AI-Powered Rules as Code

Government Briefing | April 23, 2025

%



GEORGETOWN UNIVERSITY

DBN Community Norms + Code of Conduct

- Being kind and respectful to each other
- Making space and taking space
- Assuming good will
- Staying curious beeck center.org/dbncode
- Being mindful and reflective
- Being open and honest
- Taking the best, leaving the rest

beeckcenter.org/dbncode

The Digital Benefits Network is dedicated to providing a harassment-free experience for everyone, We do not tolerate harassment of participants in any form.

If you have any questions or concerns, please email us at digitalbenefits@georgetown.edu.

Agenda

Welcome

Research Briefing

• Highlights from our research publication using AI for RaC translation for SNAP and Medicaid (started as part of the Policy2Code Prototyping Challenge)

Open Discussion









Digital Benefits Network

Supports government in delivering public benefits services and technology that are accessible, effective, and equitable in order to ultimately increase economic opportunity.

SNAP | WIC | Medicaid/CHIP | TANF | Basic Income | UI | Child Care

digitalbenefitsnetwork.org



About the Massive Data Institute

At Georgetown's McCourt School of Public Policy, the **Massive Data Institute (MDI)** is an interdisciplinary research institute that connects experts across computer science, data science, public health, public policy, and social science to tackle societal scale issues and impact public policy in ways that improves people's lives through responsible evidence-based research.



mdi.georgetown.edu



OECD Working Papers on Public Governance No. 42 Cracking the code: Rulemaking for humans and machines

OECD

Rules as Code Definition

"an official version of rules (e.g., laws and regulations) in a machine-consumable form, which allows rules to be understood and actioned by computer systems in a consistent way."

> - Organisation for Economic Co-operation and Development

> > beeckcenter.org/oecdRAC



Benefits Eligibility Rules as Code



Source: Beeck Center for Social Impact + Innovation and the Massive Data Institute, Georgetown University, 2025 GEORGETOWIX beeckcenter D

Digital Benefits NETWORK



Our Rules as Code Work

Community	beeckcenter.org/dbncommunity			
Research	Initiatives / Events			
Upcoming: State Eligibility & Enrollment Systems	Quarterly: Ongoing roundtable sessions			
2025: AI-Powered Rules as Code	2025: Program design for prototyping open syntax &			
2024: Cross-Sector Insights Report from our Rules	code library			
as Code Community	2024: Policy2Code Prototyping Challenge + Demo Day			
2023: Exploring Rules Communication	beeckcenter.org/Policy2Code			
2022: Benefit Eligibility Rules as Code	2022: Rules as Code Demo Day			
beeckcenter.org/dpolicy				

GEORGETOWX beeckcenter Digital Benefits NETWORK MASSIVE DATA

AI-Powered Rules as Code

Experiments with Public Benefits Policy

Can generative AI be used to expedite the translation of policies into software code for implementation in public benefits eligibility and enrollment systems under a Rules as Code approach?

beeckcenter.org/AI_RaC



Digital Benefits

EORGETOWEX UNIVERSITY Court School of Public Policy

Policy2Code Prototyping Challenge



Summer meetups: June - August 2024

12 teams

Demo Day: September 2024

beeckcenter.org/Policy2Code





GEORGETOWX UNIVERSITY McCourt School of Paster Paster McCourt School of Paster Paster

Policy2Code Prototyping Challenge



Hoyas Lex Ad Codex

3 Georgetown Centers Beeck Center

Massive Data Institute

Center for Security and Emerging Technology (CSET)

Faculty, students, staff collaboration

beeckcenter.org/Policy2Code



beeckcenter





AI-Powered Rules as Code: Research Overview





GEORGETOWEX UNIVERSITY MoCourt School of Public Public INSTITUTE

AI-Powered Rules as Code: Research Overview

Experiment 1	Experiment 2	Experiment 3	Experiment 4	
Asking eligibility questions to chatbots and API	Asking eligibility questions to API guided by policy document	Prompting to write benefits eligibility rules using a template with and without policy documents	Prompting to generate software code for benefits eligibility	
How well can LLM chatbots answer general SNAP and Medicaid eligibility questions based on their training data and/or resources available on the internet? What factors affect their responses?	How well can an LLM generate accurate, complete, and logical summaries of benefits policy rules when provided official policy documents?	How well can an LLM extract machine-readable rules from unstructured policy documents in terms of output relevance and accuracy? How does the use of a structured rules template impact an LLM's ability to produce relevant and accurate output?	How effectively can an LLM generate software code to determine eligibility for public benefits programs?	



beeckcenter social impact + innovation





Experiment 1 Methodology

Ask eligibility questions and assess the outputs



Generate prompts across different benefits scenarios and levels of specificity (e.g. individual eligibility criteria, employment, household, income).

