

ARTICLE

GENERATIVE MISINTERPRETATION

JAMES GRIMMELMANN,* BENJAMIN L.W. SOBEL,† &
DAVID STEIN‡

ABSTRACT

In a series of provocative experiments, a loose group of scholars, lawyers, and judges has endorsed generative interpretation: asking large language models (LLMs) like ChatGPT and Claude to resolve interpretive issues from actual cases. With varying degrees of confidence, they argue that LLMs are (or will soon be) able to assist—or even replace—judges in performing interpretive tasks like determining the meaning of a term in a contract or statute. A few go even further and argue for using LLMs to decide entire cases and to generate opinions supporting those decisions.

We respectfully dissent. In this Article, we show that LLMs are not yet fit for use in judicial chambers. Generative interpretation, like all empirical methods, must bridge two gaps to be useful and legitimate. The first is a reliability gap: are its methods consistent and reproducible enough to be trusted in high-stakes, real-world settings? Unfortunately, as we show, LLM proponents' experimental results are brittle and frequently arbitrary. The second is an epistemic gap: do these methods measure what they purport to? Here, LLM proponents have pointed to (1) LLMs' training processes on large datasets, (2) empirical measures of LLM outputs, (3) the rhetorical persuasiveness of those outputs, and (4) the assumed predictability of algorithmic methods. We show, however, that all of these justifications rest on unstated and faulty premises about the nature of LLMs and the nature of judging.

The superficial fluency of LLM-generated text conceals fundamental gaps between what these models are currently capable of and what legal interpretation requires to be methodologically and socially legitimate. Put simply, any human or computer can put words on a page, but it takes something more to turn those words into a legitimate act of legal interpretation. LLM proponents do not yet have a plausible story of what that “something more” comprises.

* Tessler Family Professor of Digital and Information Law, Cornell Law School and Cornell Tech. Author order alphabetical.

† Assistant Professor, University of Wisconsin Law School.

‡ Assistant Professor, Vanderbilt Law School. We owe thanks for helpful suggestions and discussions to Yonathan Arbel, Micheal Beauvais, Aislinn Black, Katrina Geddes, David Hoffman, Florencia Marotta-Wurgler, Margaret Lemos, Frank Pasquale, Morgan Ricks, Joel Sobel, Christopher Sprigman, Katherine Jo Strandburg, and Kathryn Woolard. We also thank participants at the NYU Information Law Institute Fellows Workshop, the Internet Law Works-in-Progress Conference, and Vanderbilt Law School. Sobel gratefully acknowledges the support of a grant from the Law & AI Lab at the University of Miami School of Law. All opinions and errors are our own. This Article may be freely reused under the terms of the Creative Commons Attribution 4.0 International license, <https://creativecommons.org/licenses/by/4.0/>.

I.	INTRODUCTION	231
II.	GENERATIVE INTERPRETATION	236
	A. <i>The Proponents' Common Premises</i>	238
	B. <i>LLMs as Interpretive Tools</i>	239
	1. <i>Yonathan Arbel and David Hoffman</i>	239
	2. <i>Judge Kevin Newsom</i>	242
	3. <i>Judge Joshua Deahl</i>	244
	C. <i>LLMs as Adjudicators</i>	245
	1. <i>Adam Unikowsky</i>	245
	2. <i>Legal Tech Entrepreneurs</i>	247
	D. <i>Varied Justifications</i>	248
III.	LLMs GENERATE TEXT BY PREDICTING IT	249
IV.	IS GENERATIVE INTERPRETATION RELIABLE?	252
	A. <i>Famiglio: Model Instability</i>	253
	B. <i>Famiglio: Prompting Instability</i>	255
	C. <i>Snell: Instability through Randomness</i>	258
	D. <i>Katrina Canal Breaches: Implementation Instability</i>	262
	E. <i>Methodological Instability</i>	267
V.	IS GENERATIVE INTERPRETATION EPISTEMICALLY JUSTIFIED?	271
	A. <i>Relying on the Training Process?</i>	272
	1. <i>Deductive Arguments</i>	272
	2. <i>Informal Arguments</i>	274
	B. <i>Relying on Empirical Measurements?</i>	276
	1. <i>Internal Robustness</i>	277
	2. <i>External Calibration</i>	280
	3. <i>Protocol Specification</i>	282
	4. <i>No Easy Fixes</i>	283
	C. <i>Relying on Persuasiveness?</i>	284
	1. <i>Generative Interpretation as Abduction</i>	284
	2. <i>Rhetoric ≠ Adjudication</i>	286
	3. <i>Why Text Matters—And Doesn't</i>	290
	4. <i>The Risk of Projection</i>	294
	D. <i>Relying on Predictability?</i>	295
VI.	REWORKING GENERATIVE INTERPRETATION	297
	A. <i>How AI-Generated Artifacts Can Acquire Normative Weight</i>	297
	1. <i>Rhetorical Value</i>	297
	2. <i>Experimental Validation</i>	299
	3. <i>Adversarial Testing</i>	301
	B. <i>Two Examples</i>	301
	1. <i>Trademark Surveys</i>	302
	2. <i>Corpus Linguistics</i>	303
	3. <i>Implications</i>	305
VII.	CONCLUSION	307

