

# Strengthening the Reliability, Resilience, and Integrity of Public-Benefit Information Systems: A Literature Review of Data-Driven Approaches in U.S. State Workforce and Public-Benefit Agencies

<sup>1</sup>Favour Onyebuchi Favor, <sup>2</sup>Hannah Adeniji

<sup>1</sup>University of Illinois Springfield

fudeh2@uis.edu

<sup>2</sup>University of Illinois System | Office of Medicaid Innovation

hanpeters700@gmail.com

---

## ARTICLE INFO

Received: 01 Oct 2023

Revised: 18 Nov 2023

Accepted: 28 Nov 2023

## ABSTRACT

State workforce and public-benefit agencies administer the programs that millions of Americans rely on in moments of acute need, among them unemployment insurance (UI), the Supplemental Nutrition Assistance Program (SNAP), Temporary Assistance for Needy Families (TANF), and Medicaid. The information systems that deliver these services are often decades old, built for ordinary rather than crisis-level demand, and fragile under stress. The COVID-19 pandemic made this fragility plain. Claim volumes surged, aging systems faltered, and tens of billions of dollars were paid improperly even as eligible workers waited weeks for assistance. This paper reviews the governmental, scholarly, and technical literature on the reliability, resilience, and integrity of these systems and on the supervised and unsupervised data-mining methods proposed to support them. The review is organized around four themes: the fragility of legacy benefit infrastructure; resilience and cybersecurity; payment integrity and fraud; and the applied analytics used to address these problems, including demand forecasting, anomaly detection, and claimant segmentation. The literature points to a common lesson. Data-driven models offer real operational value, but that value is limited by predictive accuracy, by the uneven cost of false positives, by how well models hold up when conditions change, and by the procedural safeguards that surround their use. The paper closes with implications for treating reliability, resilience, and integrity as related design goals rather than separate ones.

**Keywords:** unemployment insurance; public-benefit administration; legacy information systems; resilience; payment integrity; fraud detection; machine learning

---

## . Introduction

Public-benefit and workforce programs serve as the operational core of the U.S. social safety net, and the software that runs them is, for practical purposes, critical public infrastructure. Much of that software predates the people who now depend on it. A substantial number of state UI systems still run on COBOL, a programming language first standardized in the late 1950s, inside mainframe environments built around the assumptions of an earlier era, including hard limits on the number of claims that could be

processed in a single day (The Century Foundation [TCF] et al., 2022; WHY, 2020). When demand departs sharply from those assumptions, the systems do not fail gently. They fail outright.

The pandemic provided a clear stress test. As businesses closed in the spring of 2020, UI claims rose to levels without precedent, and at least one state governor publicly appealed for volunteer COBOL programmers to keep an overwhelmed UI system running. That same state had already abandoned a prior modernization effort at a cost of roughly \$350 million (WHY, 2020). Analyses of the period stress a single point: the code itself still worked, but it had been built under outdated design constraints and was maintained by a shrinking pool of programmers, often with little documentation (TCF et al., 2022; WHY, 2020).

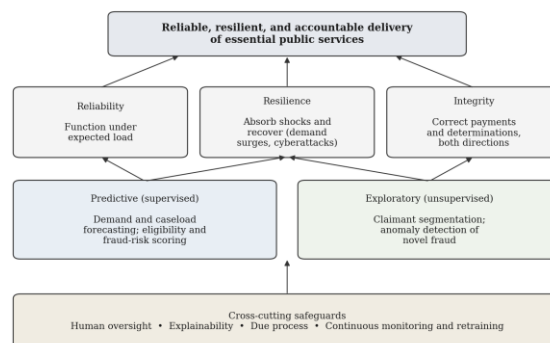
This review takes up the lens named in its title, reliability, resilience, and integrity, because that lens mirrors how the federal oversight community frames the problem. The U.S. Government Accountability Office (GAO) placed UI on its High-Risk List in 2022 and built its assessment around program design, infrastructure, and integrity risks (GAO, 2022). Reliability concerns whether a system performs its intended function under expected conditions. Resilience concerns whether it can absorb shocks, whether a demand surge or a cyberattack, and recover. Integrity concerns whether the funds and determinations moving through it are correct, both in catching wrongful payments and in not wrongly penalizing legitimate claimants. This paper synthesizes what is known across these three dimensions and assesses the role that data-driven methods can and cannot play.