"Who is eligible for Medicaid in Texas?"



GEORGETOWN UNIVERSITY beeckcenter

Evaluators score

Diaital Benefits

NETWORK

NSTITUTE

Input prompts into

Experiment 1 Results Summary

Ask eligibility questions and assess the outputs



- No major difference in the three models' performance.
- Perform reasonably well on state-specific /current and relevance, but big drop off for completeness



social impact + innovation





Experiment 1 Results Summary

Ask chatbots eligibility questions and assess the outputs



Average Scores By state and program

 Medicaid responses scored slightly higher than SNAP





Digital Benefits

GEORGETOWX UNIVERSITY McCourt School of Paulice Paulice

Experiment 1 Materials

Digital	Library Topics ~	Trending About Get Involved	a Q						
Benefits Use Cases	Experiment 1 Materia	ls							
Conclusion Key Takeaways Potential for Future	Access Experiment 1 Materials								
Experimentation Get in Touch	Spreadsneet: Prompts with scores Spreadsheet: Cumulative scores Rubric Experiment 1 Bubric	Spreadsheet: Prompts with scores and comments Spreadsheet: Cumulative scores Rubric							
Citation									
Thank You	Current and state-specific applicability	Is the information from the response current and state-specific	Score 1-5						
Appendix	Completeness	Is the response thorough and does it cover all elements requested in the prompt?	Score 1-5						
Experiment 1 Materials Experiment 2 Materials Experiment 3 Materials	Relevance	Is the response focused on the question, without adding irrelevant or unnecessary details?	Score 1-5						
Experiment 4 Materials	The following serves as an overview	of the assessment rubric and explains v	vhat each score means:						
Footnotes	Current, and state-specific information S points: All information provided is accurate, current, and state-specific								

- Spreadsheet: Prompts with scores and comments
- Spreadsheet: Cumulative scores
- Rubric







Experiment 2 Methodology

Collect policies, feed them to the LLM and assess the outputs



Extract and consolidate policy documents.



Generate prompts for specific goals

(e.g. plain language summaries, logic requests, eligibility criteria)

State Websites Policy Manual PDFs



"Based on this document, what are the eligibility criteria for SNAP in Texas?"



Build a RAG environment to capture knowledge about state specific policies for LLM to use.



Evaluators score responses 1-5 points based on developed rubric.



beeckcenter







Experiment 2 Assessment of Policy Documents



Extract and consolidate policy documents.





Key Challenges:

- Separate PDF documents (Georgia)
- Scanned documents make text unsearchable and unusable for LLMs (California)
- Inconsistent formats or interactive content that's hard to download (Alaska, Pennsylvania)

Observations:

- Option to download full PDF, easy access (Michigan, Texas)
- Separated or scanned non-readable documents (Georgia, California)
- Good navigation, but harder to extract as a single PDF (Alaska, Pennsylvania)





Experiment 2 Prompt Development and Response Assessments



Generate prompts for

specific goals (e.g. plain language summaries, logic requests, eligibility criteria)

"Based on this document, what are the eligibility criteria for SNAP in Texas?"

State	Policy	Question on Topic	Prompt	State	Policy	Question on Topic	Prompt
PA	SNAP	"Plain language" summary	Please provide a plain language summary of the eligibility criteria for Medicaid in Pennsylvania found in the provided document.	PA	SNAP	Eligibility Scenarios	I am a single adult who cannot work due to a medical condition in Pennsylvania. Based on the provided document, what would determine if I am eligible for SNAP?

beeckcenter

social impact + innovation





Experiment 2 Results Summary

Collect policies, feed them to the LLM and assess the outputs

Prompt Effectiveness for SNAP and Medicaid in Pennsylvania



 Generally returns accurate results; however, scores drop off significantly for completeness and relevance.



How did your state perform? See the report!

