



HAL
open science

Can GPT-3 Perform Statutory Reasoning?

Andrew Blair-Stanek, Nils Holzenberger, Benjamin van Durme

► **To cite this version:**

Andrew Blair-Stanek, Nils Holzenberger, Benjamin van Durme. Can GPT-3 Perform Statutory Reasoning?. ICAIL 2023: Nineteenth International Conference on Artificial Intelligence and Law, Jun 2023, Braga, Portugal. pp.22-31, 10.1145/3594536.3595163 . hal-04360081

HAL Id: hal-04360081

<https://hal.science/hal-04360081>

Submitted on 21 Dec 2023

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Can GPT-3 Perform Statutory Reasoning?

Andrew Blair-Stanek
ablair-stanek@law.umaryland.edu
University of Maryland School of Law
Baltimore, Maryland, USA

Nils Holzenberger
nils.holzenberger@telecom-paris.fr
Télécom Paris
Institut Polytechnique de Paris
Palaiseau, France

Benjamin Van Durme
vandurme@jhu.edu
Johns Hopkins University
Baltimore, Maryland, USA

ABSTRACT

Statutory reasoning is the task of reasoning with facts and statutes, which are rules written in natural language by a legislature. It is a basic legal skill. In this paper we explore the capabilities of the most capable GPT-3 model, text-davinci-003, on an established statutory-reasoning dataset called SARA. We consider a variety of approaches, including dynamic few-shot prompting, chain-of-thought prompting, and zero-shot prompting. While we achieve results with GPT-3 that are better than the previous best published results, we also identify several types of clear errors it makes. We investigate why these errors happen. We discover that GPT-3 has imperfect prior knowledge of the actual U.S. statutes on which SARA is based. More importantly, we create simple synthetic statutes, which GPT-3 is guaranteed not to have seen during training. We find GPT-3 performs poorly at answering straightforward questions about these simple synthetic statutes.

CCS CONCEPTS

• **Applied computing** → Law; • **Computing methodologies** → **Natural language generation**; **Reasoning about belief and knowledge**.

KEYWORDS

natural language processing, reasoning, law, statutes, GPT-3

ACM Reference Format:

Andrew Blair-Stanek, Nils Holzenberger, and Benjamin Van Durme. 2023. Can GPT-3 Perform Statutory Reasoning?. In *Nineteenth International Conference on Artificial Intelligence and Law (ICAIL 2023)*, June 19–23, 2023, Braga, Portugal. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3594536.3595163>

1 INTRODUCTION

There has been great excitement about the emergent capabilities of large language models such as GPT-3 [3] and ChatGPT on legal tasks. For example, Bommarito and Katz [2] show that GPT-3 performs at a passing rate on evidence and torts questions on the U.S. multistate bar exam. In Choi et al. [4], law professors had ChatGPT take their final exams and blindly graded its output along with their human students' answers; it achieved low but passing grades.

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

ICAIL 2023, June 19–23, 2023, Braga, Portugal

© 2023 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 979-8-4007-0197-9/23/06...\$15.00

<https://doi.org/10.1145/3594536.3595163>

Section 1001. Definition of Rolang.

(a) General rule

The term "rolang" means-

(1) any soultratesly, or

(2) any parkinse.

(b) Soultratesly

The term "soultratesly" means-

(1) any oxideney, or

(2) any chastiles.

(c) Parkinse

The term "parkinse" means-

(1) any portle, or

(2) any frestes.

Alexis is a portle. Is section 1001(b) applicable to Alexis? Let's think step by step.

Section 1001(a) defines rolang as either soultratesly or parkinse.

Section 1001(b) defines soultratesly as either oxideney or chastiles.

Section 1001(c) defines parkinse as either portle or frestes.

Since Alexis is a portle, section 1001(b) is applicable to Alexis.

Figure 1: Example of GPT-3 incorrectly answering a simple question about a short synthetic statute. The text generated by GPT-3 is highlighted.

We set out to test how the most powerful GPT-3 model currently exposed through public APIs would perform on one of the most basic tasks required of lawyers, *statutory reasoning*. Statutes are legal rules written by legislative bodies in natural language. Statutory reasoning is the application of those rules to facts that are also in natural language [13]. We first consider the StAtutory Reasoning Assessment (SARA) dataset, a handcrafted computational benchmark covering nine sections of the U.S. tax code [8]. We find that straightforward application of GPT-3 leads to significant gains compared to prior work on this task, but that GPT-3 makes a number of clear mistakes (§3).

We then ask why this could be, starting with the question of how much GPT-3 already knows about the current U.S. Code (§4). We find that it has some imperfect knowledge of the U.S. Code. Then we detail (§5) an experiment based on novel synthetic statutes to measure how well GPT-3 can answer simple statutory-reasoning questions when prompted with a single fact (see Figure 1 for an example). When presented with statutory language that is easily interpretable by humans, we find that GPT-3 will regularly make mistakes. Our results suggest statutory reasoning as an area of

interest for new AI research, as a challenge to motivate future improvements in large language models.

2 RELATED WORK

Statutory reasoning. The SARA dataset [8] is a benchmark for statutory reasoning. Efforts to solve SARA through machine reading methods have had limited success [9]. SARA’s curated statutes, facts, and questions were designed to avoid semantic ambiguity. But semantic ambiguity also poses challenges for statutory reasoning, such as the U.S. Supreme Court case *Nix v. Hedden* about whether tomatoes were “vegetables” in a customs statute. There have been several approaches for using NLP to handle statutory reasoning involving such semantic ambiguity [18, 19]. The task of determining whether one statute entails another one, called *statute entailment*, is the closest legal NLP task to statutory reasoning. The yearly COLIEE challenge [16] includes statute entailment. BERT-based models [5] have been most successful on COLIEE’s statute entailment, until they were recently surpassed by GPT-3-based models [23]. Legal judgment prediction is another task related to statutory reasoning, and is also a significant challenge for current machine reading models [1, 12, 26].

Prompting GPT-3. GPT-3 is a large-scale language model, which has shown a surprising ability to solve NLP tasks without the need to fine-tune its parameters [3]. This is achieved by providing GPT-3 with a set of natural language instructions, specifying what task is to be solved, along with a few input-output examples for the task. This prompt-based approach has been further refined, in particular by augmenting the expected output with human-written natural language explanations of how that output was derived from the input [22]. Kojima et al. [11] show that GPT-3 can be prompted to produce these explanations itself, providing a substitute for human annotation. In particular, the phrase “Let’s think step by step” was found, out of a number of prompts, to maximize performance. We prompt GPT-3 in multiple ways, taking inspiration from Wei et al. [22] and Kojima et al. [11].