I. INTRODUCTION

“This isn’t right. It’s not even wrong.”¹

“Generative interpretation” is all the rage.² Instead of asking a human judge to interpret a contract, a statute, or other legal text, generative interpretation assigns the task to a large language model (LLM). The idea itself is not new; versions of computer-assisted interpretation have been bouncing around in legal scholarship for over a decade. But the launch of ChatGPT in Fall 2022, followed by rapid improvements and the launch of impressive competitors like Claude, Llama, and Gemini, has led some to argue that the future is now.

In a series of provocative experiments, scholars, judges, and lawyers have fed legal materials into LLMs to generate plausible-seeming answers to real-life interpretive questions. Yonathan Arbel and David Hoffman’s *Generative Interpretation* uses a series of case studies to argue that LLMs can interpret contractual terms well enough that they are ready to assist courts and parties in contract litigation.³ In a pair of concurring opinions, Judge Kevin Newsom has experimented with asking LLMs to interpret key terms from the cases before him.⁴ And in a series of blog posts, appellate lawyer Adam Unikowsky has fed entire case files into Claude, asking it to generate judicial opinions.⁵

These LLM proponents are joined by their admirable curiosity and willingness to experiment with new methods, and by their cautious optimism about generative interpretation. Their bottom lines are broadly similar:

¹ Attributed to Wolfgang Pauli. Michael Shermer, *Wronger Than Wrong*, SCI. AM. (Nov. 1, 2006), <https://www.scientificamerican.com/article/wronger-than-wrong> [<https://perma.cc/YHLS-39NY>].

² See generally Yonathan Arbel & David A. Hoffman, *Generative Interpretation*, 99 N.Y.U. L. REV. 451 (2024) [hereinafter *Generative Interpretation*] (introducing the concept of “generative interpretation”).

³ See *id.* at 458 (“[G]enerative interpretation is good enough for many cases that currently employ more expensive . . . methodologies.”).

⁴ See *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1224–25 (11th Cir. 2024) (Newsom, J., concurring) (describing the use of LLMs to determine the ordinary meaning of “landscaping”); *United States v. Deleon*, 116 F.4th 1260, 1272–75 (11th Cir. 2024) (Newsom, J., concurring) (describing the use of LLMs to determine the ordinary meaning of “physically restrained” and comparing the outputs of different models).

⁵ See Adam Unikowsky, *In AI We Trust*, ADAM’S LEGAL NEWSLETTER (June 8, 2024), <https://adamunikowsky.substack.com/p/in-ai-we-trust> [<https://perma.cc/7NSQ-VDM9>] [hereinafter *In AI We Trust I*]; Adam Unikowsky, *In AI We Trust, Part II*, ADAM’S LEGAL NEWSLETTER (June 16, 2024), <https://adamunikowsky.substack.com/p/in-ai-we-trust-part-ii> [<https://perma.cc/BF5E-Q27X>] [hereinafter *In AI We Trust II*]; Adam Unikowsky, *A Brief History of the Confrontation Clause*, ADAM’S LEGAL NEWSLETTER (June 26, 2024), <https://adamunikowsky.substack.com/p/a-brief-history-of-the-confrontation> [<https://perma.cc/5PEL-7Z4Z>] [hereinafter *A Brief History*]; Adam Unikowsky, *Automating Criminal Appeals*, ADAM’S LEGAL NEWSLETTER (Sep. 18, 2024), <https://adamunikowsky.substack.com/p/automating-criminal-appeals> [<https://perma.cc/JZ3X-KMFA>] [hereinafter *Automating Criminal Appeals*]; Adam Unikowsky, *Ignore the Future*, ADAM’S LEGAL NEWSLETTER (Oct. 21, 2024), <https://adamunikowsky.substack.com/p/ignore-the-future> [<https://perma.cc/GQN4-LF8H>] [hereinafter *Ignore the Future*].