Source: Beeck Center for Social Impact + Innovation and the Massive Data Institute, Georgetown University, 2025



Digital Benefits

GEORGETOWX UNIVERSITY McCourt School of Public Policy INSTITUTE

Experiment 2 Common Errors Found in Responses

Accuracy

Common errors:

- Misinterpreting criteria (like age or group eligibility)
- Confusing recent vs. still-applicable older policies.

Completeness

Common errors:

- Missing details for eligibility,
- Focus on admin information (timelines, verification)

Relevance

Common errors:

- Off-Topic Answers
- Irrelevant pieces of Information





beeckcenter



Experiment 2 Materials

	Experiment 2 Ma	terials	
Conclusion	Access Experiment 2 M	aterials	
Key Takeaways			
Potential for Future			
Experimentation	 Spreadsheet: Prompts with 	n scores and comments	
	 Spreadsheet: Cumulative s 	cores	
Get in Touch	Rubric		
	Experiment 2 Pubric		
Citation	Experiment 2 Rubite		
Thank You	Accuracy	Is the information from the response accurate?	Score 1-5
Appendix	Completeness	Is the response thorough and doe it cover all elements requested in the prompt?	Score 1-5
Appendix Experiment 1 Materials	Completeness	Is the response thorough and doe it cover all elements requested in the prompt?	Score 1-5
Appendix Experiment 1 Materials Experiment 2 Materials	Completeness	Is the response thorough and doe it cover all elements requested in the prompt?	Score 1-5

- Spreadsheet: Prompts with scores and comments
- Spreadsheet: Cumulative scores
- Rubric







Experiment 3 Methodology

Rules generation





Prompt LLM to generate rules under different conditions



```
program: SNAP
  state: "<STATE>"
  eligibility:
    income_limits:
      gross: "<AMOUNT>"
      net: "<AMOUNT>"
    asset_limit: "<AMOUNT>"
    citizenship: ["<STATUS>"]
    work_requirements:
10
      age_range: ["<AGE>", "<AGE>"]
11
      hours: "<HOURS>"
12
13
  benefits:
14
    max_allotment: "<AMOUNT>"
15
    standard_deduction: "<AMOUNT>"
16
17
  certification:
18
    period: "<MONTHS>"
19
20
21 reporting:
    simplified: "<BOOLEAN>"
22
```

Simplified SNAP Rules Template









Experiment 3 Experimental Design

Develop rules template, prompt LLM to generate rules with and without policy documents Approach Approach



beeckcenter

Diaital Benefits

NETWORK

DATA

NSTITUTE

Court School of Public Policy

Experiment 3 Results Summary

Develop rules template, prompt LLM to generate rules with and without policy documents

Average Rules Generation Accuracy

By state and category as percentages

	Georgia		Penr	nsylvania	Texas		
	LLM	LLM+RAG	LLM	LLM+RAG	LLM	LLM+RAC	
Income	0	100	0	100	0	100	
Asset Limit	0	100	0	50	0	100	
Categorical Eligibility	75	100	100	100	100	100	
Citizenship & Residency	100	100	100	100	0	100	
Monthly Allotment	20	100	0	100	20	40	
Work Requirements	100	100	100	100	100	100	
Deductions	50	100	21	64	29	100	
Certification Period	50	0	50	100	50	100	

- Plain prompting yielded unreliable generation accuracy for state-specific rules.
- RAG enhanced the alignment of generated rules with policy documents.
 - Structured templates are essential for rules extraction.
 - Rules generation can provide a structure for code generation.





