case studies have facilitated knowledge sharing, these efforts remain limited and vulnerable to funding cuts due to their peripheral roles within federal agencies.

C. The Federal Precursor to Modernization Through Evaluation

State experimentation with AI, amid legal uncertainty and limited federal oversight, can be seen as one end of a spectrum of approaches to government modernization. It has allowed for diverse and fast experimentation, but scant independent evaluation has made it difficult to build credibility and ensure due process. At the other end of the spectrum are the reform efforts of the IRS, perhaps the most closely watched and criticized agency by Congress and the press.²⁸¹ The IRS's implementation of phone-based call centers provides a helpful illustration of an approach to modernization that places evaluation at its core.

In 1974, the IRS changed the way that millions of Americans interacted with the agency by introducing a telephone tax assistance system, eliminating the need to resolve questions in-person or over the mail.²⁸² The system quickly scaled to employ thousands of "assistors" across thirty-two call centers, handling nearly twenty million taxpayer calls annually.²⁸³ Doubt about the efficacy of the program, and its cost, soon followed. In 1978, the House Government Operations Committee asked GAO to evaluate the effectiveness of the new telephone assistance program; specifically, "(1) the extent to which taxpayers' telephone calls are answered by IRS, (2) how long it takes for calls to be answered, and (3) the extent to which taxpayers are receiving

281. See generally CHARLES O. ROSSOTTI, MANY UNHAPPY RETURNS 7–23 (2005).

282. *Internal Revenue Service's Taxpayer Assistance Programs, Monitoring of Commercial Tax Return Preparers and Tax Form Simplification Efforts, Hearing Before the H. Subcomm. on Gov't Operations*, 95th Cong. 2–30 (1978) (statement of Victor Lowe, Director, General Government Division, General Accounting Office); U.S. GOV'T ACCOUNTABILITY OFF., GGD-75-69, TELEPHONE ASSISTANCE TO TAXPAYERS CAN BE IMPROVED 5 (1975). Telephone assistance not only increased convenience for taxpayers, eliminating the need to visit an in-person IRS office, but it also initially allowed the IRS to answer an estimated 33% more people than with in-person service. *Id.* at 25.

283. U.S. GOV'T ACCOUNTABILITY OFF., TAX SYSTEM MODERNIZATION: FURTHER TESTING OF IRS' AUTOMATED TAXPAYER SERVICE SYSTEMS IS NEEDED 2 (1991) [hereinafter 1991 GAO REPORT].

correct answers to their inquiries.”²⁸⁴ The initial audit, based on a small sample size, provided promising results: about 96% of calls were answered; on average it took one minute to get in contact with an employee; and about 87% of the IRS’s responses were correct.²⁸⁵

Over the next decade, GAO conducted semiregular audits of the accessibility and accuracy of the program, gradually increasing and improving the quality of tax law questions asked as well as its sampling methodology.²⁸⁶ While the IRS generally agreed that GAO’s questions were reasonable, it often had critiques of the survey design and methodology.²⁸⁷ These disagreements spurred the IRS to develop its own internal evaluation system, the Integrated Test Call Survey System (ITCSS), in 1988.²⁸⁸ The ITCSS aimed to better reflect taxpayers’ real-world inquiries and improve benchmarking accuracy. By 1989, GAO had validated ITCSS reliability, aligning its findings with the IRS’s own assessments.²⁸⁹ This system remains in place today, with improved accuracy rates—by 2023, the IRS reported 90.4% accuracy in responses to tax law inquiries.²⁹⁰ The IRS stands distinct from DOL and USDA in this regard; while state agencies implementing both programs operate UI and SNAP telephone assistance systems that provide critical clarity to applicants, there is currently still no comprehensive system evaluating the accessibility and accuracy of these systems—it took independent researchers making 2,000 calls to government offices in 2021 to determine that less than 40% of calls to UI agencies reached a live representative.²⁹¹

284. Lowe, *supra* note 282, at 1–2.

285. However, taxpayers referred to a specialist—i.e., those with a more difficult question—only received a correct answer 79% of the time. *Id.* at 1.

286. U.S. GOV’T ACCOUNTABILITY OFF., GAO-89-30, TAX ADMINISTRATION: ACCESSIBILITY, TIMELINESS, AND ACCURACY OF IRS’ TELEPHONE ASSISTANCE PROGRAM (1989).

287. *Id.* at 25, 34.

288. *Id.* at 25, 34–35. GAO repeatedly alleged that the “IRS has agreed that our questions were reasonable”; however, it also claimed that “no one knows what constitutes a representative set of [evaluation] questions nor are we certain that a representative set of questions can be developed.” *Id.* at 25.

289. U.S. GOV’T ACCOUNTABILITY OFF., GAO-90-36, TAX ADMINISTRATION: MONITORING THE ACCURACY AND ADMINISTRATION OF IRS’ 1989 TEST CALL SURVEY 6 (1990), <https://www.gao.gov/assets/ggd-90-36.pdf> [<https://perma.cc/FR5M-PXVU>].

290. INTERNAL REVENUE SERV., INTERNAL REVENUE SERVICE DATA BOOK 2024 (2024) [hereinafter IRS DATA BOOK], <https://www.irs.gov/pub/irs-pdf/p55b.pdf> [<https://perma.cc/P3VU-KGZR>].

291. Oeindrila Dube, Sendhil Mullainathan & Devin G. Pope, *A Note on the Level of Customer Support by State Governments: A Mystery-Shopping Approach* (Becker Friedman Inst. Working

Crucially, the IRS's evaluation system highlighted the interplay between human and machine systems in service delivery. Early evaluations identified assistor error rates as a significant challenge, prompting the IRS to explore automation as a means of improving accuracy.²⁹² The IRS piloted two automation projects in 1989: (1) an "automated taxpayer service system" which replaced assistors' paper reference materials with a computer database that helped automate researching taxpayer questions, and (2) an "expert system" which dynamically guides assistors through questions to ask taxpayers.²⁹³ While preliminary results suggested a 21% improvement in accuracy with the expert system, subsequent GAO audits flagged flaws in the pilot's design, namely that the IRS "did not consider whether assistors used the Expert System properly, or at all."²⁹⁴ As a result, OMB denied the IRS's budget request to continue funding the automation experiment, insisting on further testing first.²⁹⁵ These critiques underscore the complexities of designing, benchmarking, and evaluating systems that integrate human and automated components—issues that remain central to contemporary debates about AI governance. Moreover, the IRS's acknowledgement of the shortcomings of the status quo and embrace of automation stand out as an early example of technology being deployed to enhance accuracy rather than merely efficiency.

The IRS's telephone assistance program exemplifies how rigorous evaluation can serve as a foundation for government modernization. It demonstrates that the government can—and does—conduct robust evaluations to improve service delivery. It was spurred by close congressional scrutiny, which is often triggered by taxpayer complaints, but its continued success also depended on the interplay of statutory support, institutional will, and capable personnel to drive the initiative forward.²⁹⁶

Paper No. 2021-89, 2021), https://bfi.uchicago.edu/wp-content/uploads/2021/07/BFI_WP_2021-89.pdf [<https://perma.cc/XNX3-GKET>]. The lack of call center evaluation for UI seems particularly egregious when considering that the UI call volume is at least comparable to, if not greater than, IRS call volume. For example, California's UI agency received millions of calls *per week*—as many as 9.7 million—in 2021; the entire IRS received 37 million calls in 2023 and answered 27.3 million of them. IRS DATA BOOK, *supra* note 290, at iv, 21, 24; *California Call Center Data Dashboard*, CAL. EMP. DEV. DEP'T (May 2022), <https://edd.ca.gov/site-assets/files/pdf/phone-calls-dashboard-052622.xlsx> [<https://perma.cc/Y37V-D7SX>].

292. 1991 GAO REPORT, *supra* note 283, at 2, 12, 14.

293. *Id.* at 2–3.

294. *Id.* at 4.

295. *Id.* at 8, 10.

296. The 1998 Restructuring Act allowed these new measures to take a more central role in determining employee performance. Key leaders eager to improve customer service, such as Commissioner Rossotti, were likely also critical in the ongoing development of the program. ROSSOTTI, *supra* note 281.

The lessons from this case extend to modern challenges in evaluating AI and human-machine interaction. Benchmarking must be an ongoing, iterative process, not a one-time exercise. Such ongoing evaluations are central to programmatic improvements. Evaluations should focus not only on outcomes—such as successful audits or total tax receipts—but also on understanding the dynamics between human users and automated systems. In other words, one cannot understand impact looking at the technical system alone. Finally, systems like the IRS’s telephone assistance program show that government agencies are and have long been capable of conducting rigorous, credible evaluations, provided they invest in the necessary structures and expertise. Such systems are indispensable as we consider how to evaluate emerging technologies like AI-driven decision tools, ensuring they align with values of accuracy, fairness, and public trust.

IV. THE LEGAL AND POLICY NECESSITY OF EVALUATION

State experimentation with AI is likely to only grow. The question for policymakers is how to manage the risk of these new technologies—such as benefits falsely flagged as fraudulent or incorrect advice about eligibility from a chatbot—while still reaping the benefits of sorely needed innovation in under-resourced agencies. This section examines the legal and policy imperatives for rigorous evaluation of AI systems in benefits administration. We argue that existing legal frameworks, including the Administrative Procedure Act, due process requirements, and procurement regulations, create an implicit mandate for agencies to evaluate AI tools before deployment. We then outline principles for meaningful evaluation and provide an example framework for assessing AI chatbots. Finally, we explore potential avenues for formalizing and implementing a more explicit evaluation mandate.

A. Evaluation of AI Systems Is Legally Mandated

1. Adopting AI Without Evaluation Is Likely a Violation of Procedural Due Process

The Due Process Clause demands that government agencies evaluate AI systems that meaningfully affect receipt and termination of benefits before deployment to avoid unconstitutional deprivations of benefits. Terminating public benefits without proper procedural safeguards violates constitutionally protected property interests.²⁹⁷ When agencies adopt AI systems without rigorous evaluation, they risk systematic violations of these due process protections.

The *Mathews v. Eldridge* framework requires courts to balance three factors when determining adequate procedural safeguards: (1) the private interest

297. *Goldberg v. Kelly*, 397 U.S. 254, 264–65 (1970).

affected; (2) the risk of erroneous deprivation and value of additional safeguards; and (3) the government's interest, including administrative costs.²⁹⁸ The second factor—the risk of erroneous deprivation—is particularly critical for AI adoption. Without evaluation, agencies cannot quantify error rates, false positives in fraud detection algorithms, or the accuracy of chatbot-provided eligibility information. This lack of baseline measurement makes meaningful *Mathews* balancing impossible and endangers due process by imposing unknown and unmeasured risks of wrongful benefit denials.