Reasoning with GPT-3. Exploring GPT-3’s ability to reason has yielded improvements on multiple benchmark datasets. Lu et al. [14] collect a dataset of common sense and science questions, pairing each question with background information and an explanation of the answer. The additional context, serving as a chain-of-thought, marginally improves GPT-3’s performance. Zhou et al. [27] explicitly decompose reasoning problems, forcing GPT-3 to solve a complex problem in multiple increments. Khot et al. [10] show GPT-3 can be prompted to decompose problems into simpler steps, on both artificial tasks and open-domain question-answering. Zelikman et al. [24] exploit a language model’s ability to generate chains of reasoning: GPT-J [21] is made to generate chains of reasoning for a set of questions and answers, and is further fine-tuned on those chains that lead to the correct answer. Adding structure to the process of natural language generation, Zhang et al. [25] run a symbolic reasoning engine on the output of a language model, yielding improvements on semi-synthetic benchmarks meant to test reasoning abilities. Thorburn and Kruger [20] evaluate GPT Neo’s argumentative abilities, by making it suggest claims or reasons to support a claim, with no particular emphasis on logical reasoning.

Data comes from Kialo, a collection of arguments curated by online users. Evaluations show that reasons and claims generated by GPT Neo are less coherent than those generated by humans. Saparov and He [17] systematically study GPT-3’s ability to reason, using synthetic, controlled data, and metrics that measure the internal coherence of chains of reasoning. The authors show that GPT-3 has the ability to perform individual steps of reasoning, but tends to lose track when producing a proof with multiple steps. Results further show that GPT-3 leverages its background knowledge about the world to draw conclusions, so that GPT-3’s ability to reason decreases as the overlap between the target domain and its real-world knowledge decreases.

GPT-3 on legal tasks. Hendrycks et al. [7] use GPT-3 to solve legal multiple-choice questions, with data collected from online sources. Results show that the largest GPT-3 model performs significantly better than random, but still far less than expert human performance. Yu et al. [23] experiment with various ways of prompting GPT-3 for statute entailment. The authors find that GPT-3 does much better than previous BERT-based models. In particular, prompts designed with inspiration from legal reasoning work best. Bommarito and Katz [2] test GPT-3’s ability to answer multiple choice questions, involving short paragraphs of context, from the U.S. multistate bar exam. Prompting GPT-3 to rank the available answer choices gets close to a passing score on some subject areas. In the medical domain, Gutierrez et al. [6] study GPT-3’s abilities on two biomedical information extraction tasks. The authors show that fine-tuning a BERT model on a small training set consistently outperforms few-shot learning with GPT-3.

3 SARA

3.1 Data and Experimental Setup

SARA stands for StAtutory Reasoning Assessment. It consists of nine statutory sections drawn from the U.S. tax code (i.e. title 26 of the U.S. Code) and 376 hand-crafted cases that state simple facts, involving taxpayers named Alice, Bob, Charlie, and Dan, and ask one question that can be answered by applying some of the nine sections to the facts. Obviously, the nine statutes, the facts, and the questions are all in natural language. But SARA also includes a translation of the statutes, facts, and questions to Prolog. Solving the questions via Prolog results in 100% accuracy.

Of the 376 cases, 100 are pure “tax cases” and ask for how much tax an individual had to pay, with answers typically in the many thousands of dollars. Because GPT-3 performs poorly in reasoning with such numbers, we do not consider these 100 cases. All the remaining 276 cases ask for either “Entailment” or “Contradiction”. Of these, 193 (121 train, 72 test) involve numbers. The remaining 83 (55 train, 28 test) involve no numbers at all. For all our tests, we break out our results separately on the 72 test cases with numbers and the 28 test cases with no numbers.

In all experiments throughout this paper, for our calls to GPT-3, we set the temperature to 0.0 and the top_P to 1.0.¹ These settings serve two purposes. First, they maximize reproducibility. Second,

¹All our code is at github.com/BlairStanek/gpt-statutes.

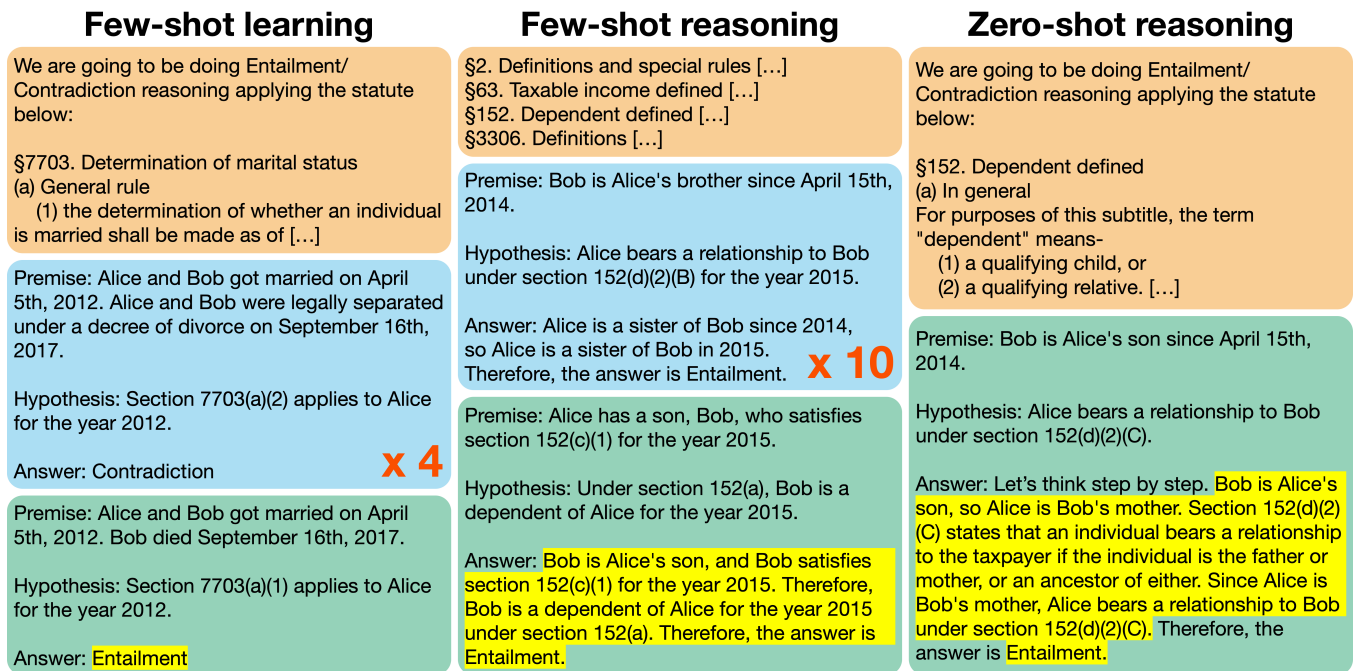


Figure 2: Prompts used in Section 3 to pose SARA test cases to GPT-3. Top boxes, in orange, contain statutes; depending on the setting, statutes may or may not be included. Example cases are in blue; in zero-shot there are no example cases. At the bottom, in green, are test cases. Text highlighted in yellow is generated by GPT-3. If GPT-3's first response is not clear, we prompt a second time with "Therefore the answer is", following Kojima et al. [11].

they minimize the risk that GPT-3 will wander off topic or hallucinate. In SARA, every question is annotated with a single correct answer.