Experiment 3 Materials

œ	experiment 3_generated_rules 🕁 🗈 🖉 File Edit View Insert Format Data Tools	b Extensions Help				Ū		Share -
	९ ५ ८ वि ही 100% → \$ % .0ੵ .0약	123 Arial 🝷	- 10 + B Z -S	<u>A</u> <u>A</u> = 5	ē - ≣ - ↓	• ə • <u>A</u> • cə <u>+</u>	ω Υ 🖷 - Σ	^
A1	✓ ∫fic							
	A	В	С	D	E	F	G	н
1			Georgia			Pennsylvan	ia	
2	Rule	GPT-40 V	GPT-40 + RAG 🔹	ground truth 💌	GPT-40 •	(GPT-40 + RAG V	ground truth 💌	GPT-40
3	Eligibility_Income_Gross_Monthly_Limit_1_Person	1473	1580	1580	1473	1580	1580	
4	Eligibility_Income_Gross_Monthly_Limit_2_Person	1984	2137	2137	1984	2137	2137	
5	Eligibility_Income_Gross_Monthly_Limit_3_Person	2495	2694	2694	2495	2694	2694	
6	Eligibility_Income_Gross_Monthly_Limit_Additional_Person	511	557	557	511	557	557	
7	Eligibility_Income_Net_Monthly_Limit_1_Person	1133	1215	1215	1133	1215	1215	
8	Eligibility_Income_Net_Monthly_Limit_2_Person	1526	1644	1644	1526	1644	1644	
9	Eligibility_Income_Net_Monthly_Limit_3_Person	1920	2072	2072	1920	2072	2072	
10	Eligibility_Income_Net_Monthly_Limit_Additional_Person	394	429	429	394	429	429	
11	Eligibility_Asset_Limit_Standard	2500	2750	2750	2500	2750	5500	
12	Eligibility_Asset_Limit_Elderly_Disabled	3750	4250	4250	3750	9000	9000	
13	Eligibility_Categorical_Eligibility_Applies	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE	TRUE
14	Eligibility_Categorical_Eligibility_Qualifying_Programs	TANF, SSI	TANF, SSI, Work Support Payments, TANF Community Outreach Services	TANF,SSI,WSP	TANF, SSI	TANF, SSI	TANF, SSI	TANF, SSI
15	Eligibility_Citizenship	US Citizen, Qualified Non-Citizen	U.S. citizen, U.S. National, eligible immigrant	US Citizen, Qualified Non-Citizen	US Citizen, Qualified Non-Citizen	U.S. Citizen, Qualified Non-Citizen	U.S. Citizen, Qualified Non-Citizen	US Citizen, Qua Alien
16	Eligibility_Residency_Requirement	Must be a resident of Georgia	Must be a resident of Georgia	Must be a resident of Georgia	Must be a resident of Pennsylvania	Must be a resident of Pennsylvania.	Must be a resident of Pennsylvania.	Must be a reside Texas
17	Eligibility_Work_Requirements_Age_Range_Min	16	16	16	16	16	16	
18	Eligibility_Work_Requirements_Age_Range_Max	59	59	59	59	59	59	
19	Eligibility_Work_Requirements_Weekly_Hours	20	20	20	20	20	20	
					Elderly,			

Spreadsheet: Contains the generated rules from both the plain prompting and RAG approaches, compared against the verified SNAP Policy values in Georgia, Pennsylvania, and Texas.







Experiment 4 Methodology

Code generation





3

Prompt LLM to generate code under different conditions

Evaluate code

Step 6: Documentation and Structure # (This is done through comments in the code) # Step 7: Rule Storage # Save the decision-making rules in a JSON file # [STRENGTH] Externalized rule storage shows understanding of maintainability rules = ("exemptions": [{"age less than": 16}, {"age between": [16, 17], "head of au": False}, {"age greater equal": 60}, {"tanf recipient": True}, {"physically unfit": True}, ("mentally unfit": True), {"pregnancy_status": True, "pregnancy_risk": True}, {"chronic homelessness": True}, {"disability status": True}, {"student_status": True, "student_enrollment": True, "student_type": ["GED", "home_school", "high_school", "online_courses", "post secondary", "seminary", "trade school"] }, {"caretaker child": True}, ("caretaker incapacitated": True), {"employment_hours_greater_equal": 30}, ("employment earnings greater equal": 7.25 * 30), {"vista_volunteer": True}, {"migrant seasonal worker": True}, {"drug treatment program": True}, {"ucb status": True}, {"homeless status": True}, ("veteran status": True), ("in foster care": True) 1, "registrants": [{"age between": [16, 59]}







Experiment 4 Experimental Design

Generate code using simple and detailed LLM prompts and iterative steps



Experiment 4 Results Summary

Generate code using simple and detailed LLM prompts and iterative steps

Performance of Designs Across Criteria

Criteria	Design 1 Simple Prompt	Design 2 Detailed Prompt	Design 3 Iterative Prompts
Variable Identification	Partial	Good	Good
Input Handling	Poor	Partially good	Poor
Output Correctness	Incorrect	Partially correct	Incorrect
Decision Making	Poor	Partially good	Poor (mechanical)
Logical Consistency	Poor	Improved	Partially consistent
Rule Coverage	Partially covered	Improved	Partially covered
Code Execution	Runs (unstable)	Runs (Improved)	Doesn't run

- At a high level, the detailed prompt design was the most successful.
- Summarized policy guidance for code generation reduces code errors.
- Modular design is particularly important in LLM workflows.