Evaluation serves due process in broader ways than merely providing evidence-based approaches to *Mathews* balancing. The evaluation process itself should help prevent the adoption of systems that increase error rates and thus the risk of erroneous deprivation. Jerry Mashaw's "managerial" conception of due process reinforces this evaluation imperative by emphasizing systemic quality assurance over individual procedural remedies in court.²⁹⁹ The broader analytic scope of evaluation appears more likely to achieve due process's mandate of systemic accuracy than reliance on individual-level exercise of procedural rights through the appeals process, which likely does little to promote systemic accuracy.³⁰⁰

Due process has already figured prominently in the most significant cases challenging agency's attempts or failures to modernize benefit systems, but these cases have stopped short of establishing system-level mandates. In the sprawling federal case against the Michigan unemployment agency and contractors for the failure of the state's automated fraud detection system, residents alleged a due process violation for termination of payments without adequate notice and hearing.³⁰¹ While the due process claims survived a motion to dismiss, plaintiffs' class certification failed due to their varied positions of deprivation, leaving just four plaintiffs to continue with the case. The Sixth Circuit ultimately held that the unemployment agency supervisor defendants were entitled to qualified immunity because they did not trespass

298. *Mathews v. Eldridge*, 424 U.S. 319, 335 (1976).

299. Jerry L. Mashaw, *The Management Side of Due Process: Some Theoretical and Litigation Notes on the Assurance of Accuracy, Fairness, and Timeliness in the Adjudication of Social Welfare Claims*, 59 CORNELL L. REV. 772 (1974). Danielle Citron also develops a notion of inquisitorial quality assurance in her work. Citron, *supra* note 20.

300. David Ames, Cassandra Handan-Nader, Daniel E. Ho & David Marcus, *Due Process and Mass Adjudication: Crisis and Reform*, 72 STAN. L. REV. 1, 23 (2020); Mashaw, *supra* note 299, at 785.

301. *Cahoo v. SAS Analytics, Inc.*, 912 F.3d 887 (6th Cir. 2019). Plaintiffs also brought state law claims relating to product reliability and gross negligence. These counts were dismissed on account of inadequate tortious damage to the plaintiffs.

on clearly established law.³⁰² The state case, however, faced fewer class certification hurdles.³⁰³ The Michigan Court of Appeals ultimately found that the agency's use of untested automated decisionmaking violated the state constitution's due process clause.³⁰⁴

Similarly, plaintiffs in Missouri brought a case against the state SNAP agency, alleging that the agency's wrongful denials of their SNAP benefits for failure to interview violated due process.³⁰⁵ The severely understaffed system effectively made it impossible to schedule an interview. The district court ordered injunctive relief in the form of monitoring the state SNAP agency's performance improvements, particularly around its ability to offer interviews and answer calls in a timely fashion.³⁰⁶ In some ways, the Missouri case is the flipside of the Michigan case. Both agencies were overwhelmed and understaffed. Michigan's unemployment agency turned to a form of auto-adjudication and found itself in violation of due process for its arbitrary assignments of fraud and inadequate notice. The Missouri agency did not try to modernize and instead stuck to its status quo of understaffed services; yet it too found itself in the same due process hole for arbitrary denials.

302. This is partly due to the fact that the court was restricted to review the more subtle due process violation that was applicable to the four remaining plaintiffs. Those plaintiffs—unlike *Goldberg* and some of those in the original class—did *not* lose ongoing benefits. Rather, they received benefits and were later “auto-adjudicate[d] . . . guilty of fraud, disqualifying them for future benefits” and clawing back the previous benefits with penalty. *Cahoo v. SAS Inst., Inc.*, 71 F.4th 401, 420 (6th Cir. 2023). They received notice of this liability. Ultimately, the court held: “The content of the two notices, the opportunity for a hearing, and the months-or years-long delay before these plaintiffs faced a deprivation distinguish this case from existing due process precedent.” *Id.* at 408.

303. The lead plaintiff had received unemployment benefits but was later determined to have been overpaid due to suspected fraud. As a penalty, his tax refund was intercepted and his wages were garnished. *Bauserman v. Unemployment Ins. Agency*, No. 333181, 2017 WL 3044120 (Mich. Ct. App. July 18, 2017).

304. *Bauserman v. Unemployment Ins. Agency*, 983 N.W.2d 855 (Mich. 2022). Status reports warned that the MiDAS system could not read large chunks of converted legacy data. Paul Egan, *Michigan Integrated Data Automated System Experiences 93 Percent Error Rate During Nearly Two Years of Operation*, GOV'T TECH. (July 31, 2017), <https://www.govtech.com/data/michigan-integrated-data-automated-system-experiences-93-percent-error-rate-during-nearly-two-years-of-operation.html> [<https://perma.cc/T2H8-TUFT>]. The state auditor's report revealed that the agency did not fully analyze and validate data to help identify payments requiring further review. MICH. AUDITOR GEN., MICHIGAN INTEGRATED DATA AUTOMATED SYSTEM (MiDAS) 20 (2016).

305. For example, one resource center at which applicants are expected to interview “has a line out the door and over a two-hour wait because it does not have sufficient staffing and/or hours.” *Holmes v. Knodell*, 733 F. Supp. 3d 775, 785 (W.D. Mo. 2024).

306. *Id.* at 800–01.

These cases underscore the idea that true procedural due process requires systemic evaluation before full-scale deployment rather than the post-hoc remedies made possible through litigation. Left to generalist courts that may know little about program administration and evaluation design, litigation is not guaranteed to focus on and improve the agency-level decisions that matter most for the average claimant.³⁰⁷ Moreover, *Mathews* balancing is an inherently case-level, fact-bound exercise. As Danielle Citron has observed, this case-level application misses the fact that AI-based tools are designed to operate at scale, and that a small change to a tool because of procedural scrutiny may impact the thousands or millions of cases to which the tool is applied.³⁰⁸ The Michigan cases, particularly the federal case, ultimately turned on the exact wording sent out in various notices, highlighting the judiciary's tendency to focus on the notice aspect of due process to the exclusion of analysis of the actual fraud detection system's design, when and why it was used, and how it was supervised.³⁰⁹ While improved notice and appeal processes are likely fine outcome incentives, due process as it was litigated in Michigan seemed to miss the more important questions: how could the agency have designed a better fraud detection system? What critical conversations with, oversight of, or incentives for contractors were missed?

Moreover, due process, as cast by *Goldberg* and *Mathews* and as litigated to this day, often places accuracy as the sole cornerstone of adjudicative quality. As Mashaw argued, the governing caselaw's narrow focus on accuracy is "unresponsive to the full range of concerns embodied in the due process

307. As Adrian Vermeule has argued, "The federal judicial system is not set up, not equipped, to engage in a sustained course of synoptic institutional engineering." Adrian Vermeule, *LAW'S ABNEGATION: FROM LAW'S EMPIRE TO THE ADMINISTRATIVE STATE* 115 (2016). Moreover, due to its rather sparse constitutional basis, the exact scope and requirements of due process are subject to inconsistent judicial interpretation. For these and other reasons, Jerry Mashaw's managerial turn in due process instead places agency's internal administrative law at the center of quality assurance, rather than externalist judicial interpretations of constitutional due process. See generally JERRY L. MASHAW, *BUREAUCRATIC JUSTICE: MANAGING SOCIAL SECURITY DISABILITY CLAIMS* (1983).

308. Citron, *supra* note 21, at 1249; David Freeman Engstrom & Daniel E. Ho, *Algorithmic Accountability in the Administrative State*, 37 *YALE J. REG.* 800, 827 (2020).

309. See *Cahoo v. SAS Inst., Inc.*, 71 F.4th 401, 412 (6th Cir. 2023) ("Yes, MiDAS's logic trees spawned internal, interim fraud findings that marked individuals for further review. But that internal step did not deprive a claimant of property. As noted, those property deprivations came months or years later, following notices of deprivation and a multi-level appeal process. The missing link, then, is a case that clearly established the inadequacy of the questionnaires, notices of determination, and appeal processes available in the run-up to the property deprivation.")

clause,” such as dignity, transparency, and social trust.³¹⁰ Decisional accuracy should be but one outcome of interest in agency evaluation plans. Ideally, enforcement of an evaluation mandate would take a broader scope to encourage the consideration of multiple outcomes including accuracy and customer satisfaction.

Taking due process seriously means taking evaluation seriously. For AI adoption specifically, this means agencies must test error rates, validate decisionmaking processes, and ensure adequate human oversight mechanisms exist. Without such evaluation, agencies cannot satisfy *Mathews* balancing, cannot ensure systemic accuracy, and cannot avoid the constitutional violations that plagued both Michigan’s automated system and Missouri’s overwhelmed manual processes.

2. *Adopting AI Without Evaluation Is Likely Arbitrary and Capricious*

The Administrative Procedure Act (APA) and state analogues provide another potential foundation for an evaluation mandate. The APA forbids agency action that is “arbitrary, capricious, an abuse of discretion, or otherwise not in accordance with law.”³¹¹ Agency actions are arbitrary and capricious if they “entirely fail[] to consider an important aspect of the problem” or “offer[] an explanation . . . counter to the evidence before the agency.”³¹² An agency must “examine the relevant data and articulate a satisfactory explanation for its action including a ‘rational connection between the facts found and the choice made.’”³¹³

Agencies deploying AI systems without evaluation appropriately scaled to risk fail to consider relevant factors such as accuracy rates across populations, disparate impact on protected classes, consequences of errors, and implementation costs versus benefits.³¹⁴ Without evaluation, agencies cannot demonstrate rational connections between AI deployment and policy objectives or show that benefits outweigh risks. However, courts are generally hesitant to affirmatively require agencies to conduct evaluations under the

310. Jerry L. Mashaw, *The Supreme Court’s Due Process Calculus for Administrative Adjudication in Mathews v. Eldridge: Three Factors in Search of a Theory of Value*, 44 U. CHI. L. REV. 28, 30 (1976).

311. 5 U.S.C. § 706(2)(A).

312. *Motor Vehicle Mfrs. Ass’n v. State Farm Mut. Auto. Ins. Co.*, 463 U.S. 29, 43 (1983).

313. *Id.* (quoting *Burlington Truck Lines, Inc. v. United States*, 371 U.S. 156, 168 (1962)).

314. It is unclear, of course, as to whether disparate impact claims can sound in the Administrative Procedure Act (APA). See Cristina Isabel Ceballos, David Freeman Engstrom & Daniel E. Ho, *Disparate Limbo: How Administrative Law Erased Antidiscrimination*, 131 YALE L.J. 370 (2021).