3.2 Three Approaches

We try three approaches to prompting GPT-3 to answer the SARA questions.

First, we use **4-shot dynamic** prompting. For each test case, we picked the 4 most similar training cases. For these 4, we gave the full text of the case, labeling the "Premise", "Hypothesis", and "Answer", the latter being either "Entailment" or "Contradiction". Because SARA is balanced, 2 of these 4 training cases had "Answer: Entailment" and the other 2 had "Answer: Contradiction". Then we appended the "Premise" and "Hypothesis" of the test case, followed by the word "Answer". This is illustrated in the left column of Figure 2.

Second, we use **zero-shot** prompting. There we use no training cases at all. We just put the test case's "Premise" and "Hypothesis" and then "Answer". This is illustrated in the right column of Figure 2.

Third, we do **10-shot chain-of-thought** prompting, using a prompt involving 10 of SARA's training examples. For each of the 10, we have the "Premise", the "Hypothesis", and "Answer", followed by a human-written chain of thought explaining the relevant reasoning and concluding with the appropriate word "Entailment" or "Contradiction". The same 10-example prompt is used for all test cases. This is illustrated in the center column of Figure 2.

GPT-3 does not always complete its answers with "Entailment" or "Contradiction". In such cases, following Kojima et al. [11], we then pass to GPT-3 our original prompt, plus GPT-3's response, plus an Answer Extraction prompt reading "Therefore, the answer (Entailment or Contradiction) is".

3.3 "Let's think step by step."

Kojima et al. [11] find substantial improvement on a range of tasks by simply adding "Let's think step by step." at the end of the prompt passed to GPT-3. They find this forces GPT-3 to explain its reasoning step by step and thus arrive at better answers. We try each of the zero-shot and 4-shot dynamic prompting experiments both with and without "Let's think step by step." A column in Table 1 indicates whether each row's result used this approach. Note that we do not use this approach for the 10-shot hand-crafted chain-of-thought, since the 10 examples already contain detailed step-by-step reasoning.

3.4 With or Without Statute

We hypothesize that GPT-3 may have seen the U.S. tax code in its training data, since U.S. federal statutes cannot be copyrighted and thus appear on multiple websites. So for each experimental setup, we try all the SARA prompts both with and without the text of the statute(s) required to solve all the hypotheses presented, as specified in the "Include Statute" column in Table 1. In 4-shot dynamic prompting and zero-shot prompting, we use the most relevant section from the statutes. For 10-shot chain-of-thought

Prompt type	Include Statute	"Let's think step by step"	Accuracy in % on cases with:		Aggregate accuracy
			numbers	no numbers	
4-shot dynamic	Y	Y	60 ± 10 (43/72)	61 ± 16 (17/28)	60 ± 8 (60/100)
	Y	N	47 ± 10 (34/72)	50 ± 16 (14/28)	48 ± 8 (48/100)
	N	Y	47 ± 10 (34/72)	50 ± 16 (14/28)	48 ± 8 (48/100)
	N	N	49 ± 10 (35/72)	64 ± 16 (18/28)	53 ± 8 (53/100)
zero-shot	Y	Y	61 ± 10 (44/72)	75 ± 14 (21/28)	65 ± 8 (65/100)
	Y	N	74 ± 9 (53/72)	64 ± 16 (18/28)	71 ± 8 (71/100)
	N	Y	49 ± 10 (35/72)	50 ± 16 (14/28)	49 ± 8 (49/100)
	N	N	38 ± 10 (27/72)	39 ± 16 (11/28)	38 ± 8 (38/100)
10-shot hand-crafted chain-of-thought	Y	N	56 ± 10 (40/72)	61 ± 16 (17/28)	57 ± 8 (57/100)
	N	N	54 ± 10 (39/72)	64 ± 16 (18/28)	57 ± 8 (57/100)
BERT-based state of the art [9]	Y	N	56 ± 10 (40/72)	68 ± 15 (19/28)	59 ± 8 (59/100)
Majority baseline	N	N	50 ± 10 (36/72)	50 ± 16 (14/28)	50 ± 8 (50/100)

Table 1: GPT-3 performance on SARA with various setups. Below the double line are non-GPT-3 approaches. We report both the 90% confidence interval and the raw number of accurate answers.

prompting, we use a set of sections most relevant to the entire test set: 2(b), 63(c)(1) and (6), 152, 3306(a), and 7703.

When we do not include the statutes, we replace all instances of "section" in the SARA premises and hypotheses with "I.R.C. section". (The U.S. tax code is formally known as the Internal Revenue Code, and "I.R.C." is the standard citation.)

3.5 Discussion

Results for the settings described above are reported in Table 1. As expected, GPT-3 is quite sensitive to the prompt setting, with large variations in performance across settings. The prompt "Let's think step by step." occasionally improves statutory reasoning performance, though not systematically. The BERT-based state of the art we report uses some amount of human-provided structure. Some of our prompt settings surpass that previous state of the art, with far less human input, on the full set of 100 entailment cases.

As measured by Welch's unequal variances t -test, our best GPT-3-based model performs significantly better than the BERT-based state of the art ($p \approx 0.038$), and than the majority baseline ($p \approx 0.0011$). With $p < 0.05$, only our best model has greater accuracy than the previous state of the art, and only our top two have greater accuracy than the majority baseline.

We note that on nine of our ten experimental setups, GPT-3 performs at least as well on the 28 cases involving no numbers than on the 72 cases involving numbers. This is expected, since most of the numbers are quite large (e.g. "In 2016, Alice's income was \$310192.") and GPT-3 has a known limitation in doing calculations with large numbers.

As expected, GPT-3 performs at chance (0.50) or worse in the zero-shot tests where there is no statute included. But we were quite surprised to examine GPT-3's step by step reasoning, which shows some (flawed) knowledge of the U.S. tax code. Here is one example, with our full prompt in regular font, and GPT-3's full response highlighted:

Premise: Bob is Alice's brother since April 15th, 2014.
Hypothesis: Alice bears a relationship to Bob under I.R.C. section 152(d)(2)(A).