LLMs are already “good enough,”⁶ and judges should seriously consider trusting them to assist with interpretive work in actual cases. They argue that generative interpretation has both quality and quantity advantages over purely human judging. Quality, because LLMs make inferences based on huge amounts of data on contemporary language use, and because automated processes can be objective in a way that fallible and biased humans cannot. And quantity, because LLMs can be deployed cheaply and at scale.

We respectfully dissent. Generative interpretation in its current form is Potemkin interpretation: an attractive facade with nothing behind it. The superficial fluency of LLM-generated text is deeply misleading. LLMs are remarkable text-producing machines, but the way they produce that text falls far short of what the standard account of adjudication demands.

Put simply, any human or computer can put words on a page, but it takes something more to turn those words into a legitimate act of legal interpretation. LLM proponents have no convincing story of what that something more entails. Currently, society treats judges’ opinions as legitimate (most of the time) because they represent a principled working forward from source materials to a persuasive legal conclusion.⁷ LLM proponents treat the persuasiveness of the output as proof that the process was principled. But the lesson of their experiments is precisely the opposite: they have shown that it is possible to produce persuasive, law-like rhetoric without going through the hard work of legal reasoning.

While LLM proponents are right to point out LLMs’ potential and to experiment with them, they are wrong to the extent that they claim that LLMs are currently fit for purpose at legal interpretation. LLMs are already good at some kinds of tasks—such as writing JavaScript code—where the output text is valued for what it does.⁸ They may even be good at some kinds of legal writing: under the right circumstances and with proper adult supervision, it is possible to imagine them helping draft contracts or briefs reliably and effectively.⁹ But legal interpretation is different because adjudication is different;

⁶ For varying definitions of “good enough,” see *infra* Part II.D.

⁷ See generally Neil Duxbury, *Faith in Reason: The Process Tradition in American Jurisprudence*, 15 CARDOZO L. REV. 601 (1993) (discussing Legal Process School’s emphasis on reason-giving).

⁸ See, e.g., Simon Willison, *Here’s How I Use LLMs to Help Me Write Code*, SIMON WILLISON’S WEBLOG (Mar. 11, 2025), <https://simonwillison.net/2025/Mar/11/using-llms-for-code> [<https://perma.cc/9N34-ZVMB>] (providing a tutorial on how to use LLMs to write code).

⁹ Compare Daniel Martin Katz, Michael James Bommarito, Shang Gao, & Pablo Arredondo, *GPT-4 Passes the Bar Exam*, 382 PHIL. TRANSACTIONS ROYAL SOC’Y, Feb. 26, 2024, at 5–8 (arguing that a leading LLM is capable of passing the bar exam), and Andrew Blair-Stanek, Donald G. Gifford, Mark Graber, Guha Krishnamurthi, Jeff Sovern, Donald B. Tobin & Michael Van P. Alstine, *AI Gets Its First Law School A+s* (forthcoming) (manuscript at 4), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5274547 [<https://perma.cc/3JL4-Z9TS>] (“AI models can now perform at an A+ level on some law school final exams.”), with Andrew Blair-Stanek, Nils Holzenberger, & Benjamin Van Durme, *BLT: Can Large Language Models Handle Basic Legal Text?*, in PROCEEDINGS OF THE NATURAL LEGAL LANGUAGE PROCESSING WORKSHOP 2024 216, 216 (“[W]e find that leading publicly available LLMs perform poorly out-of-the-box.”).

it plays a different social role and must satisfy different constraints. Human interpretation varies—somewhat more or less reliable, somewhat more or less persuasive, somewhat more or less affordable—but is fundamentally the same kind of thing. LLM interpretation is something different in kind. It is a simulacrum of interpretation, one that bears no necessary relationship to the realities of linguistic meaning.¹⁰ Perhaps this chasm can be bridged. The first, necessary step in doing so is to admit that the chasm exists.

In particular, generative interpretation must overcome two challenges to be useful and legitimate when used in the judicial process. First, it faces a *reliability* gap: LLM proponents have not yet established that their methods reliably measure anything. The examples they have proffered are intriguing and worthy of further study, but they are very far from being sufficiently replicable to be regarded as reliable.¹¹ A consumer survey with a sample size of $n=1$ and an interviewer who improvised the questions on the spot would be unreliable¹²—no matter how plausible or how dramatic the results were—because repeating the survey would lead to very different results. In a similar way, generative interpretation is unreliable if small tweaks to the prompt or the system flip the outputs. The generative interpretation experiments reported in the literature purport to be “informed conclusion[s] based on a statistical analysis of billions of texts.”¹³ Our examination of their own examples shows, however, that the results might be better described as unwarranted extrapolations from extraneous factors.