APA, and for good reason.³¹⁵ Notably, the D.C. Circuit struck down the Federal Trade Commission's controversial rule requiring two randomized clinical trials (RCTs) to substantiate disease claims in order to comply with the FTC Act's deceptive advertising provisions.³¹⁶ The court criticized the rigidity of the two-study requirement, which appeared to ignore the Commission's own acknowledgement that "the quality of studies will be more important than quantity."³¹⁷ However, the court found no fault with a singular RCT requirement, noting that "*some* RCT substantiation for disease claims directly advances, and is not more extensive than necessary to serve, the interest in preventing misleading commercial speech."³¹⁸ We certainly do not advocate that full RCTs are always—or even often—required under the

315. Imposing such a blanket restriction could arguably run afoul of *Vermont Yankee Nuclear Power Corp. v. Natural Resources Defense Council, Inc.*, 435 U.S. 519 (1978) (holding that courts cannot impose additional rulemaking procedures on administrative agencies beyond the APA). Courts are, however, comfortable finding agency action arbitrary and capricious when the agency ignores relevant studies evaluating a policy's impact. See, e.g., *WildEarth Guardians v. U.S. Bureau of Land Mgmt.*, 870 F.3d 1222, 1235 (10th Cir. 2017) (conclusions "unsupported by hard data" do not provide "information sufficient to permit a reasoned choice" and are therefore arbitrary and capricious); *Genuine Parts Co. v. EPA*, 890 F.3d 304 (D.C. Cir. 2018) (arbitrary and capricious to rely on portions of studies that supported its position while ignoring cross sections in those studies that did not support the position); *WildEarth Guardians v. Haaland*, 561 F. Supp. 3d 890 (C.D. Cal. 2021) (holding it was arbitrary to ignore study projections of habitat impact). But cf. *Nantucket Residents Against Turbines v. U.S. Bureau of Ocean Energy Mgmt.*, 100 F.4th 1 (1st Cir. 2024) (finding an agency's decision to ignore studies on whale risks when approving an offshore wind energy project in compliance with the APA).

Perhaps one of the most notable cases in this vein is *Business Roundtable v. SEC*, in which the D.C. Circuit struck down a new SEC rule on shareholders for failure to first evaluate the likely consequences of the rule. See 647 F.3d 1144 (D.C. Cir. 2011). Here, the evaluation mandate was partly driven by the APA, but in large part also driven by an SEC-specific statutory obligation to consider the effects of new rules on "efficiency, competition and capital formation"; however, a similar argument could be made for the statutes undergirding benefits systems that require agencies to make timely and accurate decisions. *Id.* at 1148.

316. *POM Wonderful, LLC v. FTC*, 777 F.3d 478, 503 (D.C. Cir. 2015).

317. *Id.*

318. *Id.* at 502. The court also acknowledges the appropriateness of "competent and reliable scientific evidence"—which may not necessarily be a randomized clinical trial (RCT)—to substantiate advertising claims. *Id.* at 489. Notably, the court's analysis was conducted under *Central Hudson* scrutiny for First Amendment violations, not the APA. However, the specific prong of analysis that the court focused on—that the government's restriction on speech "directly advance[s] the governmental interest" and "is not more extensive than is necessary to serve that interest"—is highly analogous to arbitrary and capricious review under *State Farm*. *Id.* at 501; see *Motor Vehicles Mfrs. Ass'n Inc v. State Farm Mut. Auto. Ins. Co.*, 483 U.S. 29 (1982).

APA's arbitrary and capricious standard.³¹⁹ As we detail in Part II.B, evaluation requirements should be context-specific and proportionate to the stakes involved, unlike the FTC's categorical mandate.

A broad APA-based evaluation mandate risks becoming the very kind of process-heavy impediment to innovation that has hampered government technology adoption for decades. The danger is real and historically documented: if every AI integration requires exhaustive pre-deployment analysis, agencies may simply avoid beneficial automation altogether, defaulting to dysfunctional status quo systems that themselves violate due process, as we saw in Missouri's SNAP case.³²⁰ Critics might reasonably ask why AI systems should trigger evaluation requirements when the human systems they replace often lack such systematic assessment. This asymmetry problem has real force: state UI call centers, for instance, operate without the kind of accuracy monitoring that the IRS conducts, despite providing equally critical guidance to claimants.³²¹ There is a legitimate worry that imposing evaluation requirements only on AI deployments puts a thumb on the scale in favor of inadequate status quo systems, essentially penalizing agencies for attempting to modernize.

The solution lies in recognizing that "evaluation" encompasses a spectrum of approaches rather than a monolithic requirement. At one end are simple benchmark tests—measuring whether an AI system meets basic performance metrics like accuracy scores on standardized datasets. In the middle are sandbox deployments and limited pilots that test human–AI interaction in controlled settings without full deployment. At the other end are comprehensive randomized controlled trials and ongoing quality assurance systems like those Jerry Mashaw advocated. The APA's arbitrary and capricious standard should require agencies to select evaluation methods proportionate to the risks and stakes involved, not demand the most rigorous possible assessment in every circumstance, following the calibrated balancing demanded in the due process context by *Mathews*.

To specifically address the asymmetry issue, any evaluation requirement read from the APA should emphasize ongoing assessment rather than front-loaded predeployment testing. The Biden administration's OMB memo exemplified the dangers of process-heavy predeployment requirements that

319. *But see* Michael Abramowicz, Ian Ayres & Yair Listokin, *Randomizing Law*, 159 U. PA. L. REV. 929 (2011) (advocating for broader use of RCTs in policymaking and arguing that RCTs are compatible with the APA despite their necessarily "arbitrary" nature).

320. *See* *Holmes v. Knodell*, 733 F. Supp. 3d 775, 795 (W.D. Mo. 2024).

321. *See* *OIG Oversight of the Unemployment Insurance Program*, OFF. OF INSPECTOR GEN. (June 5, 2025), <https://oig.dol.gov/doloigoversightwork.htm> [<https://perma.cc/JHF7-U6RR>].

may bear little relationship to real-world performance.³²² A chatbot that performs well in laboratory testing may fail when confronted with the creative and contextual questions that real claimants ask. Conversely, an adjudication assistance tool that seems concerning in theory may prove more accurate and more fair than overwhelmed human adjudicators in practice. Predeployment evaluation, while valuable, can only go so far in predicting how AI systems will perform when they encounter the full complexity of real-world benefit administration.

To prevent the APA evaluation requirement from becoming a source of administrative paralysis, several limiting principles should apply. First, agencies should have safe harbors for pilot programs and limited deployments that affect small numbers of claimants and include robust fallback procedures. The mere act of testing an AI system on a subset of cases should not trigger the full apparatus of predeployment evaluation, particularly when agencies can demonstrate adequate safeguards and monitoring. Second, agencies that can demonstrate they lack adequate human baselines—such as those with call centers that do not track accuracy rates—should have flexibility to establish baseline measurements concurrent with AI deployment rather than requiring pre-existing benchmarks that may not exist. Third, the APA standard should focus on the reasonableness of an agency’s evaluation approach rather than demanding specific methodologies. An agency that conducts basic accuracy testing and implements ongoing monitoring should satisfy the arbitrary and capricious standard, even if it forgoes more sophisticated evaluation methods that might be appropriate for higher-risk deployments. The goal is ensuring reasoned decisionmaking based on evidence, not mandating particular analytical techniques.

3. Adopting AI Without Evaluation Is Likely Contrary to Procurement Law and Policy

One last implicit legal imperative for evaluation lies in procurement requirements. As our case studies illustrate, most government AI systems are not developed in-house: they are purchased. As such, these purchases are subject to procurement law.

Several elements of the Federal Acquisition Regulation (FAR) militate in favor of prepurchase evaluation of AI systems wherever possible.³²³ FAR

322. Comment Letter on OMB–2023–0020 from Daniel E. Ho, Percy Liang, Timothy O’Reilly, Jennifer Pahlka, Todd Park, DJ Patil, et al., Proposed Memorandum for the Heads of Executive Departments and Agencies: “Advancing Governance, Innovation, and Risk Management for Agency Use of Artificial Intelligence” (Dec. 4, 2023), https://dho.stanford.edu/wp-content/uploads/OMB_Open_Innovation.pdf [<https://perma.cc/XJS5-M8TE>].

323. The Federal Acquisition Regulation (FAR) itself does not directly apply to states, but many state procurement codes are modeled after FAR. FAR 9.102(b)(1) (2025).

Part 7 requires written acquisition plans to “describe the test program” for each major phase of system acquisition, ensuring agencies formally plan pilot tests before full implementation.³²⁴ FAR Subpart 9.3 allows agencies to require “first article testing” when acquiring new or unproven products—precisely what AI systems represent—allowing agencies to require initial prototypes to be tested against performance criteria before committing to full deployment.³²⁵

Specifically for IT contracting, FAR Part 39 encourages agencies to apply “continuous collection and evaluation of risk-based assessment data” and “prototyping prior to implementation,” with AI systems representing quintessential high-risk IT requiring such evaluation.³²⁶ This Part also mandates modular contracting for IT acquisitions “to the maximum extent practicable,” requiring “delivery, implementation, and testing of workable systems . . . in discrete increments,” directly supporting pilot deployments and phased testing of AI capabilities.³²⁷ More broadly, FAR Part 46 requires that agencies conduct inspection and testing to ensure “supplies or services . . . tendered by contractors meet contract requirements” before acceptance, empowering agencies to subject AI systems to rigorous acceptance testing for accuracy, security, and bias.³²⁸

Beyond the FAR, the Federal Information Technology Acquisition Reform Act requires non-Defense agencies to obtain approval by their Chief Information Officer (CIO) prior to entering any IT contract.³²⁹ For high-risk IT investments, CIOs must identify the causes of risk, the extent to which these causes can be addressed, and the “probability of future success,” actions that are somewhat meaningless for AI procurement without comparative evaluation against baselines.

OMB’s interpretation of acquisition principles goes even further in creating an obligation to evaluate AI tools predeployment. In seeking to “acquire the best solutions at lower cost to the taxpayer,” OMB counsels agencies to “pay careful attention to vendor sourcing, data portability, and long-term interoperability to avoid significant and costly dependencies on a single vendor.”³³⁰ Testing and evaluation is situated as the first step in selecting an AI

324. FAR 7.105(b)(13) (2025). Per FAR 7.102(a), acquisition planning applies to “all acquisitions.”

325. FAR 9.302–9.303 (2025).

326. FAR 39.102 (2025).

327. FAR 39.103 (2025).

328. FAR 46.102 (2025).

329. 40 U.S.C. § 11319(C)(i)(I).