Answer: Let's think step by step. **First, I.R.C. section 152(d)(2)(A) states that a taxpayer is allowed to claim a dependent if the dependent is a "brother or sister by whole or half blood." Since Bob is Alice's brother, this would qualify him as a dependent under I.R.C. section 152(d)(2)(A). Therefore, the hypothesis is correct.**

Recall that no statutory text was provided in the prompt, so we see that GPT-3 thinks it knows about the U.S. tax code. But its knowledge is incorrect, since section 152(d)(2)(A) actually states that an individual is a qualifying relative if they are "(A) A child or a descendant of a child" of the taxpayer. The closest match is in the following subparagraph, section 152(d)(2)(B), which lists "(B) A brother, sister, stepbrother, or stepsister." In other words, GPT-3 went awry in two ways: first, it mistakenly referred to the wrong subparagraph; and second, it confused half-siblings with step-siblings. As a result of this and similar results, we run experiments, described below in Section 4, to determine how much of the U.S. Code GPT-3 knows.

GPT-3's tendency to mistakenly refer to the wrong part of a statute is not limited to prompts that do not include the text of the statutes. It happens even when we include the statute in the prompt itself. An example is shown in Figure 3. Given the facts in the prompt, GPT-3 correctly reasons that Alice is Bob's stepmother. But then it confuses section 152(b)(2)(D) with (b)(2)(C) and gives the wrong answer. Because of errors like these, we set out in Section 5 to test GPT-3's ability to reason on synthetic statutes that it has never seen before.

3.6 Does GPT-3 know about SARA?

SARA was first made available on the internet in 2020, so it is possible that SARA was in GPT-3's training data. This, in turn,

...

§152. Dependent defined

...

(d) Qualifying relative

...

(2) Relationship For purposes of paragraph (1)(A), an individual bears a relationship to the taxpayer described in this paragraph if the individual is any of the following with respect to the taxpayer:

(A) A child or a descendant of a child.

(B) A brother, sister, stepbrother, or stepsister.

(C) The father or mother, or an ancestor of either.

(D) A stepfather or stepmother.

...

Premise: Charlie is Bob’s father since April 15th, 1995. Dorothy is Bob’s mother. Alice married Charlie on August 8th, 2018.

Hypothesis: Alice bears a relationship to Bob under section 152(d)(2)(C).

Answer: In 2018, Alice is the wife of Charlie. Charlie is Bob’s father. Therefore, Alice is Bob’s stepmother. Therefore, the answer is Entailment.

Figure 3: Example of SARA getting the wrong answer, because it confused section 152(d)(2)(C) with (d)(2)(D). This example came from the 10-shot chain-of-thought prompting with the statutes included. We elided much of the prompt for brevity.

might have biased the results of these experiments. To test this possibility, we first asked GPT-3 "What is the StAtutory Reasoning Assessment (SARA)?" and it gave an answer that it was a multiple choice test for "law students, lawyers, and other legal professionals." (We know of no such test.) We then tried passing 20 randomly selected SARA cases followed by the text "Where does the text above come from?" None of the answers remotely implicated SARA. The modal answer (5 of 20) was that the language was from the IRS website.

4 GPT-3’S AWARENESS OF THE U.S. CODE

As discussed earlier, GPT-3’s responses to zero-shot SARA prompts without any statutory text demonstrate that it has some, possibly imperfect, knowledge of the U.S. tax code. In this section we probe what GPT-3 knows about the tax code and U.S. Code in general. The U.S. Code or U.S.C. is the official compilation of permanent U.S. federal statutes. It is organized into 53 numbered titles, each corresponding to a particular topic. Title 26, also called the Internal Revenue Code, is the U.S. tax code.

We analyze GPT-3’s knowledge of the U.S. Code in two directions. One is to give GPT-3 the full text of a section of the U.S. Code and ask it to identify its origin, including title and section. The other is to ask GPT-3 to recite the text of a particular section and compare its answer to the actual text.

4.1 Given statutory text from the U.S. Code, GPT-3 can often identify the title number

For each of the 53 titles of the U.S. Code, we randomly selected 10 sections whose bodies are between 100 and 1000 words. The lower limit of 100 words was set to exclude statutes with little semantic

content, such as 26 U.S.C. §8003, which in its entirety reads "The Joint Committee shall elect a chairman and vice chairman from among its members." The upper limit of 1000 words was set to avoid running into GPT-3’s hard limit of 4000 total tokens, with each word potentially consisting of multiple tokens.

For each section, we passed a prompt consisting of the text of the statute (**not** including the section’s name, number, or title), followed by a newline and the text "Where is the text above from?". Once we had a response, following Kojima et al. [11], we passed the original prompt and response, plus the text "So is it from the U.S. Code? The answer (Yes or No) is". If the answer was Yes, we passed all the previous prompts and responses concatenated, plus the text "What title of the U.S. Code is it from? The answer (arabic numerals) is". If we got the correct answer, we passed all the previous prompts and responses concatenated, plus the text "What section of title N of the U.S. Code is it from? The answer (arabic numerals) is" where N is the proper title. The results from our 530 tests are in Table 2.

Of the 265 sections for which GPT-3 correctly identified the title but could not correctly identify the section, it identified a numerical section for 250 of the 265. Of these 250, for 36 the section was off by just one (e.g. GPT-3 identified text as from §104 when it was actually from §103). For exactly 100 sections, GPT-3 was off by nine or less (e.g. GPT-3 identified text as from §161 when it was actually from §170).

Even many of the 64 prompts to which GPT-3 responded "No" when asked "So is it from the U.S. Code? The answer (Yes or No) is" are understandable. For example, GPT-3 said that the text of 34 U.S.C. §10707 is not from the U.S. Code. Rather, GPT-3 responded, "The text above is from the Comprehensive Opioid Abuse Grant Program Evaluation Act of 2016." The section actually was added by the Comprehensive Addiction and Recovery Act of 2016. GPT-3 got the title of the Act slightly wrong and entirely missed its subsequent codification into the U.S. Code.

Of particular relevance to understanding GPT-3’s performance on SARA, we separately asked GPT-3 to identify the text of the nine curated statutes in SARA. GPT-3 correctly identified the title and exact section for eight of the nine. (It incorrectly identified 26 U.S.C. §3301 as §3111.) This high performance makes sense, given that the sections used in SARA are central to individuals’ tax calculations and thus more likely to be talked about extensively in the training corpus for GPT-3.

Not from U.S. Code	12% (64)
Wrong title	25% (133)
Right title, wrong section	50% (265)
Title and section correct	13% (68)

Table 2: GPT-3’s answers when given the text of the body (but not the name, section number, or title number) of a section from the U.S. Code.