Second, generative interpretation faces an *epistemic* gap: LLM proponents have not yet established that their methods measure what they purport to. Even if those methods were robust and replicable, it would all be for naught if they were measuring the wrong thing. A study to determine the surface temperature of the sun by surveying breakfast-cereal consumers would be epistemically invalid, even if it had a rigorously specified survey protocol, interviewed an immense sample size, and passed every statistical test for significance—simply because *what people think the sun’s surface temperature is* is not a good measure of *what the sun’s surface temperature actually is*. In a similar way, generative interpretation is epistemically invalid if its text predictions do not correspond to the linguistic judgments of the relevant community of speakers and listeners. This gap is most obvious for generative adjudication, in which LLMs are tasked with generating entire opinions, but

¹⁰ See generally JEAN BAUDRILLARD, *SIMULACRA AND SIMULATION* (Sheila Faria Glaser trans., U. Michigan Press 1994) (1981) (offering classic account of simulacra).

¹¹ Cf. FED. R. EVID. 702(c) (requiring that expert testimony be “the product of reliable principles and methods”); *Daubert v. Merrell Dow Pharms., Inc.*, 509 U.S. 579, 589 (1993) (confirming that the Federal Rules of Evidence require that “all scientific testimony or evidence admitted [be] . . . reliable.”).

¹² Cf. Gail M. Sullivan & Anthony R. Artino, *How to Create a Bad Survey Instrument*, 9 J. GRAD. MED. EDUC. 411, 412 (2017) (listing common issues in survey design).

¹³ *Generative Interpretation*, *supra* note 2, at 458.

we will show that it is present for the seemingly more constrained task of generative interpretation as well.

The reliability gap and the epistemic gap are closely related. Indeed, we think that the beginning of wisdom is to recognize that both gaps must be bridged for generative interpretation to be useful. Some of the most thoughtful attempts to deal with the reliability gap raise the biggest epistemic problems, and vice versa. Arbel and Hoffman's precise measurements of the inner workings of an LLM, for example, help make generative methods more reliable—but they beg the question of why those measurements are a good epistemic measure of meaning.¹⁴ Conversely, Unikowsky's *res ipsa loquitur* willingness to accept LLM outputs as authoritative because they read persuasively gives a (debatable) response to the epistemic question, but it does nothing to answer the reliability question.¹⁵ Taking both of these gaps seriously *at the same time* is essential, and it is here that we think the LLM proponents have the most work ahead of them.

Our arguments dovetail with recent work by Brandon Waldon, Nathan Schneider, Ethan Wilcox, Amir Zeldes, and Kevin Tobia. In a forthcoming paper, they persuasively show that arguments for LLMs as authoritative “‘super-judges’ of ordinary meaning rest[] upon fundamental myths about how these tools work.”¹⁶ They focus on debunking the factual assumptions about LLMs that generative-interpretation proponents have invoked to justify their purported expertise on ordinary meaning.¹⁷ Our complementary critique addresses distinct empirical and jurisprudential problems with generative-interpretation proposals. On the empirical side, we offer detailed, under the hood demonstrations that generative interpretation produces arbitrary and unpredictable results in a manner that belies its proponents' optimistic claims. On the jurisprudential side, we ask what it would take for LLMs to possess the legal-interpretive authority that generative-interpretation proponents claim they do and show why, as a jurisprudential matter, the proponents' accounts fall short of justifying such authority.

This Article proceeds in five Parts. **Part I** surveys the recent wave of writings from LLM proponents. We divide them roughly into two camps: those who are interested in *generative interpretation*, using LLMs to ascertain the meaning of words and phrases, and those who are interested in *generative adjudication*, using LLMs to resolve entire disputes.

¹⁴ See, e.g., *id.* at 484–85 (examining probability of outputs from LLMs to help gauge the use of “English in its most public and common sense.”).

¹⁵ See, e.g., *In AI We Trust I*, *supra* note 5.

¹⁶ Brandon Waldon, Nathan Schneider, Ethan Wilcox, Amir Zeldes & Kevin Tobia, *Large Language Models for Legal Interpretation? Don't Take Their Word for It*, 114 GEO. L.J. (forthcoming) (manuscript at 22), <https://papers.ssrn.com/abstract=5123124> [<https://perma.cc/2JVM-CMHU>].

¹⁷ See *id.* at 22–39.