330. OFF. OF MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, MEMORANDUM M-25-22: DRIVING EFFICIENT ACQUISITION OF ARTIFICIAL INTELLIGENCE IN GOVERNMENT 1 (2025).

vendor; the memo requires agencies to “test proposed solutions to understand the capabilities and limitations of any offered AI system,” including by creating a “testing environment” on government systems where possible.³³¹ Moreover, OMB requires agencies to contractually allow for regular “monitor[ing] and evaluat[ion]” of “performance, risks, and effectiveness” of the applicable AI system.³³² These procurement-specific guidelines echo OMB’s general guidance on federal AI use, which counsels agencies to “conduct[] ongoing testing and validation on AI model performance . . . including by testing in real-world conditions” and “consider[] contractual terms that prioritize the continuous improvement, performance monitoring, and evaluation of effectiveness of procured AI.”³³³

States have also released guidance around generative AI that explicitly mandates evaluation. For example, California requires agencies to “prepare data inputs and test models adequately” before deploying a solution in order to “experiment,” “gather feedback,” and “correct outcomes to reduce bias and inaccurate information.”³³⁴ For IT services contracts over \$500,000, California’s procurement code requires departments to perform a “post evaluation”; however, the required form merely asks (with answers in yes or no checkboxes) whether the contracted work was completed on time, whether it was within budget, and whether it “fulfill[ed] all other requirements of the contract including quality standards.”³³⁵ New York specifically calls out use of an AI “without thorough testing to confirm accuracy” as “unacceptable.”³³⁶ However, these state guidance documents are often less

331. *Id.* at 9.

332. *Id.* at 10. The memorandum goes even further in encouraging agencies to act based on regular monitoring of AI system performance: for example, if a new system version fails to meet performance standards, agencies should require vendors to “roll-back to a previous version.” *Id.*

333. OFF. OF MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, MEMORANDUM M-25-21: ACCELERATING FEDERAL USE OF AI THROUGH INNOVATION, GOVERNANCE, AND PUBLIC TRUST 8–9 (2025).

334. CAL. GOV’T OPERATIONS AGENCY, STATE OF CALIFORNIA GENAI GUIDELINES FOR PUBLIC SECTOR PROCUREMENT, USES AND TRAINING 9 (2024).

335. CAL. PUB. CONT. CODE § 12102.3 (West 2025); CAL. DEP’T OF GEN. SERVS., PROCUREMENT DIV., STATE CONTRACTING MANUAL § 21-6: POST EVALUATION FOR IT SERVICES CONTRACTS – 2100.6 (PROCUREMENT DIV.), <https://www.dgs.ca.gov/PD/Resources/SCM/TOC/21/21-6> [<https://perma.cc/65P8-SKZW>]; CAL. DEP’T OF GEN. SERVS. & CAL. DEP’T OF TECH., POST EVALUATION FOR IT SERVICES CONTRACTS, STD. FORM 971 (2020).

336. N.Y. OFF. OF INFO. TECH. SERVS., NEW YORK STATE INFORMATION TECHNOLOGY POLICY NO. NYS-P24-001: ACCEPTABLE USE OF ARTIFICIAL INTELLIGENCE TECHNOLOGIES 7 (2025).

comprehensive than those at the federal level,³³⁷ as are the state procurement codes underlying them.³³⁸

B. *What Meaningful Evaluation Looks Like*

Statutory and constitutional principles militate for pre-adoption evaluation of AI systems in benefits. But what does meaningful evaluation look like in practice? Drawing from the IRS's decades-long experience with systematic evaluation of its telephone assistance program and emerging best practices in AI governance, several core principles should guide agencies' evaluation efforts.

1. *Core Principles*

Scale evaluation rigor with potential impact. Not all AI applications require the same level of scrutiny. The spectrum of AI integration demands graduated evaluation approaches that match the level of oversight to the degree of automation and potential harm.³³⁹ This requires agencies to draw from an evaluation "toolkit" that encompasses different methodologies depending on the intervention's complexity and risk profile. Simple OCR systems that help digitize documents may require only basic accuracy testing using benchmark datasets (such as F1 scores on test sets) and periodic spot-checks through random sampling. Chatbots providing eligibility guidance merit comprehensive evaluation including sandbox testing with realistic user scenarios, accuracy assessment across different query types and user populations, and structured interactions with claims examiners to understand how the tool affects

337. For example, the only mention of testing or evaluation in New York's AI IT policy comes in a list of unacceptable examples. *Id.* at 7.

338. While state procurement codes are generally modeled on FAR as well as the American Bar Association's Model Procurement Code, their testing and evaluation requirements are somewhat less explicit. For example, state procurement codes often have a "responsible bidder provision," which requires vendors to demonstrate "trustworthiness, as well as quality, fitness, capacity, and experience to satisfactorily perform," which could be construed to require meaningful AI evaluation in order to make these determinations. CAL. PUB. CONT. CODE § 1103. California requires its Department of Technology to oversee IT projects by "[e]valuating" them "based on the business case justification, resources requirements, proposed technical solution, project management, oversight and risk management approach, and compliance with statewide strategies, policies, and procedures." CAL. GOV. CODE § 11546. However, none of these provisions comes close to creating a clear textual mandate for continuous evaluation.

339. See Olivia Martin, Faiz Surani, Kit Rodolfa, Amy Perez & Daniel E. Ho, *The Spectrum of AI Integration: The Case of Benefits Adjudication*, in ARTIFICIAL INTELLIGENCE: LEGAL ISSUES, POLICY & PRACTICAL STRATEGIES 45, 51 (2024).

workflow. Auto-adjudication systems demand the most rigorous evaluation protocols from the toolkit: comprehensive bias testing, field experiments or randomized controlled trials to measure real-world performance, ongoing accuracy monitoring through quality assurance frameworks (following models like Mashaw or the Government Performance and Results Act's (GPRA's) approach to performance measurement) and regular human review of algorithmic decisions compared to established benchmarks.³⁴⁰ Agencies must mix and match evaluation tools based on the specific intervention, ensuring the evaluation approach is proportionate to both the degree of automation and the potential consequences of system failure. This graduated, toolkit-based approach ideally prevents bureaucratic paralysis while ensuring adequate protection where stakes are highest.

Establish human baselines wherever possible. A critical component of meaningful evaluation is comparing AI performance against the existing human-operated system, not against abstract benchmarks or vendor promises. As our chatbot audit revealed, a tool may appear to work well in isolation but perform poorly relative to human call center staff, or vice versa. For benefits agencies, this means measuring current human adjudicator accuracy rates, call center response times, and case processing speeds before implementing AI tools, and conducting parallel testing to determine whether AI improves, maintains, or degrades these metrics.

Measure what matters beyond accuracy. While decisional accuracy is paramount, particularly given *Mathews's* focus on minimizing erroneous deprivation, evaluation should encompass the full range of values that procedural due process seeks to protect where possible.³⁴¹ The IRS evaluates not just accuracy of telephone advice, but also timeliness, accessibility, and customer satisfaction.³⁴² Benefits agencies, where feasible, should similarly track whether AI tools improve or hinder claimants' ability to navigate the system, understand their rights, and receive timely decisions. This includes measuring downstream effects: does a chatbot that provides technically accurate but overly complex information actually help claimants complete applications successfully? Do fraud detection systems reduce improper payments without creating insurmountable barriers for legitimate claimants?

Emphasize continuous monitoring over one-time testing. Predeployment testing, however rigorous, cannot always predict long-term performance. AI systems degrade over time due to data drift, changing legal requirements and models,

340. See Mashaw, *supra* note 299, at 774–76; GPRA Modernization Act of 2010, Pub. L. No. 111-352, sec. 2, § 306(a), sec. 4, § 1116, 124 Stat. 3866.

341. *Mathews v. Eldridge*, 424 U.S. 319, 335 (1976).

342. See IRM 21.10.1 (Dec. 20, 2022).

and evolving user needs.³⁴³ The IRS's ongoing evaluation system exemplifies this principle—rather than relying on initial validation, the agency continuously measures and reports performance metrics. Benefits agencies must similarly embed evaluation into their operations, conducting regular audits of AI outputs and maintaining human expertise to assess system performance even as automated tools handle increasing caseloads.

These principles recognize both the potential of AI to improve benefits administration and the constitutional imperative to ensure that technological modernization enhances rather than undermines procedural fairness. Meaningful evaluation provides the bridge between these goals, enabling agencies to harness AI's benefits while maintaining the accuracy, timeliness, and dignity that due process demands.

2. *Example: An Evidence-Based Adoption Cycle for AI Chatbots*

The first step to adopting an AI tool is collecting data to establish a status-quo, often human, baseline.³⁴⁴ For evaluating a chatbot, this could involve measuring the number of calls received per week, call topics, call duration, and client satisfaction for a sample of calls. It is also essential to assess the accuracy of advice given by human agents. For example, the IRS today evaluates client touch points for customer accuracy (“giving the correct answer with the correct resolution”), regulatory accuracy, procedural accuracy, professionalism, and timeliness.³⁴⁵ The IRS records calls and reviews a random sample, reporting that 91.4% of tax law answers and 89.2% of account answers (such as questions about balance due) were accurate from 1.6 million calls in 2023.³⁴⁶ The IRS is distinct from the DOL in conducting such frequent and thorough evaluations of its call centers. Unlike the IRS, the DOL does not mandate regular evaluations of UI call centers, focusing instead on states' timeliness and accuracy in determining benefits, leaving call center

343. See Martin, Surani, Rodolfa, Perez & Ho, *supra* note 339, at 51.

344. Program statutes, such as the Food and Nutrition Act, do create baseline evaluation requirements, as discussed *supra* Part I.A. The Government Performance and Results Act (GPRA), and its 2010 update, also require that agencies collect and report data related to their “performance” goals, which connect to the agency-designed strategic plan. GPRA Modernization Act § 1, 124 Stat. at 3,866 (2011). For further discussion of GPRA's requirements, see *infra* Part IV.B. However, as this Article has shown, agencies are not always conducting performance evaluation on key areas that affect their interactions with citizens, such as the quality of their call centers.

345. IRM 21.10.1, *supra* note 342; DATA BOOK, *supra* note 290, at 24.

346. INTERNAL REVENUE SERV., INTERNAL REVENUE SERVICE DATA BOOK 2023 13–15 (2023), <https://www.irs.gov/pub/irs-prior/p55b--2024.pdf>.

evaluations optional for state administrators.³⁴⁷ Ideally, agencies would also track downstream outcomes such as application completion, grant rates, and appeal rates to understand claimants' full experience.