### 1.1 Guiding Questions

Recasting the predictive and exploratory inquiry that often motivates applied data-mining work, this review is organized around two questions. The first is a predictive question: can data-driven models reliably predict the quantities that benefit agencies must act on in near real time, such as demand and caseload, eligibility, and the risk that a payment is improper, with enough accuracy to support operational decisions without harming the people those systems serve? The second is an exploratory question: how can unsupervised techniques such as cluster analysis and anomaly detection surface homogeneous groups of claims, or instead statistically unusual ones, in ways that improve oversight, service design, and equitable decision-making? The two questions follow the supervised and unsupervised structure of data-mining research, in which supervised models forecast a target of interest while unsupervised methods reveal hidden structure. Figure 1 sets out the framework that organizes the review.

Figure 1

A Framework Linking Data-Driven Methods, Cross-Cutting Safeguards, and the Three Pillars of Public-Benefit System Performance



*Note.* Predictive (supervised) and exploratory (unsupervised) methods support the three pillars only when paired with the cross-cutting safeguards shown at the base of the figure.

### 2. Reliability: Legacy Systems and the Limits of Brittle Infrastructure

The most heavily documented body of work concerns reliability, the question of whether benefit systems can perform their core function under realistic load. The recurring finding is that they often cannot, for reasons that are structural rather than incidental. The proximate cause most often cited is legacy technology. A large share of state UI claims systems depend on COBOL, the surrounding hardware and integrations were never designed for surge volumes, and the workforce able to maintain this code keeps shrinking as veteran programmers retire (TCF et al., 2022; WHY, 2020). The deeper difficulty, as analysts point out, is that these systems were built under constraints, such as a ceiling on daily claims, that no longer match reality, and that institutional knowledge of them is thin (WHY, 2020).

Modernization is therefore widely recommended, yet the literature is candid about its difficulty. State efforts have run over budget and been abandoned, and even successful migrations must eventually be repeated as today's mainstream technologies become tomorrow's legacy (WHY, 2020). The most constructive synthesis comes from a joint report by The Century Foundation, the National Employment Law Project, and Philadelphia Legal Assistance, which documents how modernization reshapes the claimant experience and argues that user-centered design and implementation matter as much as the underlying platform (TCF et al., 2022). The reliability literature thus recasts modernization from a one-time procurement into a continuing capability that an agency must sustain, with adaptability and surge capacity treated as primary requirements.

### 3. Resilience: Absorbing Shocks and Recovering Function

Where reliability concerns ordinary operation, resilience concerns behavior under shock. The literature treats two shock types as primary: demand surges and cyberattacks. The two are related, since resilience in both cases is the capacity to absorb a disturbance, keep delivering essential function, and recover. The demand dimension appears in Sections 2 and 5; the security dimension has a growing literature of its own.

Government systems increasingly appear in the literature as high-value targets. The Federal Bureau of Investigation (2022) reported that local government entities were the second most frequently victimized group, behind only academia, among ransomware incidents reported to it during 2021, and warned that the public's dependence on government-run utilities, emergency services, and benefit programs makes these entities attractive to attackers. Industry survey data tell a similar story. Roughly six in ten state and local government organizations reported a ransomware attack in 2021, up from about one third the year before, and most of those attacked had data encrypted (Sophos, 2022). Tracking of disclosed incidents counted on the order of one hundred state or municipal entities affected in 2022 (Emsisoft, 2023).

The National Institute of Standards and Technology offers the most widely adopted structuring framework. Its Cybersecurity Framework organizes resilience around five functions: identify, protect, detect, respond, and recover. Its ransomware profile applies these functions to ransomware specifically and includes approaches for recovering when an attack on data integrity succeeds (NIST, 2022). The framework's idea of comparing a current profile against a target profile gives benefit agencies a transferable way to make resilience measurable rather than aspirational. The practical literature ties these threads back to legacy infrastructure. Brittle, poorly understood systems are at once hard to scale under demand and

hard to defend or patch under attack, so resilience and reliability investments tend to reinforce one another (NIST, 2022; WHYY, 2020).