Part II provides a brief overview of LLM technology. We do not rehash the history of LLMs or describe the technology in detail. Instead, we highlight a few points that will be relevant to our discussion.

Part III considers the reliability gap. Here, our methodological critique shows, with representative examples, that LLM proponents' results are brittle and often arbitrary. Small changes to an LLM, its settings, or its input can dramatically change the nature of its outputs, often in ways that would alter the outcome of cases. These are not small quibbles. Legal interpretation is a large and complex task, and the space of possible experimental setups is correspondingly large and unconstrained. The fact that LLMs are so easy to use—at the simplest, just type a question in a box—makes it easy to overlook the assumptions and choices involved in using them.

Part IV then considers the epistemic gap. We taxonomize four types of arguments that LLM proponents have made to establish the legitimacy of their use in adjudication. Respectively, they are based on (1) an LLM's *training process* and the immense bodies of text it is trained on, (2) *empirical validation* of the LLM's outputs in comparison with other methods, (3) the *rhetorical persuasiveness* of the text the LLM produces, or (4) the purported *predictability* of LLM-generated outputs. In each case, we point out the unstated and mistaken assumptions on which these arguments rest. In particular, we argue that LLMs highlight the distinction between the *process* that makes adjudication legitimate—the hard work of legal reasoning that judges engage in—and the *artifacts* the process generates—like definitions of disputed terms and opinions justifying an outcome. LLMs sever the connection between the two, because they are capable of producing plausible-looking opinions through a process that is fundamentally alien to traditional legal reasoning. Indeed, they force us to consider the prospect that the superficial fluency of legal text is no longer a meaningful guarantee of anything substantive.

Part V considers what it might take to make LLMs practice-ready for judicial chambers. Even if generative interpretation and adjudication are different and harder, we are more optimistic about the use of LLMs in other aspects of legal practice, and for certain uses in judicial practice. Some scholars have argued that judging is an intrinsically human activity that cannot be outsourced, in whole or part, to any algorithmic system.¹⁸ Our claim is

¹⁸ See, e.g., Amin Ebrahimi Afrouzi, *John Robots, Thurgood Martian, and the Syntax Monster: A New Argument Against AI Judges*, 37 CAN. J. L. & JURIS. 369, 370 (2024) (arguing that “even if AI reaches the same outcome as human judges,” it will be “deficient in rationale”); Kiel Brennan-Marquez & Stephen E. Henderson, *Artificial Intelligence and Role-Reversible Judgment*, 109 J. CRIM. L. & CRIMINOLOGY 137, 140 (2019) (“In some contexts, those who exercise judgment should be vulnerable, in reverse, to its processes and effects.”); Ian Kerr & Carissima Mathen, *Chief Justice John Roberts is a Robot* (U. of Ottawa Working Paper, 2019) (manuscript at 39–40), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3395885 [<https://perma.cc/FL9R-HPBX>] (arguing that AI cannot be a judge because “legal reasoning cannot be reduced to mere functional capabilities”). But see, e.g., Eugene Volokh, *Chief Justice Robots*, 68 DUKE L.J. 1135, 1158 (2019) (arguing for human-algorithmic equivalence in judging). An even stronger version of the claim is that certain decisions must be made by *specific people* and

narrower. We take no position on whether adjudication requires the exercise of specifically moral faculties or whether computer systems, including LLMs, are capable of having and exercising those faculties. Instead, our view is that adjudication is a fundamentally social process, one that must meet certain criteria of rationality to serve its function of legitimating the exercise of authority. Thus, in this Part we describe the work required to make generative interpretation meet the necessary criteria of rationality. In particular, we compare it to two other empirical interpretive methods—trademark surveys and legal corpus linguistics—one of which has successfully established itself in adjudication and the other of which is attempting to do so.

Finally, in a brief **Conclusion**, we consider the broader implications of our analysis. Perhaps the problem facing generative interpretation is not that LLMs aren't good enough at language to be useful to judges, but that they are *too good*. Proponents of generative interpretation cite the persuasiveness of LLM-generated text as proof of that text's legitimacy. But the lesson of their experiments is precisely the opposite. They show that it is possible to produce persuasive, law-like rhetoric to support a conclusion without going through the hard work of legal reasoning to validate that conclusion. Their work shows that we should be more suspicious of words on a page in the age of AI, not less.