4.2 Given a U.S. Code citation, GPT-3 can recite plausible but incorrect statutory language

We just saw that GPT-3, when given statutory text from the U.S. Code, is decent at identifying it as such (89% correct) and even identifying its title number (63% correct). We now turn to the flip side: given a citation from the U.S. Code, can GPT-3 recite the text?

For 10 randomly selected statutes in each of the 53 titles of the U.S. Code, we prompted GPT-3 with "The text of __ U.S. Code section __ is:", with the title number in place of the first blank and the section number in place of the second.²

We never saw GPT-3 exactly recite statutory text. It always provides an answer in the same style as the U.S. Code, using Congress' formal, dry style and often organizing into numbered subsections, paragraphs, subparagraphs, etc. Sometimes GPT-3's recitation of a section gets the gist of the section correct, but omits some statutory language and inserts some nonexistent statutory language. But many times GPT-3 recites text that gives the gist of an entirely incorrect section.

We aimed to quantify how closely GPT-3 gets at least the gist of sections. One natural possibility is BLEU, a score used to evaluate machine translations against reference human-generated high-quality translations [15]. But standard BLEU strongly penalizes translations that are shorter than the human-generated reference text, and GPT-3 often recites only the key portions of a statute. So, we use the metric of **unpenalized-BLEU**, which is the standard BLEU score without the brevity penalty. Thus, like BLEU, unpenalized-BLEU ranges from 0 (worst) to 100 (best).

Across the 530 sections we asked GPT-3 to recite, we got a mean unpenalized-BLEU score of 7.11 and median of 4.52. Only 28 of the 530 have an unpenalized-BLEU above 20. For machine translation standards, these are fairly poor results.

Given that GPT-3 often recites roughly correct statutory language for a section different than the one requested, we calculated the unpenalized-BLEU score of the recited text against all the sections in the correct title. For example, for 17 U.S.C. §101 we computed the unpenalized-BLEU score of GPT-3's recitation of the section against all 147 sections contained in title 17. We can then rank the actual U.S.C. sections against GPT-3's prediction. With this, we can measure the rank of the expected section, and thus recall@k. Recall@1 is only 1.5% (8 of the 530 sections sampled) and recall@5 is only 4.7% (25/530).

These results led us to evaluate whether GPT-3's recited text reflected the section number at all. For each of the 530 samples, we normalized the rank of the correct answer within the title by subtracting 1 from the rank and dividing by the number of sections in the title, minus 1. If the correct answer was ranked first for a given section, its normalized recall would thus be 0.0; if GPT-3 is paying no attention whatsoever to the requested section number, we would expect the normalized rank to average 0.5. We observed a mean of 0.48 and median of 0.46, which suggests that GPT-3 pays little attention to the requested section number.

²Recall that throughout our experiments we use temperature=0 and top_P=1.0, which ensures minimal creativity.

5 GPT-3 STRUGGLES WITH SIMPLE REASONING ON SYNTHETIC STATUTES

We know of no technique to determine how much of GPT-3's flawed knowledge of the U.S. Code impacts its performance on SARA, which is based on nine sections drawn from the U.S. Code. To test GPT-3's ability to reason on statutes it has definitely never seen before, we wrote code to create entirely synthetic statutes, using the same numbering style (i.e. subsection, paragraph, etc.) as the U.S. Code. We then systematically prompt GPT-3 with the basic statutory reasoning task of giving it one fact and asking it to determine whether a particular subsection applies.³ These statutory-reasoning prompts are similar to SARA, involving self-contained statutes, facts, then simple questions. Also, we take inspiration from the setting of Saparov and He [17].

5.1 Creation of Synthetic Statutes

An example of one of our synthetic statutes appeared in Figure 1. To avoid any ambiguity, the statutes involve repeated application of the most basic logical form $A \implies B$, such as *parkinse* \implies *rolang*. In other words, the statutes define terms (rather than, say, imposing duties). Many real-world statutes consist of precisely this sort of term defining, such as 26 U.S.C. §152(a)(1) which provides *qualifying child* \implies *dependent*.

The statutes are defined in terms of *depth* and *width*. The width is the number of times, for each term B , a relation of $A \implies B$ is stated. The depth is the maximum number of times this logical form can be applied. For example, the synthetic statutes in Figures 1 and 4 are both 2-wide and 2-deep. Meanwhile, the synthetic statute in Figure 5 is 2-wide and 3-deep. Internally, our code represents these statutes using trees.

Our synthetic statutes are much simpler than many sections of the U.S. Code. They involve no negation or ambiguity, and they are entirely balanced, with all leaf nodes being the same depth from the root, and all non-leaf nodes having the same number of children. In contrast, many sections of the U.S. Code have depth 5, with subsections, paragraphs, subparagraphs, clauses, and subclauses, such as 26 U.S.C. §2.

To ensure we are testing GPT-3 on synthetic statutes it has definitely never seen before, we take two approaches to generating the terms being defined. The first is to use English **nonces** generated by a nonce generator; examples include *roland*, *parkinse*, and *oxideney*. The second is to use random **ids** consisting of one letter plus the same numeral repeated twice; examples include *s88*, *f77*, *m55*, and *a22*.

5.2 Zero-Shot Experimental Setup

For zero-shot prompting involving our synthetic statutes, as shown in Figure 1, the prompt we give to GPT-3 consists of four concatenated strings:

- the synthetic statute
- a single fact in the form " is a ", as with "Alexis is a portle." The name is randomly chosen from a set of 30 names

³In this section, as throughout this paper, all our calls to GPT-3 have temperature at 0.0 and the top_P at 1.0 to maximize reproducibility and minimize hallucination and wandering off topic.

phrasing	example	accuracy (%)
Is S applicable to N ?	Is section 1001(b) applicable to Alexis?	77 (2303/3000)
Does S apply to N ?	Does section 1001(b) apply to Alexis?	74 (2220/3000)
Does S apply to N , making her/him a T ?	Does section 1001(b) apply to Alexis, making her a rolang?	54 (540/1000)
Does S apply to make N a T ?	Does section 1001(b) apply to make Alexis a rolang?	58 (576/1000)
Is N a T because of S ?	Is Alexis a rolang because of section 1001(b)?	57 (565/1000)
Is N a T owing to S ?	Is Alexis a rolang owing to section 1001(b)?	54 (535/1000)
Is N a T as per S ?	Is Alexis a rolang as per section 1001(b)?	52 (519/1000)

Table 3: Results from different phrasing styles of the same question, all measured on nonce 2-wide, 2-deep synthetic statutes. The person is N , and S is the subsection. T is the top-level term being defined in the synthetic statute. Note that in all the synthetic prompts we generate, person N is a T ; the only question is whether that happens through S . Note that questions including T uniformly perform worse than the two without. We ran an extra 2000 experiments to choose between the top two phrasings.