The next step is to identify an appropriate tool based on which elements of the human baseline are in greatest need of improvement and through user consultation. For example, if most calls are status inquiries, a simplified online portal might be a better investment than a generative AI chatbot. Once a tool is chosen, agencies should run a pilot, testing it on a small, representative user subset.³⁴⁸ Data should be collected on the number of chats initiated, topic of the chat, time per chat, and if available, a measure of client satisfaction. As with the human baseline, a random sample of chatbot responses should be reviewed for accuracy. During the pilot, data collection for the human call center should continue.

Agencies can then compare chatbot performance with human benchmarks but should avoid hasty conclusions. Leveraging the full scope of data is critical to prevent overly simplistic comparisons. We outline examples of how careful analysis of chatbot impact on call volume and assistance accuracy can provide insights that might lead to more informed decisions.

Call Volume. A common goal of chatbots is to reduce burdens on overwhelmed call centers, especially during economic downturns, when high call volumes lead to long wait times and poor claimant experiences. Some states report declines in call center volume after introducing chatbots, but these comparisons are often skewed by background pandemic-induced time trends, as many chatbots were launched during peak demand.³⁴⁹ Future comparisons should adjust for factors like claim volume, seasonality, or broader time trends.

Even with adjustments, agencies should consider the possibility of a “net widening effect.”³⁵⁰ For example, in a pre-chatbot world, if 10,000 individuals called a center each month and chatbots perfectly substituted for calls, we would expect call center volume to drop and chatbot usage to rise equivalently. However, in reality, not everyone who needs help calls due to barriers like limited hours, language issues, or stigma. In this scenario, there might

347. *Id.* at 21–23.

348. Most chatbots surveyed did not appear to have run limited pilots, though they have been run and not publicly announced. The Indiana chatbot states that it is in “beta release” but there is no indication that it was first launched at a smaller scale.

349. Davidson, *supra* note 262.

350. Net widening, often discussed in criminal justice, refers to policies like probation that aim to reduce juveniles in the system but instead expand the number of individuals under its control. *See, e.g.,* Scott H. Decker, *A Systemic Analysis of Diversion: Net Widening and Beyond*, 13 J. CRIM. JUST. 207 (1985).

be 15,000 people with questions, but only 10,000 callers. Chatbots could attract both those who previously called and those who did not, creating broader reach.

An agency observing no change in call center volume after launching a chatbot should not assume failure. The chatbot may have expanded access, reaching claimants who were previously unserved. While this net widening effect may not align with the agency's original goal, it can support broader objectives like increasing eligible applications or reducing follow-up work by clarifying initial questions.³⁵¹

Accuracy. No public data currently exists on the accuracy rates of chatbots or human call centers serving SNAP or UI claimants. However, as this Article's audit shows, chatbots can err. While collecting accuracy data and comparing chatbots to a human baseline is essential, evaluations should not stop there. People may ask different types of questions to chatbots and humans; for example, users might feel freer to ask complex or bizarre questions to chatbots, or reserve such questions for human agents. To make meaningful comparisons, agencies should adjust for differences in query type and complexity. A good starting point could mirror the IRS's segmentation of account-based (e.g., claim status) and eligibility (e.g., complex policy) questions to avoid misjudging relative competencies, though data collection for the human baseline will have to mirror this stratification.

If a chatbot shows a higher error rate after adjusting for query composition, agencies should consider factors beyond their institutional risk tolerance. For example, a higher error rate might become more concerning if chatbots handle a significant share of overall volume, including new users, increasing the absolute impact of errors.³⁵² The level of trust users place in chatbot advice is another key factor. Blank and Osofsky suggest that personalized, instant, and unqualified guidance—like that from chatbots—may lead to higher reliance compared to human agents, though this remains empirically uncertain.³⁵³ Trust likely varies by demographic, making it crucial for agencies to conduct their own studies, such as user questionnaires during pilot testing.³⁵⁴ If

351. Another potential result is more informed claimants such that the need for advice might generally go down.

352. However, there remains the question of whether interacting with the higher error rate chatbot is still better for new claimants than acting with no information at all.

353. Blank & Osofsky, *supra* note 229, at 221.

354. See, e.g., Tae Hyun Baek & Minseong Kim, *Is ChatGPT Scary Good? How User Motivations Affect Creepiness and Trust in Generative Artificial Intelligence*, *TELEMATICS & INFORMATICS*, Sept. 2023 (finding that over 80% of those surveyed trusted ChatGPT and that trust was correlated with perceived personalization); Martin Amoozadeh, David Daniels, Daye Nam, Aayush Kumar, Stella Chen, Michael Hilton, et al., *Trust in Generative AI Among Students: An*

users indeed trust chatbot advice significantly more than human representatives, a higher chatbot error rate would be far less acceptable.

User Satisfaction. While evaluating volume and accuracy rates may cover a narrow conception of the due process balancing implicit in *Mathews*, agencies can and should evaluate the impact of a tool on other core constitutional values like dignity. Although these might seem like lofty goals, carefully designed user surveys and focus groups could help further agencies' understanding of how an AI tool impacts claimant's conception of the program's accessibility, clarity, and empathetic communication.

* * *

As these example findings illustrate, evaluating a new chatbot can lead to different adoption pathways. The simplest decision is if the chatbot outperforms the human baseline in accuracy and either reduces call center volume (a substitution effect) or serves more people without increasing call volume (a net widening effect). The chatbot is a clear win in that case. However, states will likely face tradeoffs between accuracy and the number of customers served. The final decision should focus on how to best serve residents, with this evaluation exercise showing how disaggregated data can guide more informed choices.

Moreover, different AI adoptions will call for different evaluation designs. The implementation of a fraud detection system may focus even more on error rates, especially false positives, as well as demographic disparities and downstream appeal success rates. An AI system assisting adjudicators may instead measure decision quality metrics, time-to-decision, and consistent application across workers.

C. Prospective Implementations of an Evaluation Mandate

We argue that the demands of existing legal structures militate agencies to conduct evaluations when making decisions about significant new integrations of AI systems. Yet admittedly, none of these structures forms a textually explicit and legally binding evaluation mandate, especially in the sense of including a comprehensive human baseline. We explore four prospective

Exploratory Study, SIGCSE 2024: PROC. OF THE 55TH ACM TECH. SYMP. ON COMPUT. SCI. EDUC. V. 1 67–73 (2024) (surveying university students and finding high variance in levels of trust in generative AI). Trust in chatbots has also been a frequent feature of consumer polls. See, e.g., *Poll: Most Who Use Artificial Intelligence Doubt AI Chatbots Provide Accurate Health Information*, KFF (Aug. 15, 2024), <https://www.kff.org/health-misinformation-and-trust/press-release/poll-most-who-use-artificial-intelligence-doubt-ai-chatbots-provide-accurate-health-information/> [<https://perma.cc/CU5Y-WXP9>] (finding most adults are not confident in the accuracy of health information from chatbots).

embeddings of an evaluation mandate: in litigation, in agency guidance and rulemaking, in statute, and in procurement requirements.

1. *Impact Litigation*

As this Article has shown, the courts have stepped into retrospectively address the damages of an AI system gone wrong in Michigan and are also beginning to hold SNAP agencies accountable for their failure to meet constitutional and statutory obligations. These instances build on a long tradition of the courts filling a central role in institutional reform—through what Abram Chayes famously termed “public law litigation”—catalyzing restructuring of school systems, prisons, police departments, and public housing authorities.³⁵⁵ There have long been doubts, however, about Chayes’s optimistic view of the judge’s role in overseeing such sweeping litigation both in terms of institutional competence and separation of powers.³⁵⁶ Judges are trained to adjudicate facts and law, not manage sprawling agencies through the ongoing oversight and monitoring often embedded in injunctive relief. Charles Sabel and William Simon reinterpret Chayes’s vision of public law cases as experimentalist instantiations of “destabilization rights,” or rights to disenfranchise an institution that has systematically failed to meet its obligations.³⁵⁷

One avenue for embedding an evaluation mandate is through impact litigation that explicitly requires agencies to present and defend performance metrics as part of meeting their due process obligations. Judges, after all, are better positioned than agencies to articulate desirable outcomes that are consistent with due process values—they regularly design and implement performance measures in other types of institutional reform cases—and they can defer to agencies’ expertise in how best to satisfy these measures. Indeed, this is precisely what the judge in the Missouri SNAP case did when he ordered the agency to measure and report data on average call wait times and the percentage of denials based on failure to complete timely interviews.³⁵⁸ This approach resonates with the requirement in *Mathews* that agencies demonstrate that they have meaningfully accounted for the risks and complexities of a given AI system. In doing so, due process litigation could move from a backward-looking, detail-bound posture to a forward-focused mechanism of decentralized, court-assisted institutional improvement.

355. Abram Chayes, *The Role of the Judge in Public Law Litigation*, 89 HARV. L. REV. 1281, 1281 (1976); Charles F. Sabel & William H. Simon, *Destabilization Rights: How Public Law Litigation Succeeds*, 117 HARV. L. REV. 1015, 1016 (2004).

356. See Sabel & Simon, *supra* note 355, at 1018.

357. *Id.* at 1055.

358. *Holmes v. Knodell*, 733 F. Supp. 3d 775, 800–01 (W.D. Mo. 2024).

2. *Agency Guidance and Rulemaking*

A more comprehensive way for such an evaluation mandate to be imposed would be through the day-to-day policy guidance that federal agencies issue to guide states in their program administration. This especially makes sense if the evaluation mandate is seen as an interpretation of existing federal law, as we argued above. For example, DOL's demarcation of the inherently government boundary to merit staffing, especially as it pertains to automation, has been expressed through UIPLs—the agency's favored format of policy guidance for UI administration. In one such 1995 letter, the agency describes the directives as “stat[ing] or clarify[ing] the Department's position, particularly with respect to the Department's interpretation of the minimum Federal requirements for conformity or compliance.”³⁵⁹ Unsurprisingly, the agency then asserts that such directives fall within the APA's exception to notice-and-comment requirements for “interpretive rules, general statements of policy, or rules of agency organization, procedure or practice.”³⁶⁰ Yet despite this claim to being an interpretive rule, the directive explicitly asserts to the states that these directives do, in fact, have legal effect. This already seems in tension with APA caselaw that generally suggests that rules that are binding and strip away agency discretion—which seems to be implied by “legal effect”—require notice-and-comment.³⁶¹

Despite the 1995 DOL letter's claims to the contrary, and despite obvious speed and efficiency gains, policy guidance seems to legally be a poor vehicle for an evaluation mandate which, by definition, necessitates some form of binding action on the part of the states. The obvious pivot would then seem to be formal notice-and-comment rulemaking, which would have a binding effect on state agencies. While notice-and-comment rulemaking is admittedly much slower than issuing guidance,³⁶² it may still be quicker and more predictable than legislative change. Moreover, the exact articulation of an

359. U.S. DEP'T OF LAB., EMP. & TRAINING ADMIN., UIPL NO. 01-96 (1995), https://oui.doleta.gov/dmstree/uipl/uipl96/uipl_0196.htm [<https://perma.cc/5S7A-CGTC>].