Experiment 4 Materials

```
# Step 6: Documentation and Structure
# (This is done through comments in the code)
# Step 7: Rule Storage
# Save the decision-making rules in a JSON file
# [STRENGTH] Externalized rule storage shows understanding of maintainability
rules = {
    "exemptions": [
        {"age less than": 16),
        {"age between": [16, 17], "head of au": False},
        {"age greater equal": 60},
        {"tanf recipient": True},
        {"physically unfit": True},
        {"mentally unfit": True},
        {"pregnancy status": True, "pregnancy risk": True},
        {"chronic homelessness": True},
        {"disability status": True},
        {"student_status": True, "student_enrollment": True, "student_type": ["GED", "home_school", "high_school", "online_courses",
"post secondary", "seminary", "trade school"]},
        {"caretaker child": True},
        {"caretaker incapacitated": True},
        {"employment hours greater equal": 30},
        {"employment_earnings_greater_equal": 7.25 * 30},
        {"vista volunteer": True},
        {"migrant seasonal worker": True},
        {"drug treatment program": True},
        {"ucb status": True},
        {"homeless status": True},
        {"veteran status": True},
        {"in foster care": True}
   1.
    "registrants": [
        {"age between": [16, 59]}
```

- Prompts: Includes the system and user prompts used for code generation. For the iterative approach, there are three sets of prompts for each stage, as outlined in the experiment details.
- Generated Code: Contains the code results from each approach, with our manual annotations highlighting strengths [STRENGTH] and weaknesses [WEAKNESS].

Diaital Benefits

NETWORK

GEORGETOWN

McCourt School of Public Policy

INSTITUTE

UNIVERSITY



beeckcenter

Key Takeaways: Rules as Code Generation

- LLMs can help support the Rules as Code pipeline.
- Humans in the loop: Accuracy and equity considerations must outweigh efficiency in high-stakes benefits systems.
- Authoritative sources: It is possible to improve the performance of the benefits-related responses by pointing LLMs to authoritative sources like policy manuals.

- State governments can make it easier for LLMs to use their policies by making them digitally accessible.
- Asking LLMs to write policy code directly leads to inconsistent code quality.
- Code generation can be improved by using an LLM with RAG to generate machine-readable policy rules, but requires a manually-curated template.





Key Takeaways: Impacts on People

• Mixed LLM results have a direct impact on people seeking or receiving benefits, **risking incorrect information** when they ask generative AI models questions about programs like SNAP and Medicaid.

• When AI models provide incorrect information, they often do so in a **confident tone,** which can mislead those without subject expertise.







Future Experimentation

- Repeat methodologies for other states, programs, policies or LLMS.
- Repeat for different programming languages for outputs.
- Pair between policy experts and software engineers to program and evaluate.
- Explore specific program or policy LLM

- Compare rules against existing systems
- Explore extracting code from legacy systems.
- Compare code writing efforts to policy analysis efforts.
- Design a prototype toolkit to test different rules and code generation prompts for SNAP.





AI-Powered Rules as Code



Full report

Summary + Key Takeaways

Rubrics Prompts Results Scores Code Samples

beeckcenter.org/AI_RaC

beeckcenter

social impact + innovation

Source: Beeck Center for Social Impact + Innovation and the Massive Data Institute, Georgetown University, 2025





GEORGETOWAX UNIVERSITY McCourt School of Packies Packey

AI-Powered Rules as Code

Experiments with Public Benefits Policy

Questions?







Open Discussion

GEORGETOWN UNIVERSITY



R

ີ່

Apps

Guide Questions

- How have you been using AI with policy and program implementation?
- Are there any new wins? New challenges?
- Do you have questions for this group?



Survey Rate the Briefing





Thank you!



@digitalbenefits.bsky.social



Digital Benefits Network

rulesascode@georgetown.edu