(15 female, 15 male).⁴ The second part of the fact (e.g. portle) is always a term present somewhere in the synthetic statute.

- a question about whether a subsection applies (e.g. "Is section 1001(c) applicable to Alexis?"). This subsection is randomly chosen, although it is never a leaf, meaning it is never a subsection taking up only a single line. We chose the phrasing "Is __ applicable to __?" after trying a variety of different ways of phrasing the question, with the results shown in Table 3.
- the text "Let's think step by step.", where Kojima et al. [11] found that adding this at the end of the prompt maximizes performance on a variety of reasoning tasks.

We prompt GPT-3 with this string to get an initial response. To derive a definitive answer we then call GPT-3 a second time with the prompt being the original prompt, plus the initial response, and the text "Therefore, the answer (Yes or No) is". This follows Kojima et al. [11] and is often necessary to force GPT-3 to give a clear answer. Examples of this prompting can be seen in Figure 4.

To test whether GPT-3 has a problem specifically with reasoning over statutes, our code also generates semantically-identical sentence versions of each synthetic statute. For each defined term in the synthetic statute, our code creates one sentence with the exact same definitional language as the corresponding part of the statute. These sentences are numbered as "Sentence 1", "Sentence 2", etc. for reference in determining whether a particular sentence is applicable. For each zero-shot problem posed to GPT-3, we test the semantically-identical sentence version of the same statute and question. An example appears at the bottom of Figure 4. For these sentence versions, the prompt again consists of four parts:

- the numbered sentences expressing the same definitions as the synthetic statute.
- the same fact as in the statutory version (e.g. "Alexis is a portle.")
- a question about whether a sentence applies (e.g. "Does sentence 3 apply to Alexis?").

- the text "Let's think step by step."

We run all tests using both nonce and ids terms. We also run all tests using both the statute version of the definitions and the semantically-identical sentence-based versions. We balance all tests with equal numbers having positive groundtruth (e.g. section 1001(b) is not applicable) and negative groundtruth (e.g. section 1001(b) is applicable).

The results are in Table 4. Our primary finding is that GPT-3 performs poorly, with performance around 78% for even the simplest statutes, like that in Figure 1, that are 2-wide, 2-deep. Performance declines further as the statutes get wider or deeper.

term	width	depth	accuracy on statutes (%)	accuracy on sentences (%)
nonce	2	2	78 (779/1000)	79 (793/1000)
ids	2	2	78 (778/1000)	83 (834/1000)
nonce	3	2	70 (698/1000)	72 (718/1000)
ids	3	2	68 (682/1000)	69 (692/1000)
nonce	4	2	69 (688/1000)	72 (723/1000)
ids	4	2	66 (663/1000)	69 (691/1000)
nonce	2	3	75 (754/1000)	75 (747/1000)
ids	2	3	77 (771/1000)	66 (662/1000)
nonce	3	3	75 (374/500)	74 (371/500)
ids	3	3	70 (350/500)	66 (331/500)

Table 4: Zero-shot accuracy of GPT-3 on answering whether a statutory section applies to Alice, given a single fact about Alice. The column "term" indicates whether the statutory terms were nonces like "ansgivath" and "propial" or ids like "f55" and "q11". The rightmost column is GPT-3's accuracy answering questions based on sentences that are semantically-identical to the synthetic statute, as seen in Figure 4.

⁴<https://www.ssa.gov/oact/babynames/decades/names2000s.html>

Section 1001. Definition of Infarber.

(a) General rule

The term "infarber" means-

- (1) any parentiable, or
- (2) any digirderasters.

(b) Parentiable

The term "parentiable" means-

- (1) any packle, or
- (2) any amperseced.

(c) Digirderasters

The term "digirderasters" means-

- (1) any artion, or
- (2) any irtityrating.

Emma is an artion. Is section 1001(b) applicable to Emma? Let's think step by step.

Section 1001 defines the term "infarber" as either a "parentiable" or a "digirderasters". Emma is an artion, which is a type of digirderasters. Therefore, section 1001(b) is applicable to Emma.

Therefore, the answer (Yes or No) is Yes.

Sentence 1: The term "infarber" means any parentiable or any digirderasters.

Sentence 2: The term "parentiable" means any packle or any amperseced.

Sentence 3: The term "digirderasters" means any artion or any irtityrating.

Emma is an artion. Is sentence 2 applicable to Emma? Let's think step by step.

Sentence 2 states that the term "parentiable" means any packle or any amperseced. Emma is an artion, which is a type of digirderaster. Therefore, sentence 2 is not applicable to Emma.

Therefore, the answer (Yes or No) is No.

Figure 4: Top: an example of GPT-3 incorrectly answering a zero-shot question about a 2-wide, 2-deep statute. Bottom: GPT-3 correctly answering the semantically-identical question where the synthetic statute was written in sentence form. Text generated by GPT-3 is highlighted. Note that sentence 1 corresponds exactly to section 1001(a), sentence 2 to 1001(b), and sentence 3 to 1001(c).

We notice comparable performance between statutes with terms that are nonces like *portle* and statutes with terms that are ids like *m77*. This makes sense, since both types of synthetic statutes were designed to be entirely novel to GPT-3.

We also observe comparable performance whether the statute is presented in statute version or semantically-identical sentence-based versions. GPT-3 correctly answers 6537 of the 9000 statute questions but 6562 of the 9000 sentence-based versions.

GPT-3's errors are overwhelmingly false positives, meaning it concludes that a section or sentence applies when it actually does not. In the 9000 nonce runs, there were 2272 errors. Of these, 2204

were false positives and 61 false negatives. (There were also 7 where GPT-3 did not give a definitive answer.)

Section 1001. Definition of Bowlerly.

(a) General rule

The term "bowlerly" means-

- (1) any waitormenteed, or
- (2) any kiterrupider.

(b) Waitormenteed

(1) General rule

The term "waitormenteed" means-

- (A) any redeba, or
- (B) any dischieviders.

(2) Redeba

The term "redeba" means-

- (A) any ersubs, or
- (B) any pushotyptopses.

(3) Dischieviders

The term "dischieviders" means-

- (A) any nookede, or
- (B) any chastiles.

(c) Kiterrupider

(1) General rule

The term "kiterrupider" means-

- (A) any bruselers, or
- (B) any fashiple.

(2) Bruselers

The term "bruselers" means-

- (A) any legimetar, or
- (B) any exematess.

(3) Fashiple

The term "fashiple" means-

- (A) any tanded, or
- (B) any goghts.