360. 5 U.S.C. § 553(b)(A). In doing so, it cites a 9th Circuit case in which a UIPL was affirmed to be an interpretive rule under the requirements of the APA. *Rivera v. Becerra*, 714 F.2d 887 (9th Cir. 1983).

361. See *United States v. Texas*, 579 U.S. 547, 547 (2016); *Cnty. Nutrition Inst. v. Young*, 818 F.2d 943, 950 (D.C. Cir. 1987). The 1995 UIPL concedes that the D.C. Circuit has already found part of a UIPL to constitute a substantive rule due to its imposition of “an obligation on the states not found in the statute itself.” *Cabais v. Egger*, 690 F.2d 234, 239 (D.C. Cir. 1982).

362. See Mark Seidenfeld, *A Table of Requirements for Federal Administrative Rulemaking*, 27 FLA. ST. U. L. REV. 533, 536–37 (2000) (identifying up to 109 steps an agency must complete before issuing a final rule).

evaluation mandate raises some difficult questions that might benefit from agency-specific contextual knowledge and expertise. For example, the optimal scoping of an evaluation mandate is non-obvious. The scoping that aligns most clearly with the OMB M-Memo would be to only mandate an evaluation when making a decision on AI adoption.³⁶³ But a mandate with such timing—even if it did include comparison against a human baseline—would asymmetrically add costs to considering AI adoption and thereby put a thumb on the scale, incentivizing the status quo. A rule could balance this by mandating regular evaluation for all baseline systems that significantly affect benefit determinations. However, this would impose even greater costs on the agency.³⁶⁴

3. *New Legislation*

Legislation is uniquely positioned to address the limitations of rulemaking in this context and the issues of funding. Moreover, a statutory mandate helped transform the IRS's telephone evaluation system from a set of reports required for congressional oversight to a core part of the agency's performance evaluation system. For much of the twentieth century, the main evaluation paradigm of the IRS was compliance and enforcement, with traditional measures like audit quotas and case closures used to assess performance.³⁶⁵ The IRS Restructuring and Reform Act of 1998 (RRA), coming on the tails of GPRA in 1993, shifted this focus to quality service and taxpayer experience. The RRA explicitly prohibited the service from using the "records of tax enforcement results" to evaluate employees or even set employee goals.³⁶⁶ Instead, the RRA mandated that the agency "use the fair

363. The M-Memo also calls for ongoing monitoring and evaluation; however, its most full-throated evaluation requirements all lie within the section of practices to be followed "before" using new rights-impacting AI. OFF. OF MGMT. & BUDGET, EXEC. OFF. OF THE PRESIDENT, MEMORANDUM M-24-10: ADVANCING GOVERNANCE, INNOVATION, AND RISK MANAGEMENT FOR AGENCY USE OF ARTIFICIAL INTELLIGENCE (2024), <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> [<https://perma.cc/WB6V-GJKE>].

364. Rulemaking has other downsides. Without congressional appropriation power, rules that impose new mandates operate in a zero-sum world of budgeting tradeoffs. There also may be increased vulnerability to judicial review under *Loper Bright*, as any basis for an evaluation mandate in existing law is likely to be dependent on agency interpretation. *Loper Bright Enters. v. Raimondo*, 603 U.S. 369, 404 (2024).

365. Barry Bozeman, *Risk, Reform and Organizational Culture: The Case of IRS Tax Systems Modernization*, 6 INT'L PUB. MGMT. J. 117, 134 (2003).

366. 26 U.S.C. § 7804.

and equitable treatment of taxpayers by employees as one of the standards for evaluating employee performance.”³⁶⁷

The IRS began incorporating this legislative mandate in its 1999 launch of a “Balanced Measurement System” that sought to measure employee performance based on customer satisfaction (defined as the provision of “accurate and professional services to . . . customers in a courteous, timely manner”), employee satisfaction, and business results (“a productive quantity of work in a quality manner and . . . meaningful outreach to all customers”).³⁶⁸ In 2002, the Balanced Measurement System evolved into the current Embedded Quality Review System,³⁶⁹ which evaluates the service on a number of metrics, including call wait times and accuracy of advice given.³⁷⁰

The Social Security Act, in comparison, only requires that DOL affirm that state agencies administer benefits in a way “reasonably calculated to insure full payment of unemployment compensation when due.”³⁷¹ Through regulation, the DOL established a quality control system to fulfil this mandate, which mandates that states “[c]omplete prompt and in-depth case investigations to determine the degree of accuracy and timeliness . . . with respect to benefit determinations, benefit payments, and revenue collections.”³⁷²

A clearer statutory articulation of the expectation to evaluate potential AI systems against a human or status quo baseline seems promising, just as the RRA’s expectation-setting triggered a transformation of how the IRS conducted performance evaluations. While evaluation is generally a bipartisan issue—the Evidence Act received healthy support from both sides of the aisle—there is a risk that any legislation that survives the political process would be an unfunded mandate. The unfunded nature of the Evidence Act’s evaluation mandates, for example, have sometimes been cited as a reason that agencies’ implementations of its provisions have been so slow.³⁷³

367. *Id.* (Note. The IRS Restructuring and Reform Act of 1998 also required the IRS to develop an employee training plan to train employees on how to provide quality customer service).

368. INTERNAL REVENUE SERV., ORGANIZATIONAL PERFORMANCE MANAGEMENT AND THE IRS BALANCED MEASUREMENT SYSTEM 2–3 (1999), <https://www.irs.gov/pub/irs-pdf/p3561.pdf> [<https://perma.cc/2YEJ-JGU7>].

369. INTERNAL REVENUE SERV., EMBEDDED QUALITY REVIEW SYSTEM (EQRS) – PRIVACY IMPACT ASSESSMENT (2008), <https://www.irs.gov/pub/irs-pia/eqrs-pia.pdf> [<https://perma.cc/KX3F-BAL9>]. The National Quality Review System is a subset of the Embedded Quality Review System.

370. IRS Declining Service, *supra* note 237.

371. 42 U.S.C. § 503(a). One other option would be to update the language of extant merit staffing requirements in legislation like the Social Security Act and the Food and Nutrition Act.

372. 20 C.F.R. § 602 (2024).

373. Carmen Robinson, *Implementing the Evidence Act – What’s Next?*, FED. NEWS NETWORK

To avoid this fate, new legislation could draw inspiration from successful collaborative evaluation models, even those created through administrative action. The Manpower Demonstration Research Corporation (MDRC), a nonprofit, nonpartisan organization created by the Ford Foundation and a group of federal agencies in 1974 offers a promising blueprint for how legislation could structure evaluation infrastructure.³⁷⁴ MDRC was founded to run an ambitious, five-year \$50-million demonstration project to test whether employment programs improved outcomes for disadvantaged populations.³⁷⁵ MDRC has since become known for combining rigorous impact and implementation research with on-the-ground operational expertise to deliver policy-relevant findings to decisionmakers. While MDRC emerged through foundation-government partnership rather than statute, its model suggests key elements that legislation should incorporate: dedicated funding for independent evaluation entities, requirements for collaboration between government agencies and research institutions, and mechanisms for translation findings into policy improvements.

Building on this model, legislation could establish regional AI evaluation centers modeled on MDRC's structure but with statutory authority and guaranteed funding streams. The Workforce Innovation Fund's approach—explicitly funding third-party evaluations and creating a National Evaluation Coordinator to synthesize findings across grantees—demonstrates how federal legislation can scale collaborative evaluation.³⁷⁶ By combining RRA's mandatory evaluation requirements with MDRC's collaborative, independent structure that can incorporate private philanthropic support, new legislation could

(Mar. 28, 2024), <https://federalnewsnetwork.com/commentary/2024/03/implementing-the-evidence-act-whats-next/> [<https://perma.cc/45LZ-QZAZ>]; *The Promise of the Foundations for Evidence-Based Policymaking Act and Proposed Next Steps*, RESULTS FOR AMERICA (2019), <https://results4america.org/wp-content/uploads/2019/09/Evidence-Act-Proposed-Next-Steps-FINAL.pdf> [<https://perma.cc/756A-FUMJ>] (calling for the next version of the Evidence Act to require agencies to set aside a portion of their program funds for evaluations so evaluations are better resourced).

374. *Our History*, MANPOWER DEMONSTRATION RES. CORP., <https://www.mdrc.org/about/history> [<https://perma.cc/CTR9-PPFY>] (last visited Aug. 26, 2025); see also JUDITH M. GUERON & HOWARD ROLSTON, *FIGHTING FOR RELIABLE EVIDENCE* 38 (2013) (comparing the early MDRC to “an intermediary corporation, reflecting its position as a bridge between the interests of many diverse parties—public and private, national and local—and described its key strengths as flexibility, speed, and independence”).

375. MANPOWER DEMONSTRATION RES. CORP., *supra* note 374.

376. ABT ASSOCS., *THE WORKFORCE INNOVATION FUND (WIF): A SYNTHESIS REPORT ON EVALUATION FINDINGS AND EXPERIENCES* (2020), https://www.dol.gov/sites/dolgov/files/ETA/publications/ETAOP2022-25_WIF_Revised_Final_Synthesis_Report_%20Appendices.pdf [<https://perma.cc/XX9W-FLKT>].

ensure that AI evaluation in benefits administration is both rigorous and sustainable, with dedicated resources to match its ambitious scope.

4. Procurement

A final, but perhaps most promising, possible vehicle for an evaluation mandate would be to embed it within procurement requirements, an idea broadly outlined in the OMB's AI procurement memo.³⁷⁷ However, the memo lacks three critical elements for success: a requirement to compare AI performance to a human baseline where possible, longevity across administrations, and resources to fulfill evaluation requirements.

As discussed earlier, a human baseline is critical to defining satisfactory AI system performance. For example, improper payments for state UI agencies often exceed 20%, peaking at 45% during the pandemic.³⁷⁸ Meanwhile, UI call centers face issues like long wait times and low connection rates (fewer than 40% of calls reach live representatives).³⁷⁹ Procurement evaluations must compare both human-led and AI systems to contextualize adoption decisions, as discussed in our example evidence-based adoption cycle.