II. GENERATIVE INTERPRETATION

LLM proponents have made two different kinds of proposals for incorporating generative AI into adjudication. Some of them, more modestly, have suggested using generative AI as an *interpretive tool* that provides targeted answers to discrete questions about the meaning of a legal text. Others, more ambitiously, have suggested using generative AI as an *adjudicator* that answers broader questions about the proper resolution of a case. They are united by a belief that LLMs can perform some tasks traditionally performed by judges well enough that the legal system should think seriously about incorporating them into actual judicial practice. They differ in how much of the judicial task they believe should be delegated to LLMs. And even within these two broad camps, approaches towards the proper use of LLMs are highly diverse.

These are not completely new ideas. For years, commentators have debated whether decision-making by AIs could be an acceptable substitute for decision-making by human judges.¹⁹ These conversations span many legal

that any delegation of the reasoning or justification even to *other people* is problematic. See Bridget C.E. Dooling, *Ghostwriting the Government*, 109 MARQ. L. REV. (forthcoming 2026) (manuscript at 29), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5200672 [<https://perma.cc/PG5G-G8HC>] (arguing that a “duty to reason” restricts delegation both to humans and to generative AI).

¹⁹ See, e.g., Afrouzi, *supra* note 18; Brennan-Marquez & Henderson, *supra* note 18; Volokh, *supra* note 18.

fields, and encompass not just judicial adjudication but also the immense volume of administrative decisions. Common themes include accuracy,²⁰ consistency,²¹ and explainability.²² In the first instance, these debates have primarily been about the quality of the *decisions* themselves, rather than about the quality of the *explanations* offered to justify those decisions. In other words, these have largely been debates about *classification* or *predictive* AI, rather than about *generative* AI.

Generative AI brings two new things to the table. First, it turns hypotheticals about AIs capable of writing opinions that read like they were written by a human into actualities. That means the question of whether the judicial system should actually use those AIs is far more pressing than it was even a few years ago.²³ When proponents like Adam Unikowsky claim that LLMs are good enough for judges to use *now*, they are speaking in the present tense and the indicative mood.²⁴

Second, generative AI opens up the possibility of generative interpretation: giving an LLM an interpretive task that is significantly smaller than deciding an entire case. This is not actually a task where the generative nature of the AI is essential to its use for interpretation. One could imagine training a classification AI to disambiguate between cases where the last-antecedent canon or the series-qualifier canon is more appropriate.²⁵ Such an AI might have an architecture and training process completely different from the deep-learning methods used to train LLMs; it might be prompted with a passage of text and output a single numerical value expressing the degree to which the passage resembled other passages in which one canon rather than the other applied. *Generative* AI, on the other hand, has seemed like it might work for interpretation because it has already worked for so many other tasks. That is, LLMs are such powerful and versatile tools that they can be (and have been) used, off the shelf or with only slight modifications, for a variety of classification tasks.²⁶ If they are truly jacks of all trades, perhaps legal interpretation is one of those trades.

²⁰ See, e.g., Danielle Keats Citron, *Technological Due Process*, 85 WASH. U. L. REV. 1249, 1301–08 (2008) (providing framework to evaluate technological accuracy in law).

²¹ See, e.g., Cary Coglianese & Lavi Ben Dor, *AI in Adjudication and Administration*, 86 BROOK. L. REV. 791, 828 (2021) (“[M]achine-learning tools . . . could potentially reduce inconsistencies and other foibles that permeate human judgment.”).

²² See, e.g., Margot Kaminski, *The Right to Explanation, Explained*, 34 BERKELEY TECH. L.J. 189 (2019); Andrew D. Selbst & Solon Barocas, *The Intuitive Appeal of Explainable Machines*, 87 FORDHAM L. REV. 1085 (2018).

²³ It also changes the tenor of discussions of the effects these AIs will have. Actual experience with LLMs makes the conversation significantly more concrete. See, e.g., Richard M. Re, *Artificial Authorship and Judicial Opinions*, 92 GEO. WASH. L. REV. 1558, 1582–85 (2024) (offering predictions about effects of widespread use of AI by courts and litigants).

²⁴ See *In AI We Trust II*, *supra* note 5.

²⁵ See generally Adam Crews, *The So-Called Series-Qualifier Canon*, 116 NW. U. L. REV. ONLINE 198, 208–09 (2021) (comparing the two canons).

²⁶ See, e.g., Rebecca M. M. Hicke & David Mimno, *T5 Meets Tybalt: Author Attribution in Early Modern English Drama Using Large Language Models*, 2023 COMPUTATIONAL HUMANS. RSCH. CONF. 1 (using LLM to identify authorship of text passage).