Hannah is a chastiles. Is section 1001(c)(3) applicable to Hannah? Section 1001(c)(3) says that fashiple means any tanded or any goghts. Hannah is none of these, so section 1001(c)(3) does NOT apply to Hannah.

Alyssa is a goghts. Is section 1001(c)(3) applicable to Alyssa? Section 1001(c)(3) says that fashiple means any tanded or any goghts. Alyssa is a goghts, so section 1001(c)(3) does apply to Alyssa.

Nicholas is a pushotyptopses. Is section 1001(c)(2) applicable to Nicholas? Section 1001(c)(2) says that redeba means any ersubs or any pushotyptopses. Nicholas is a pushotyptopses, so section 1001(c)(2) does apply to Nicholas.

Figure 5: Two-shot prompt with GPT-3 giving an incorrect answer. This is a 2-wide, 3-deep statute. Note that the two examples given to GPT-3, involving Hannah and Alyssa, are correct. GPT-3's answer shows it has incorrectly looked at section 1001(b)(2), not (c)(2).

width	depth	two-shot accuracy (%)	zero-shot accuracy (%)
2	2	100 (1000/1000)	78
3	2	98 (982/1000)	70
4	2	97 (972/1000)	69
2	3	87 (874/1000)	75
3	3	81 (405/500)	75

Table 5: Two-shot accuracy of GPT-3, when the prompt has two example questions with reasoning and correct answers, as in Figure 5. The leftmost column reproduces the comparable zero-shot accuracy from Table 4.

5.3 Two-Shot Experimental Setup

Given GPT-3’s poor performance on the zero-shot statutory reasoning discussed above, we turned to whether GPT-3 performs better with two-shot reasoning. In two-shot reasoning we give two correct examples (one answered Yes, the other answered No), before posing the actual question. All two-shot experiments were run with synthetic statutes where the terms were nonces. (We used neither semantically-identical sentences, nor ids like g11.) A two-shot example involving a 2-wide, 3-deep statute appears in Figure 5.

Both the correctly answered questions are of the form “__ is a __. Does section __ apply to __?”, followed by a correct explanation and answer. The Yes-answered question comes first with 50% probability. Both questions relate to the same section, which never overlaps with the section asked about in the third question. All three questions involve different randomly-selected names (e.g. Hannah, Alyssa, Nicholas) and randomly-selected terms (e.g. chastiles, goghts, pushotyptopses).

We see uniformly better performance with the two-shot model than with the zero-shot model (Table 5). With two examples, GPT-3 seems capable of handling 100% of 2-wide, 2-deep statutes, such as those shown on our first page in Figure 1. As in the zero-shot setting, GPT-3’s accuracy tends to decrease as the statutes’ depth and width increase. GPT-3 still performs quite poorly (81%) on 3-wide, 3-deep statutes, which themselves pale in comparison to the complexity of many sections in the U.S. Code.

CONCLUSION

Being able to identify whether a specific subsection of an unfamiliar statute is applicable to a given set of facts is one of the most basic skills required of a lawyer. We find that, given a very simple 2-wide, 2-deep synthetic statute paired with a single fact and single question (*zero-shot*), GPT-3 had 78% accuracy, raising doubts about GPT-3’s ability to handle basic legal work. Providing two examples of correct statutory reasoning (*two-shot reasoning*) improves GPT-3’s performance, but it still achieves only (81%) on 3-wide, 3-deep synthetic statutes, which are far less complex than what is found in the U.S. Code. This poor performance on synthetic statutory reasoning allows us to understand our results applying GPT-3 to the SARA dataset. These results, at 71% accuracy, are better than the previous state of the art, but leave significant room for improvement.

We hope this work motivates further research into improving the performance of large language models on statutory reasoning. We also look forward to testing new large language models, like GPT-4, on statutory reasoning.

ACKNOWLEDGMENTS

This work has been supported by the U.S. National Science Foundation under grant No. 2204926. Any opinions, findings, and conclusions or recommendations expressed in this article are those of the authors and do not necessarily reflect the views of the National Science Foundation.

REFERENCES

- [1] Sheng Bi, Zhiyao Zhou, Lu Pan, and Guilin Qi. 2022. Judicial knowledge-enhanced magnitude-aware reasoning for numerical legal judgment prediction. *Artificial Intelligence and Law* (2022), 1–34. <https://doi.org/10.1007/s10506-022-09337-4>
- [2] Michael Bommarito, II and Daniel Martin Katz. 2022. GPT Takes the Bar Exam. *CoRR* abs/2212.14402 (2022). <https://doi.org/10.48550/arXiv.2212.14402>
- [3] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language Models are Few-Shot Learners. In *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6–12, 2020, virtual*, Hugo Larochelle, Marc Aurelio Ranzato, Raia Hadsell, Maria-Florina Balcan, and Hsuan-Tien Lin (Eds.). <https://proceedings.neurips.cc/paper/2020/hash/1457c0d6bfc4967418bf8ac142f64a-Abstract.html>
- [4] Jonathan H. Choi, Kristin E. Hickman, Amy Monahan, and Daniel B. Schwarcz. 2023. ChatGPT Goes to Law School. *SSRN* abs/2212.14402 (2023). <https://doi.org/10.2139/ssrn.4335905> SSRN:2212.14402
- [5] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*. Association for Computational Linguistics, Minneapolis, Minnesota, 4171–4186. <https://doi.org/10.18653/v1/N19-1423>
- [6] Bernal Jimenez Gutierrez, Nikolas McNeal, Clayton Washington, You Chen, Lang Li, Huan Sun, and Yu Su. 2022. Thinking about GPT-3 In-Context Learning for Biomedical IE? Think Again. In *Findings of the Association for Computational Linguistics: EMNLP 2022, Abu Dhabi, United Arab Emirates, December 7–11, 2022*, Yoav Goldberg, Zornitsa Kozareva, and Yue Zhang (Eds.). Association for Computational Linguistics, 4497–4512. <https://aclanthology.org/2022.findings-emnlp.329>
- [7] Dan Hendrycks, Collin Burns, Steven Basart, Andy Zou, Mantas Mazeika, Dawn Song, and Jacob Steinhardt. 2021. Measuring Massive Multitask Language Understanding. In *9th International Conference on Learning Representations, ICLR 2021, Virtual Event, Austria, May 3–7, 2021*. OpenReview.net. <https://openreview.net/forum?id=d7KBjmL3GmQ>
- [8] Nils Holzenberger, Andrew Blair-Stanek, and Benjamin Van Durme. 2020. A Dataset for Statutory Reasoning in Tax Law Entailment and Question Answering. In *Proceedings of the Natural Language Processing Workshop 2020 co-located with the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD 2020), Virtual Workshop, August 24, 2020 (CEUR Workshop Proceedings, Vol. 2645)*, Nikolaos Aletras, Ion Androutsopoulos, Leslie Barrett, Adam Meyers, and Daniel Preottiuc-Pietro (Eds.). CEUR-WS.org, 31–38. <http://ceur-ws.org/Vol-2645/paper5.pdf>
- [9] Nils Holzenberger and Benjamin Van Durme. 2021. Factoring Statutory Reasoning as Language Understanding Challenges. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing, ACL/IJCNLP 2021, (Volume 1: Long Papers), Virtual Event, August 1–6, 2021*, Chengqing Zong, Fei Xia, Wenjie Li, and Roberto Navigli (Eds.). Association for Computational Linguistics, 2742–2758. <https://doi.org/10.18653/v1/2021.acl-long.213>
- [10] Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2022. Decomposed Prompting: A Modular Approach for Solving Complex Tasks. *CoRR* abs/2210.02406 (2022). <https://doi.org/10.48550/arXiv.2210.02406>
- [11] Takeshi Kojima, Shixiang Shane Gu, Machel Reid, Yutaka Matsuo, and Yusuke Iwasawa. 2022. Large Language Models are Zero-Shot Reasoners. *CoRR* abs/2205.11916 (2022). <https://doi.org/10.48550/arXiv.2205.11916>