#### 4. Integrity: Improper Payments, Fraud, and the Hazards of Automated Adjudication

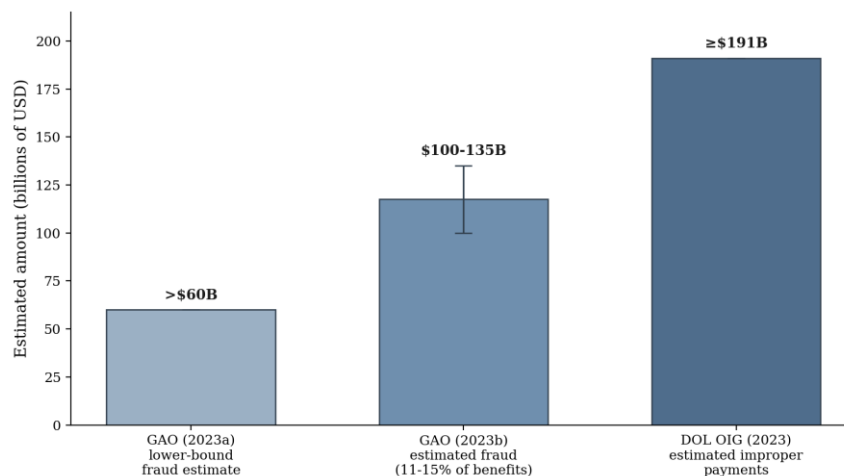
The integrity literature is the most quantitatively developed and, for the purposes of this paper, the most cautionary, because it documents failure in both directions: money paid that should not have been paid, and people penalized who should not have been penalized.

##### 4.1 The Scale of Improper Payments and Fraud

The size of the improper-payment problem is large and contested. The GAO (2023b) estimated that fraud accounted for roughly 11 to 15 percent of pandemic UI benefits, which implies losses on the order of \$100 to \$135 billion between April 2020 and May 2023, while cautioning that the true figure will likely never be known with certainty. An earlier analysis put fraud at more than \$60 billion as a conservative lower bound and faulted the U.S. Department of Labor for not having adopted an antifraud strategy aligned with the agency's own Fraud Risk Framework. That report also noted that estimated improper payments in the regular UI program had risen well above their roughly \$8 billion pre-pandemic level (GAO, 2023a). The Department of Labor's Office of Inspector General applied a 21.52 percent improper-payment rate to about \$888 billion in pandemic UI expenditures and estimated that at least \$191 billion may have been paid improperly, with a significant share attributable to fraud (U.S. DOL OIG, 2023). The placement of UI on the High-Risk List shows that payment integrity is a structural design requirement rather than an afterthought (GAO, 2022). Figure 2 sets these estimates side by side and notes the differing categories they measure.

Figure 2

Selected Federal Estimates of Pandemic Unemployment Insurance Fraud and Improper Payments, 2022 to 2023



*Note.* The estimates differ in method and scope. The GAO figures estimate fraud, whereas the DOL OIG figure estimates improper payments, a broader category that includes but is not limited to fraud. The GAO (2023b) range revises and supersedes the earlier GAO (2023a) lower-bound fraud estimate rather than adding to it, so the two GAO values are not independent observations. Sources: GAO (2023a, 2023b); U.S. DOL OIG (2023).

### 4.2 The Other Side of Integrity: Wrongful Determinations

If the fraud literature argues for stronger detection, an equally important body of work warns against pursuing detection in ways that harm claimants. The clearest example is Michigan's MiDAS system. Between October 2013 and September 2015, the Michigan Integrated Data Automated System adjudicated tens of thousands of fraud cases by algorithm with little or no human review. A later review found that of roughly 40,195 cases decided by the system alone, about 85 percent were wrong (Charette, 2021). The human consequences were severe, including wage garnishment, ruined credit, bankruptcies, and lost homes, and the episode led to extended litigation, a halt to fully automated adjudication, statutory reform, and a multimillion-dollar settlement (Benefits Tech Advocacy Hub, n.d.; Charette, 2021; Time, 2020; Undark, 2020). The lasting lesson is that a system built on the presumption that claimants are guilty until proven innocent, running on incomplete or corrupted data and lacking meaningful human oversight or due process, can cause widespread harm (Time, 2020; Undark, 2020). Integrity, then, must be defined on both sides, reducing both wrongful payments and wrongful denials, and the uneven and often irreversible cost of false accusations should weigh heavily in any automated decision.

## 5. Applied Data-Mining Methods: Forecasting, Detection, and Segmentation

The fourth body of work concerns the analytic methods proposed to advance the three pillars. It maps onto the supervised and unsupervised division and onto the model families common across applied data mining, including regression and tree-based models, neural networks, clustering, and anomaly detection (Ngai et al., 2011).