However, as the Evidence Act illustrates, quality evaluation requires time and resources.³⁸⁰ Since most state AI tools are developed under contract, setting aside part of the project budget for comprehensive evaluations could ensure adequate resourcing. This approach could draw on pay-for-performance frameworks used in medical technology procurement where procurement is costly, and payoffs are uncertain.³⁸¹ One proposal ties contract prices to the evidence quality of the product's efficacy. In the AI context, vendors might compete to provide the strongest proof of improvement against human baselines, validated by agencies. Agencies could also negotiate refunds for tools where performance deviates significantly from initial evaluations. As AI grows more advanced, state AI contracts have increasingly shifted from smaller, more specialized vendors (e.g., Fast Enterprises in Michigan) to larger vendors with more diverse expertise (e.g., Google and Deloitte) who

377. See *supra* Part IV.A.3.

378. See *Unemployment Insurance Payment Accuracy by State*, U.S. DEP'T OF LAB., <https://www.dol.gov/agencies/eta/unemployment-insurance-payment-accuracy> [<https://perma.cc/F9DA-NQA5>].

379. See Dube, Mullainathan & Pope, *supra* note 291.

380. Foundations for Evidence-Based Policymaking Act of 2018, Pub. L. No. 115-435, sec. 101, § 312, 135 Stat. 5529 (2019).

381. See Marianne Hamilton Lopez, Gregory W. Daniel, Nicholas C. Fiore, Aparna Higgins & Mark B. McClellan, *Paying for Value from Costly Medical Technologies: A Framework for Applying Value-Based Payment Reforms*, 39 HEALTH AFFS. 1018, 1024 (2020).

ought to be more capable in meeting these requirements, though their incentives may not align to do so on their own.

To further center evaluation in procurement, agencies should embed testing and comparative analysis into the very structure of AI contracts. Agencies could issue multiple small “trial” awards and require vendors to demonstrate their system’s performance against an institutional baseline (and possibly against each other’s solutions) before moving on to a larger, longer-term contract. This head-to-head testing not only clarifies which vendor is best able to meet an agency’s performance benchmarks, but it also provides critical data for honing the measurement frameworks that will govern the final contract. Agencies might enlist independent evaluators—something already contemplated by scholars in nonprocurement settings³⁸²—to audit these head-to-head trials, thereby adding an additional layer of credibility to the results.

Procurement must also address data ownership and portability. If contracts lock agencies into proprietary systems, evaluation findings may be difficult to act upon. Requiring vendors to provide interoperable data formats or export rights can enable smoother transitions to alternatives if performance benchmarks are not met. By incorporating evaluation, testing, and data portability into contracts, agencies can create a procurement process that ensures rigorous oversight, supports due process, and serves the public interest.

* * *

While procurement represents a powerful and immediate vehicle for embedding evaluation requirements, it cannot stand alone as a comprehensive solution. Procurement-based mandates excel during system development and initial deployment, creating competitive incentives for vendors to demonstrate superior performance and enabling agencies to make evidence-based selection decisions. However, the operational lifecycle of AI systems extends far beyond the initial procurement period. Systems require ongoing monitoring, periodic re-evaluation against evolving baselines, and continuous adjustment as both technology and agency needs change. The temporary nature of most vendor relationships means that procurement mandates, while crucial, must be complemented by other institutional mechanisms—whether through litigation, rulemaking, or legislation—that can ensure sustained evaluation practices throughout a system’s operational lifetime. A robust evaluation framework thus requires multiple, reinforcing vehicles that collectively address the full spectrum of AI system development, deployment, and long-term governance.

382. See, e.g., Citron, *supra* note 20; see CATHY O’NEIL, *WEAPONS OF MATH DESTRUCTION: HOW BIG DATA INCREASES INEQUALITY AND THREATENS DEMOCRACY* (1st ed. 2016).

D. *Reinterpreting Merit Staffing for the AI Age*

1. *The Limitations of the Current Framework*

While merit staffing requirements were conceived as a response to “incompetent and politically dominated personnel” administering benefits in the early twentieth century,³⁸³ merit staff does not guarantee competence in the modern system. Barely half of states meet federal quality guidelines in UI claims adjudication, and 45% of SNAP denials are improperly made.³⁸⁴ This is particularly true given the modern volume and complexity of cases and the difficulty of reaching adequate staffing levels; indeed, requiring merit staffing may undermine quality.³⁸⁵

AI heightens these challenges. While merit staffing targets a host of valuable goals—fair, independent, and efficient administration—its only means of achieving these goals is to assign tasks seen as central to these purposes to merit staff. This creates no space for an impartial examination of whether automation, or some combination of staff and automated tools, may better achieve these goals.

One way some agencies seem to have interpreted merit staffing requirements in the world of AI is as requiring that merit staff make a “final decision” on every case; FNS even published a contractor report explicitly drawing this conclusion.³⁸⁶ Implementing human review only at the final decision stage is undesirable for a number of reasons. At a practical level, such late-stage review does not address issues embedded in the AI tool’s design, such as lack of explainability, making it difficult to identify and correct systemic flaws. This improperly relieves agencies of a necessary pressure to ensure that the tools they are using are interpretable and understandable, both for good design purposes as well as to allow for appellate review. Theoretically, this “final review” interpretation can create the façade of individualized due process without its substance. Instead of mandating that humans are involved at more critical steps in the design and evaluation process, it reduces agency employees to the difficult and bias-prone job of being auditors of the output of black-box AI decisions rather than fulfilling more meaningful and active roles. Within the *Mathews* framework, AI decisions rubber-stamped by humans at the end increase the risk of erroneous deprivation of private rights, whereas ongoing human involvement in design and testing could

383. *But cf.* Diana Moreira & Santiago Pérez, *Civil Service Exams and Organizational Performance: Evidence from the Pendleton Act*, 16 AM. ECON. J.: APPLIED ECON. 250, 255–56 (2024).

384. *See supra* Part I.

385. *See* RICHARDSON & PAHLKA, *supra* note 65.

386. *See supra* Part II.C; *see also supra* Part III.B.3 (stating that Wisconsin and Nevada emphasize that humans review each recommendation before they issue a decision).

enhance procedural fairness by identifying and correcting misunderstandings before decisions are made.³⁸⁷

The merit staffing paradigm provides an important set of values—impartiality, competence, and efficiency—which should remain agencies’ north star when designing systems. However, the binary of dividing tasks between merit and non-merit staff is unable to manage the benefits or risks of AI that can be integrated into benefits administration in a spectrum of ways.³⁸⁸ In the worst case, the merit staffing paradigm encourages states to overly focus on the human role at the final decision stage, which invites cursory review prone to automation bias, and deprives agencies of necessary pressures to inspect and improve their AI’s decisionmaking process.

2. *Building Toward a New Role for Merit Staff*

Agencies can and should reinterpret merit staffing requirements to embrace evaluation-driven AI integration that enhances rather than replaces human oversight. Current DOL guidance already recognizes (in some instances) that automation is permissible where “facts must lead to only one conclusion.”³⁸⁹ While a difficult line to police, it is a good start. Rather than limiting this standard to the narrow set of cases where human disagreement is literally impossible, agencies should interpret “facts must lead to one conclusion” to encompass any determination where rigorous evaluation demonstrates that AI-assisted decisionmaking produces more accurate (or somehow better) outcomes than unassisted human review. This interpretation aligns with the core purpose of merit staffing: ensuring competent, independent, and efficient administration. For example, if evaluation shows that merit staff using AI fraud detection tools achieve 15% error rates compared to 25% rates for unassisted staff, the “facts” of system performance clearly “lead to one conclusion” about which approach better serves program integrity.

This reinterpretation requires reconceptualizing the role of merit staff from individual case processors to system validators and overseers. Rather than requiring all merit staff to personally review every document and make every micro-decision, agencies should also focus merit staff expertise on system design validation, performance monitoring, and exception handling. First, merit staff should be central to designing and testing AI systems before deployment. This includes defining performance benchmarks, identifying potential failure modes, and ensuring that automated systems align with

387. *Mathews v. Eldridge*, 424 U.S. 319, 336 (1976).

388. See Olivia Martin, Faiz Surani, Kit Rodolfä, Amy Perez & Daniel E. Ho, *The Spectrum of AI Integration: The Case of Benefits Adjudication*, in *AI: LEGAL ISSUES, POLICY & PRACTICAL STRATEGIES* 45, 49–50 (2024).

389. ET HANDBOOK 301, *supra* note 30, at 9.

program goals and legal requirements. Second, merit staff should continuously monitor AI system performance against established benchmarks, including by conducting random audits of AI decisions, tracking error rates across different populations, and identifying when human intervention is needed. Finally, merit staff should maintain direct responsibility for complex cases that fall outside AI's capabilities, cases flagged for human attention, and appeals of AI-assisted decisions. This preserves human judgment for the most consequential determinations while allowing automation to handle routine processing.

CONCLUSION

Three years before signing the Social Security Act into law, President Roosevelt famously urged "bold, persistent experimentation," emphasizing the need "to take a method and try it: If it fails, admit it frankly and try another. But above all, try something."³⁹⁰ Many state agencies administering public benefits have heeded this call—facing the dual mandate to ensure civil servant oversight of eligibility determinations while delivering benefits accurately and efficiently, they have cautiously embraced automation. Whether rushing into automation without adequate testing, as Michigan's UI agency did, or clinging to a dysfunctional status quo, as Missouri's SNAP agency did, due process often hangs in the balance.

This Article offers a framework to navigate these challenges. It provides the first systematic analysis of how AI and automation intersect with the legal and historical functions of the civil service. By documenting the risks and opportunities presented by state-level automation experiments, the Article highlights the need for more deliberate oversight and proposes evaluation as the appropriate mechanism to ensure that modernization efforts align with constitutional principles of fairness and accountability.

Mathews provides agencies with a roadmap for balancing private interests, the risk of erroneous deprivation, and administrative burden. For today's era of automation, this guidance translates into rigorous testing and benchmarking—allowing agencies to determine the optimal balance between human and machine decisionmaking to minimize the risk of erroneous deprivation.

This moment of transformation presents both significant opportunities and profound risks. Properly deployed, AI can amplify the effectiveness of civil servants and improve the accuracy and efficiency of public benefits administration. Without sufficient safeguards, however, it risks undermining the very principles it seeks to support. By coupling technological innovation with rigorous evaluation, this Article charts a path forward that preserves the

390. Franklin D. Roosevelt, Governor of N.Y. Oglethorpe University Address (May 22, 1932) (transcript available in the American Presidency Project archive).

foundational values of fairness, accountability, and trust in government. Roosevelt's call to "try something" must be met with an equally resolute commitment to ensure that what we try enhances justice for the people these systems are designed to serve.