A. *The Proponents' Common Premises*

The justifications that LLM proponents offer for generative interpretation and adjudication are broadly consistent. The remainder of this Part will review arguments from some of the early proponents of LLM-powered interpretation and adjudication—from cautious experimenters like Judges Kevin Newsom and Joshua Deahl, to measured enthusiasts like professors Yonathan Arbel and David Hoffman, to unbridled optimists like practitioners Adam Unikowsky and Jack Kieffaber. These proponents differ along many dimensions, but all see LLMs as potentially more accurate, more consistent, and more efficient than current interpretive methods. The intuition behind these arguments is consistent across the board.

First, the proponents say that LLMs can be more *accurate* than human judges: the models offer higher-quality insights into the ordinary meaning of language in legal texts. They argue that generative interpretation offers a more “majoritarian reading” of text than a dictionary or a jurist’s intuitions.²⁷ The core tenet of this argument is that because the models train on what are perhaps the largest corpora of natural-language text ever assembled, and because they use statistical methods to optimize their predictions of successive words in context, their output offers probative insights about general English usage.²⁸ A related idea is that they are more democratic because their training data reflects a far broader sample than other authoritative sources, like dictionaries.²⁹

Second, the proponents argue that generative interpretation is more *consistent*: more constrained, more predictable, and more transparent.³⁰ In contrast to current interpretive methods’ “value-laden” choices about which dictionaries to consult³¹ or which interpretive canons to employ, generative interpretation offers a “cutting-edge ‘mathematization’ of language,”³² that is

²⁷ *Generative Interpretation*, *supra* note 2, at 488.

²⁸ See *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1226 (11th Cir. 2024) (Newsom, J., concurring) (“[T]he best reason to think that LLMs might provide useful information to those engaged in the interpretive enterprise” is that “the models train on a mind-bogglingly enormous amount of raw data taken from the internet” and therefore “can provide useful statistical predictions about how, in the main, ordinary people ordinarily use words and phrases in ordinary life.”); see also *Generative Interpretation*, *supra* note 2, at 513–14 (praising “superior sensitivity to actual usage,” which, the authors argue, derives from the fact that the models are trained on immense amounts of real-life examples of linguistic usage); *In AI We Trust I*, *supra* note 5 (“[W]e should come at this problem in the spirit of ‘AI has read every case ever written and in most cases, it will be more accurate than humans’ . . .”).

²⁹ See *supra* note 14 and accompanying text.

³⁰ See *Snell*, 102 F.4th at 1227–30 (Newsom, J., concurring) (“Using LLMs to facilitate ordinary meaning interpretation may actually enhance the transparency and reliability of the interpretive enterprise itself . . .”); see also *In AI We Trust I*, *supra* note 5 (“Dispersing the judicial power among so many different judges inevitably undermines predictability. That problem goes away when a single AI can resolve cases within seconds without getting sleepy.”). See generally Kimo Gandall, Jack Kieffaber & Kenny McLaren, *We Built Judge.ai. And You Should Buy It* (Working Paper), <https://papers.ssrn.com/abstract=5115184> [<https://perma.cc/D5JW-8Y7K>].

³¹ *Generative Interpretation*, *supra* note 2, at 466.

³² *Snell*, 102 F.4th at 1227 (Newsom, J., concurring).

more “transparent” and “objective.”³³ In *Snell v. United Specialty Insurance Company*,³⁴ Judge Newsom suggested that “LLM research is relatively transparent.”³⁵ “[W]e tend to take dictionaries for granted, as if delivered by a prophet,” he observed, but “the precise details of [dictionaries’] construction aren’t always self-evident.”³⁶ Moreover, judges can shop around for dictionary definitions that suit their predispositions and “rarely explain in any detail the process by which they selected one definition over others.”³⁷ By contrast, Judge Newsom argued, LLM research is comparatively transparent because we “know,” on some general level, what data LLMs “learn[]” from and because judges could disclose their full queries and models’ full answers.³⁸ Plus, LLMs are less vulnerable to hidden exercises of discretion than tools like corpus linguistics (and perhaps even dictionaries).³⁹

And third, the proponents tell us that generative interpretation is more *efficient*, so its use can improve access to justice. LLMs are widely accessible and far cheaper than actual surveys of humans’ perceptions of meaning.⁴⁰ Judicial use of a predictable, widely accessible model might let parties predict judicial outcomes to resolve disputes while avoiding court altogether.⁴¹

B. LLMs as Interpretive Tools

We start with the commentators who have proposed using LLMs for interpretation. Overall, their proposals are more concrete and specific, and thus easier to assess. Part III’s more technical critique revisits these examples in greater detail.