### 5.1 Supervised Learning: Prediction and Scoring

Supervised methods dominate two applications. The first is fraud and risk scoring, in which models train on historical labeled outcomes to estimate the probability that an application or claim is improper. Reviews of the field describe logistic regression, random forests, support vector machines, neural networks, and ensemble methods as the standard tools for this task, and they stress that the availability of high-quality, consistently labeled data is the binding constraint and that performance is judged on measures such as precision, recall, and the F1 score rather than accuracy alone (Al-Hashedi & Magalingam, 2021; Carcillo et al., 2021; Ngai et al., 2011). Government applications are documented as well. A Brookings Institution (2022) analysis catalogs algorithmic fraud-detection systems at the Internal Revenue Service and the Securities and Exchange Commission, among them a program that scores tax returns for refund fraud. The second application, demand and expenditure prediction, is covered under forecasting below.

### 5.2 Unsupervised Learning: Anomaly Detection and Segmentation

Unsupervised methods address two complementary needs. The first is detection of new fraud that no labeled rule anticipates. Clustering groups similar cases so that outliers stand out, while anomaly-detection techniques flag records that depart from normal patterns without requiring labeled examples, which allows the discovery of schemes not seen before (Carcillo et al., 2021; Ngai et al., 2011). Because these methods do not depend on past labels, the literature treats them as essential complements to supervised scoring, which can only recognize fraud that resembles what it has already seen (Carcillo et al., 2021). The second use is segmentation for service design and oversight, grouping claims or claimants into homogeneous clusters to inform staffing, outreach, and tailored processing, which answers the exploratory question posed in Section 1.1.

### 5.3 Demand Forecasting

Forecasting connects to resilience, since readiness for surges depends on anticipating load. UI claims data are themselves strong real-time indicators of labor-market conditions. An Indiana analysis found that a one percent rise in continued UI claims corresponded on average to roughly a 0.6 percent rise in monthly unemployment, which makes claims a timely input to near-term forecasts (Indiana Business Research Center, 2020). For means-tested programs, state forecasting offices use components-of-change models that break caseload into its underlying drivers and draw on the economic sensitivity of programs such as SNAP, whose participation tends to move inversely with employment (Oregon Office of Forecasting, Research and Analysis, 2019). The pandemic exposed the limits of such relationships. SNAP and TANF caseloads rose by 3.3 million between March and June 2020, their largest quarterly increase on record, yet were far less responsive to unemployment than in the pre-pandemic period, as policy choices and health conditions came to dominate (Hembre, 2023).

### 5.4 Key Methodological Lessons

Three findings recur across this literature, and together they bound what any deployment can achieve. The first is that accuracy and the cost of false positives cannot be separated. A model that looks best on an aggregate error measure may still be wrong often enough to be dangerous when its outputs trigger consequential and hard-to-reverse actions against individuals. The MiDAS error rate is the extreme illustration, and it argues for keeping a human in the loop on any adjudicatory decision and for designing models that can be explained, so that investigators can examine their reasoning (Charette, 2021). The second is that models degrade when conditions change, which is the problem of generalization. The clearest documented case is the pandemic's disruption of the long-standing relationship between labor-market conditions and benefit caseloads, a shift that would have caused a model trained on the earlier period to forecast poorly at the very moment accuracy mattered most (Hembre, 2023). This is a direct argument for continuous monitoring and retraining. The third is that interpretation of features matters as much as prediction. Knowing which variables drive a model's output is what turns a forecast into an actionable insight for program design, and in adjudicatory settings it is what makes a determination contestable and therefore lawful (Ngai et al., 2011).

## 6. Discussion: Synthesis, Tensions, and Gaps

Considered together, these four bodies of work describe a system under simultaneous pressure to become more reliable, more resilient, and more accountable, with data-driven methods proposed as a partial remedy for each. A few points stand out. The three pillars are related rather than separate. Brittle legacy infrastructure weakens reliability and resilience at the same time, and it also weakens integrity, because systems that cannot adapt invite rigid automated controls to be added on top, which is how an integrity initiative such as MiDAS turns into a source of harm (Charette, 2021; GAO, 2022). Investments that improve adaptability therefore help across all three pillars. Table 1 summarizes the reviewed literature by pillar.