APPENDIX TABLE A: RESPONSE SUMMARIES FOR CHATBOT AUDIT*

State	Scope	Low complexity query	Medium complexity query	High complexity query
Arizona ³⁹¹	Dep’t of Econ. Sec.	Could not answer	Could not answer	Could not answer
Arkansas ³⁹²	State	Could not answer	Could not answer	Could not answer
California ³⁹³	Emp. Dev. Dep’t	As expected	Not answered (same response as Q1)	Roughly as expected
Colorado ³⁹⁴	Dep’t of Lab.	As expected	Not answered (generic eligibility)	Not directly answered
Connecticut ³⁹⁵	Dep’t of Lab.	As expected	As expected	Not directly answered (generic eligibility)
Delaware ³⁹⁶	Dep’t of Lab.	Could not answer	Could not answer	Could not answer
Florida ³⁹⁷	Dep’t of Lab.	As expected	As expected	As expected
Georgia ³⁹⁸	Dep’t of Lab.	As expected	Not answered (generic eligibility)	Not directly answered (generic eligibility)

* Chatbot queries and answers listed in this table are on file with the author.

391. ARIZ. DEP’T OF ECON. SEC., *Unemployment Insurance Benefits*, <https://des.az.gov/services/employment/unemployment-individual> [https://perma.cc/6GKW-H8BE] (click on “Need Help? Ask DAVE” at the bottom right of the screen).

392. STATE OF ARK., <https://portal.arkansas.gov/> [https://perma.cc/NZ7P-TZDV] (click on the yellow icon at bottom right of the page).

393. CAL. EMP. DEV. DEP’T, *Unemployment Benefits*, <https://edd.ca.gov/en/unemployment/> [https://perma.cc/QPW8-3NSA] (click on the blue icon at bottom right of the page).

394. COLO. DEP’T OF LAB. & EMP., *Unemployment*, <https://cdle.colorado.gov/unemployment> [https://perma.cc/KA8P-YBF2] (click on the blue icon at bottom right of the page).

395. CONN. DEP’T OF LAB., *Unemployment Benefits and Claim Services*, https://portal.ct.gov/dol/unemployment-benefits?language=en_US [https://perma.cc/86QM-3Z8V] (click on grey icon at bottom right of page).

396. DEL. DEP’T OF LAB., *Division of Unemployment Insurance*, <https://labor.delaware.gov/divisions/unemployment-insurance/> [https://perma.cc/UXY6-3G42] (click on “Ask Ara” at bottom right of page).

397. FLA. COM., *Apply for Benefits*, <https://www.floridajobs.org/reemployment-assistance-service-center/reemployment-assistance/claimants/apply-for-benefits> [https://perma.cc/3FRQ-B6K3] (click on blue icon at bottom right of page).

398. GA. DEP’T OF LAB., <https://dol.georgia.gov/> [https://perma.cc/8M8J-LZES] (click on the green icon at the bottom right of the page).

Indiana ³⁹⁹	State	As expected	As expected	As expected
Iowa ⁴⁰⁰	State	Could not answer	Could not answer	Could not answer
Kansas ⁴⁰¹	State	Could not answer	Could not answer	Could not answer
Maine ⁴⁰²	Dep't of Lab.	Could not answer	As expected	Could not answer
Massachusetts ⁴⁰³	State	As expected	Redirected to overpayment	Not directly answered (generic eligibility)
Michigan ⁴⁰⁴	Dep't of Lab. and Econ. Opportunity	As expected	As expected	As expected
Mississippi ⁴⁰⁵	State	As expected	As expected	As expected
Montana ⁴⁰⁶	Dep't of Lab. and Indus.	Could not answer	Could not answer	Could not answer
New Jersey ⁴⁰⁷	Dep't of Lab.	As expected	Redirected to pandemic relief	Not directly answered (generic eligibility)

399. INDIANA STATE GOV'T, <https://www.in.gov/core/index.html> [<https://perma.cc/4BRJ-5AU9>] (click on blue “ask Indiana” icon in top right of the page).

400. STATE OF IOWA, <https://www.iowa.gov/> [<https://perma.cc/C8VW-JRW7>] (click on teal icon at the bottom left of the page).

401. STATE OF KAN., <https://portal.kansas.gov/> [<https://perma.cc/VU59-DJSW>] (click on blue icon at bottom right of page).

402. ME. DEP'T OF LAB., *Unemployment Insurance*, <https://www.maine.gov/unemployment/> [<https://perma.cc/7HZX-GJCL>] (click on blue icon at bottom right of page).

403. COMMONWEALTH OF MASS., *Apply for unemployment insurance benefits*, <https://www.mass.gov/how-to/apply-for-unemployment-benefits> [<https://perma.cc/AJ9E-VFZN>] (click on green icon at bottom right of page).

404. MICH. DEP'T OF LAB. & ECON. OPPORTUNITY, *Unemployment Insurance Agency*, <https://www.michigan.gov/leo/bureaus-agencies/uia> [<https://perma.cc/83FE-6RKZ>] (click on teal icon at bottom right of page).

405. STATE OF MISS., <https://www.ms.gov/> [<https://perma.cc/ZF2H-VGTP>] (white chatbox on the right).

406. MONT. DEP'T OF LAB. & INDUS., *Unemployment Insurance Division*, <https://uid.dli.mt.gov/> [<https://perma.cc/8G8G-RW82>] (click on blue icon at bottom right of page).

407. STATE OF N.J. DEP'T OF LAB. & WORKFORCE DEV., *Division of Unemployment Insurance*, <https://www.nj.gov/labor/myunemployment/> [<https://perma.cc/W2V9-WUP4>] (click on the icon at the bottom right of the page).

New York ⁴⁰⁸	Dep't of Lab.	Could not answer	Could not answer	Could not Answer
North Carolina ⁴⁰⁹	Dep't of Com.	As expected	As expected	Not directly answered (generic eligibility)
Ohio ⁴¹⁰	Dep't of Job and Fam. Servs.	Redirected to employer resources	Redirected to work authorization	Could not answer
South Dakota ⁴¹¹	State	As expected	As expected	As expected
Texas ⁴¹²	State	When asked about unemployment insurance, immediately connected to a live human agent	NA	NA
Utah ⁴¹³	Workforce Servs.	Could not answer	Could not answer	Could not answer
Virginia ⁴¹⁴	Emp. Comm'n	As expected	Could not answer	Could not answer
Wisconsin ⁴¹⁵	Dep't of Workforce Dev.	Roughly as expected (link but no description)	Generic response about application process with no mention of documents	Not directly answered (generic eligibility)

408. N.Y. DEP'T OF LAB., *Unemployment Insurance*, <https://dol.ny.gov/unemployment/unemployment-insurance-assistance> [<https://perma.cc/R7CL-UMWM>] (click on teal icon at bottom right of page).

409. N.C. DEP'T OF COM., *Apply for Unemployment*, <https://www.des.nc.gov/individuals/apply-ui> [<https://perma.cc/BH9C-G6LY>] (click on green icon at bottom right of page).

410. OHIO DEP'T OF JOB & FAM. SERVS., *Unemployment Insurance*, <https://jfs.ohio.gov/job-services-and-unemployment/unemployment/overview> [<https://perma.cc/D995-5BFP>] (click on blue and yellow icon at bottom right of page).

411. S.D. GOV'T CITIZEN PORTAL, <https://www.sd.gov/cs> [<https://perma.cc/6HKF-AYLU>] (click enter site; then click on yellow bird icon on right of page).

412. STATE OF TEX., *Texas by Texas*, <https://www.texas.gov/texas-by-texas/> [<https://perma.cc/6UUR-E4SP>] (click on the red icon on the bottom right of the page). Texas also recently launched what they term a "personal, portable government assistant" called "Texas by Texas (TxT)" which simplifies processes like renewing vehicle registration. *Id.*

413. UTAH WORKFORCE SERVS., *Unemployment Insurance Benefits*, <https://jobs.utah.gov/ui/home> [<https://perma.cc/VM43-BT5Y>] (click on blue icon on bottom right of page). Note that we tried to access the chatbot of various times within Utah's business hours on different days and were unable to get responses.

414. VA. EMP. COMM'N, *Unemployment*, <https://www.vec.virginia.gov/unemployment> [<https://perma.cc/U3MP-MKBT>] (click on orange icon on bottom right of page).

415. WIS. DEP'T OF WORKFORCE DEV., *Unemployment Insurance*, <https://dwd.wisconsin.gov/ui/> (click on cow icon on bottom right of the page).

APPENDIX TABLE B: RESPONSE SUMMARIES FOR GENERATIVE AI CHATBOT AUDIT**

Question	Indiana	Michigan	Florida	South Dakota	Mississippi
How much money do I have to have earned in the base period to be eligible for unemployment insurance?	Correct	Correct	Correct	Cannot verify	Incorrect: stated that you must have earned at least \$780 in the base period, whereas the requirement is to have earned \$780 in the highest quarter
I have earned \$7000 in the base period with \$3000 in one quarter. Am I eligible for unemployment insurance?	Correct	Correct	Incorrect: incorrectly asserts that \$7,000 is not greater than \$4,500	Cannot verify	Correct
If I file my initial claim for when I was unemployed between January 5, 2020, and April 4, 2020, what is my base period?	Correct	Incorrect: returns the last four quarters of the base period rather than the first four	Correct	Unanswered	Correct
What is my base period if I file my claim for when I was unemployed between October 7, 2018 and January 5, 2019?	Correct	Correct	Correct	Incorrect: temporally implausible response	Correct
How do I calculate my maximum benefit amount?	Correct	Correct	Correct (though wrong link)	Unanswered	Correct
I know for a fact that if you are really dissatisfied with your job and leave because of that, you are still eligible	Correct	Correct	Unanswered	Correct	Correct

** Chatbot queries and answers listed in this table are on file with the author.

for unemployment insurance. I left my job because my boss made me very unhappy. I'm still eligible for unemployment insurance, right?					
There was a labor dispute at my job. How does this affect my claim for unemployment insurance?	Correct	Unanswered	Correct	Correct	Correct
Under what conditions am I deemed unavailable for work, as a matter for unemployment insurance, if I'm on vacation?	Correct	Unanswered	Correct	Unanswered	Correct
How many work search activities am I required to participate in per week?	Incorrect: stated three search activities whereas only two are required. ⁴¹⁶	Correct	Correct	Unanswered	Correct
I think I got overpaid on my unemployment insurance benefits. Do I have to repay the agency if I didn't do anything wrong?	Correct	Correct	Correct	Correct	Correct

416. See IND. HANDBOOK, *supra* note 213, at 22.