- arXiv:2205.11916
- [12] Jason T. Lam, David Liang, Samuel Dahan, and Farhana H. Zulkernine. 2020. The Gap between Deep Learning and Law: Predicting Employment Notice. In *Proceedings of the Natural Legal Language Processing Workshop 2020 co-located with the 26th ACM SIGKDD International Conference on Knowledge Discovery & Data Mining (KDD 2020), Virtual Workshop, August 24, 2020 (CEUR Workshop Proceedings, Vol. 2645)*, Nikolaos Aletras, Ion Androutsopoulos, Leslie Barrett, Adam Meyers, and Daniel Preotiuc-Pietro (Eds.). CEUR-WS.org, 52–56. <http://ceur-ws.org/Vol-2645/short2.pdf>
- [13] Sarah B Lawsky. 2017. A logic for statutes. *Fla. Tax Rev.* 21 (2017), 60.
- [14] Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Tafjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to Explain: Multimodal Reasoning via Thought Chains for Science Question Answering. *CoRR abs/2209.09513* (2022). <https://doi.org/10.48550/arXiv.2209.09513> arXiv:2209.09513
- [15] Matt Post. 2018. A Call for Clarity in Reporting BLEU Scores. In *Proceedings of the Third Conference on Machine Translation: Research Papers*. Association for Computational Linguistics, Belgium, Brussels, 186–191. <https://www.aclweb.org/anthology/W18-6319>
- [16] Juliano Rabelo, Randy Goebel, Mi-Young Kim, Yoshinobu Kano, Masaharu Yoshikawa, and Ken Satoh. 2022. Overview and Discussion of the Competition on Legal Information Extraction/Entailment (COLIEE) 2021. *Rev. Socionetwork Strateg.* 16, 1 (2022), 111–133. <https://doi.org/10.1007/s12626-022-00105-z>
- [17] Abulhair Saparov and He He. 2022. Language Models Are Greedy Reasoners: A Systematic Formal Analysis of Chain-of-Thought. *CoRR abs/2210.01240* (2022). <https://doi.org/10.48550/arXiv.2210.01240> arXiv:2210.01240
- [18] Jaromír Šavelka and Kevin D Ashley. 2021. Discovering Explanatory Sentences in Legal Case Decisions Using Pre-trained Language Models. In *Findings of the Association for Computational Linguistics: EMNLP 2021*. 4273–4283.
- [19] Jaromír Šavelka and Kevin D Ashley. 2021. Legal information retrieval for understanding statutory terms. *Artificial Intelligence and Law* (2021), 1–45.
- [20] Luke Thorburn and Ariel Kruger. 2022. Optimizing Language Models for Argumentative Reasoning. In *Proceedings of the 1st Workshop on Argumentation & Machine Learning co-located with 9th International Conference on Computational Models of Argument (COMMA 2022), Cardiff, Wales, September 13th, 2022 (CEUR Workshop Proceedings, Vol. 3208)*, Isabelle Kuhlmann, Jack Mumford, and Stefan Sarkadi (Eds.). CEUR-WS.org, 27–44. <http://ceur-ws.org/Vol-3208/paper3.pdf>
- [21] Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. <https://github.com/kingoflolz/mesh-transformer-jax>.
- [22] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Ed H. Chi, Quoc Le, and Denny Zhou. 2022. Chain of Thought Prompting Elicits Reasoning in Large Language Models. *CoRR abs/2201.11903* (2022). arXiv:2201.11903 <https://arxiv.org/abs/2201.11903>
- [23] Fangyi Yu, Lee Quartey, and Frank Schilder. 2022. Legal Prompting: Teaching a Language Model to Think Like a Lawyer. *CoRR abs/2212.01326* (2022). <https://doi.org/10.48550/arXiv.2212.01326> arXiv:2212.01326
- [24] Eric Zelikman, Yuhuai Wu, and Noah D. Goodman. 2022. STaR: Bootstrapping Reasoning With Reasoning. *CoRR abs/2203.14465* (2022). <https://doi.org/10.48550/arXiv.2203.14465> arXiv:2203.14465
- [25] Hanlin Zhang, Ziyang Li, Jiani Huang, Mayur Naik, and Eric Xing. 2022. Improved logical reasoning of language models via differentiable symbolic programming. In *First Workshop on Pre-training: Perspectives, Pitfalls, and Paths Forward at ICML 2022*.
- [26] Haoxi Zhong, Yuzhong Wang, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020. Iteratively Questioning and Answering for Interpretable Legal Judgment Prediction. In *The Thirty-Fourth AAAI Conference on Artificial Intelligence, AAAI 2020, The Thirty-Second Innovative Applications of Artificial Intelligence Conference, IAAI 2020, The Tenth AAAI Symposium on Educational Advances in Artificial Intelligence, EAAI 2020, New York, NY, USA, February 7-12, 2020*. AAAI Press, 1250–1257. <https://ojs.aaai.org/index.php/AAAI/article/view/5479>
- [27] Denny Zhou, Nathanael Schärli, Le Hou, Jason Wei, Nathan Scales, Xuezhi Wang, Dale Schuurmans, Olivier Bousquet, Quoc Le, and Ed H. Chi. 2022. Least-to-Most Prompting Enables Complex Reasoning in Large Language Models. *CoRR abs/2205.10625* (2022). <https://doi.org/10.48550/arXiv.2205.10625> arXiv:2205.10625