**Table 1**

*Summary of the Reviewed Literature by Theme*

<b>Theme pillar</b>	<b>or</b>	<b>Representative sources</b>	<b>Type of evidence</b>	<b>of</b>	<b>Principal contribution</b>
Reliability: legacy infrastructure		TCF et al. (2022); WHY (2020)	Policy report; journalism		Legacy COBOL systems and a failed modernization explain the fragility exposed in 2020; modernization should be user-centered and ongoing.
Resilience: surge and cybersecurity		FBI (2022); Sophos (2022); Emsisoft (2023); NIST (2022)	Government notice; industry survey; framework		Government is a frequent ransomware target; NIST functions and current-versus-target profiling make resilience measurable.
Integrity: fraud and improper payments		GAO (2022, 2023a, 2023b); U.S. DOL OIG (2023)	Federal oversight reports		Large and contested pandemic estimates; UI sits on the GAO High-Risk List, marking payment integrity as a core requirement.
Integrity: wrongful determinations		Charette (2021); Time (2020); Undark (2020); Benefits Tech Advocacy Hub (n.d.)	Case study; journalism		Automated adjudication without human review produced mass wrongful fraud findings, showing the cost of false positives.
Applied methods		Ngai et al. (2011); Carcillo et al. (2021); Al-Hashedi & Magalingam (2021); Brookings (2022); IBRC (2020); Oregon OFRA (2019); Hembre (2023)	Reviews; applied studies		Supervised scoring, unsupervised detection and segmentation, and demand forecasting; accuracy, false-positive cost, drift, and interpretability are decisive.

*Note.* UI = unemployment insurance; GAO = U.S. Government Accountability Office; NIST = National Institute of Standards and Technology; DOL OIG = U.S. Department of Labor, Office of Inspector General; TCF = The Century Foundation; IBRC = Indiana Business Research Center; OFRA = Office of Forecasting, Research and Analysis.

The central tension is between integrity and access. The fraud literature pushes toward stronger controls, while the wrongful-determination literature pushes toward protecting legitimate claimants and preserving access. The most workable balance in the sources is procedural. Automation should triage cases and direct human attention rather than replace human judgment in consequential determinations, and it should be explainable enough to support due process (Charette, 2021; GAO, 2023a; Undark, 2020).

The literature also shows clear gaps. Rigorous, peer-reviewed evaluation of state-operated benefit systems remains thin compared with the financial-sector fraud literature, and many state analytic systems are vendor-built and opaque, which frustrates independent assessment (Charette, 2021). Standardized governance for validating, monitoring, and retraining models in benefit administration is underdeveloped, and questions of equity and bias receive less study than technical performance. Resilience is rarely made measurable. The current-versus-target profiling method offers a transferable approach, but its application to benefit systems specifically remains largely unexamined (NIST, 2022).

### 7. Implications for Practice and Research

The literature supports several design commitments for programs that aim to strengthen the reliability, resilience, and integrity of public-agency information systems. The reliability literature implies that modernization should be treated as a continuing capability, with adaptability and surge capacity set as primary requirements and user-centered design valued alongside the platform itself (TCF et al., 2022). The resilience literature implies that agencies should adopt a recognized framework such as the NIST Cybersecurity Framework, make resilience measurable through current-versus-target profiling, and treat demand forecasting as an input to capacity planning (Indiana Business Research Center, 2020; NIST, 2022).

The integrity literature implies a two-sided goal, reducing both wrongful payments and wrongful denials, pursued through analytics that triage and explain rather than adjudicate on their own, with mandatory human review, due-process safeguards, and explicit attention to the uneven cost of false positives (Charette, 2021; GAO, 2023b). The methods literature implies that any deployed model should be monitored and retrained against drift, validated for generalization beyond its training period, and designed for interpretability from the start (Hembre, 2023; Ngai et al., 2011). These commitments reinforce one another. An adaptable system is easier to secure and to correct, an explainable model is easier to govern and to defend, and a program that measures resilience and watches model performance can catch its own decline before that decline becomes a public failure.

### 8. Conclusion

Across reliability, resilience, and integrity, the evidence tells a consistent story. The information systems that U.S. state workforce and public-benefit agencies depend on are essential public infrastructure that has too often been allowed to decay, and the pandemic exposed the cost of that neglect in failed systems, billions of dollars in improper payments, and delayed relief for people in need. Data-driven methods offer real and demonstrated value, whether supervised models for forecasting and risk scoring or unsupervised methods for anomaly detection and segmentation. The literature is equally clear that this value is conditional. It is bounded by accuracy and the uneven cost of error, by the tendency of models to fail when conditions change, and by the procedural safeguards that decide whether automation serves the public or harms it. The most promising path forward treats reliability, resilience, and integrity as a single, connected design problem and keeps human judgment, explainability, and due process at the center of any system entrusted with decisions about people's livelihoods.

### References

- [1] Al-Hashedi, K. G., & Magalingam, P. (2021). Financial fraud detection applying data mining techniques: A comprehensive review from 2009 to 2019. *Computer Science Review*, 40, 100402.
- [2] Benefits Tech Advocacy Hub. (n.d.). *Michigan unemployment insurance false fraud determinations*. <https://www.btah.org/case-study/michigan-unemployment-insurance-false-fraud-determinations.html>
- [3] Brookings Institution. (2022). *Using AI and machine learning to reduce government fraud*. <https://www.brookings.edu>

- [4] Carcillo, F., Le Borgne, Y.-A., Caelen, O., Kessaci, Y., Oblé, F., & Bontempi, G. (2021). Combining unsupervised and supervised learning in credit card fraud detection. *Information Sciences*, 557, 317–331.
- [5] Charette, R. N. (2021, June 24). *Michigan's MiDAS unemployment system: Algorithm alchemy that created lead, not gold*. IEEE Spectrum. <https://spectrum.ieee.org/michigans-midas-unemployment-system-algorithm-alchemy-that-created-lead-not-gold>
- [6] Emsisoft. (2023). *The state of ransomware in the US: Report and statistics 2022*. <https://www.emsisoft.com>
- [7] Federal Bureau of Investigation. (2022, March 30). *Ransomware attacks straining local US governments and public services* (Private Industry Notification 20220330-001). Internet Crime Complaint Center. <https://www.ic3.gov/CSA/2022/220330.pdf>
- [8] Hembre, E. (2023). Examining SNAP and TANF caseload trends, responsiveness, and policies during the COVID-19 pandemic. *Contemporary Economic Policy*, 41(2), 262–281. <https://doi.org/10.1111/coep.12596>
- [9] Indiana Business Research Center. (2020). *Predicting unemployment from unemployment insurance claims*. <https://www.ibrc.indiana.edu>
- [10] National Institute of Standards and Technology. (2022). *Ransomware risk management: A Cybersecurity Framework profile* (NIST IR 8374). <https://doi.org/10.6028/NIST.IR.8374>
- [11] Ngai, E. W. T., Hu, Y., Wong, Y. H., Chen, Y., & Sun, X. (2011). The application of data mining techniques in financial fraud detection: A review and future research directions. *Expert Systems with Applications*, 38(3), 2262–2275.
- [12] Oregon Office of Forecasting, Research and Analysis. (2019). *Caseload forecast methodology*. Oregon Department of Human Services & Oregon Health Authority.
- [13] Sophos. (2022). *The state of ransomware in state and local government 2022*. <https://www.sophos.com/en-us/partner-news/2022/09/resources/the-state-of-ransomware-in-state-and-local-government-2022>
- [14] The Century Foundation, National Employment Law Project, & Philadelphia Legal Assistance. (2022). *Centering workers: How to modernize unemployment insurance technology*. <https://tcf.org>
- [15] Time. (2020, May 28). *States' automated systems are trapping citizens in bureaucratic nightmares with their lives on the line*. <https://time.com>
- [16] Undark. (2020, June 1). *Government's use of algorithm serves up false fraud charges*. <https://undark.org>
- [17] U.S. Department of Labor, Office of Inspector General. (2023). *Oversight of the unemployment insurance program*. <https://www.oig.dol.gov/doloiguioversightwork.htm>
- [18] U.S. Government Accountability Office. (2022). *Unemployment insurance: Transformation needed to address program design, infrastructure, and integrity risks* (GAO-22-105162). <https://www.gao.gov/products/gao-22-105162>
- [19] U.S. Government Accountability Office. (2023a). *Unemployment insurance: DOL needs to address substantial pandemic UI fraud and reduce persistent risks* (GAO-23-106586). <https://www.gao.gov/products/gao-23-106586>

- [20] U.S. Government Accountability Office. (2023b). *Unemployment insurance: Estimated amount of fraud during pandemic likely between \$100 billion and \$135 billion* (GAO-23-106696). <https://www.gao.gov/products/gao-23-106696>
- [21] WHYY. (2020, April 4). *Why N.J. wants coders fluent in a 60-year-old language in the middle of a pandemic*. <https://whyy.org>