1. Yonathan Arbel and David Hoffman

The leading exposition of using LLMs for legal interpretation is Yonathan Arbel and David Hoffman’s widely praised⁴² *Generative Interpretation*.⁴³ The article posits that “AI models can help factfinders ascertain ordinary meaning

³³ *Generative Interpretation*, *supra* note 2, at 455, 511.

³⁴ 102 F.4th 1208 (11th Cir. 2024).

³⁵ *Id.* at 1228 (Newsom, J., concurring).

³⁶ *Id.*

³⁷ *Id.* at 1229.

³⁸ *See id.*

³⁹ *See id.* at 1230. For a longer discussion of corpus linguistics, *see infra* note 286 and accompanying text.

⁴⁰ *See Snell*, 102 F.4th at 1228–30.

⁴¹ *See Generative Interpretation*, *supra* note 2, at 510.

⁴² *See, e.g.,* Lawrence Solum, *Hoffman & Arbel on Interpretation of Contracts by Generative AI*, LEGAL THEORY BLOG (Aug. 1, 2023) <https://lsolum.typepad.com/legaltheory/2023/08/hoffman-arbel-on-interpretation-of-contracts-by-generative-ai.html> [<https://perma.cc/U87X-FFZK>] (“A profoundly important article.”); Omri Ben-Shahar, *Towards the End of Normative Interpretation of Contracts*, JOTWELL (Nov. 2, 2023), <https://contracts.jotwell.com/towards-the-end-of-normative-interpretation-of-contracts/> [<https://perma.cc/4MS3-P4EE>] (reviewing *Generative Interpretation* positively).

⁴³ *Generative Interpretation*, *supra* note 2.

in context, quantify ambiguity, and fill gaps in parties' agreements."⁴⁴ The authors' essential argument is that large language models can help interpret contested provisions of legal instruments.

They demonstrate their approach by working through several case studies of litigated contractual disputes. In these case studies, the authors solicit responses from LLMs about the meaning of legal text, and they present these responses as "relevant" evidence of the language's "public and common" meaning.⁴⁵ "[G]enerative interpretation," they conclude, "promises an accessible, relatively predictable, tool that will help lawyers and judges interpret contracts."⁴⁶

For example, they consider a dispute about whether language in a pre-nuptial agreement referring to "the time a Petition for Dissolution of Marriage is filed" meant the time the *first* such petition was filed or the time that the *most recent* petition was filed.⁴⁷ They gave OpenAI's Davinci-003 model—a variant of the LLM used to power ChatGPT at the time—the text of the agreement and asked it, "If one of the parties files a divorce petition, withdraws it, and then a few years later a new petition is filed, what date determines the number of full years of marriage: the first filing or the second one?"⁴⁸ They then observe that the model has a 94.72% probability of outputting "second" in a sentence answering the question and only a 0.68% probability of outputting "first."⁴⁹

Arbel and Hoffman acknowledge that their article is meant as an introduction to a promising interpretive technique, rather than as an exhaustive explication of an interpretive methodology. They do not claim that LLMs surface the absolute truth of contracting parties' intentions: instead, they present generative interpretation as a "workable, workmanlike method for a resource-constrained litigation world."⁵⁰ Nor do they claim that LLMs are ready to serve as "robot judges" vested with the final say in adjudication.⁵¹ Rather, they explain that their proposal is for judges to use LLMs "as tools" to supplement their judgment, and they disclaim that techniques for querying these tools "still await a process of development, refinement, and validation."⁵²

⁴⁴ *Id.* at 451.

⁴⁵ *See id.* at 485 ("Generative interpretation . . . offers courts a better sense of the relevant probabilities if the parties were intending to use English in its most public and common sense."); *see also id.* at 489 ("Generative interpretation . . . helps us visualize a broad spectrum of meaning and quantify how likely a particular result is."); *id.* at 505 ("Models offer an approximation of general understanding that may simply not be available in any other way, and thus advance long-held goals of contract theory.").

⁴⁶ *Id.* at 509.

⁴⁷ *See id.* at 483–84.

⁴⁸ *Generative Interpretation*, *supra* note 2, at 484.

⁴⁹ *See id.*

⁵⁰ *Id.* at 458, 460–61.

⁵¹ *Id.* at 461.

⁵² *Id.* at 455, 461.

Alongside these disclaimers, however, Arbel and Hoffman consistently characterize LLM outputs as accurate, probative evidence of words' legal meanings.⁵³ They further advertise that "generative interpretation is good enough for many cases that currently employ more expensive, and arguably less certain, methodologies."⁵⁴ Whether or not the authors intended to present generative interpretation as a method ready for primetime use,⁵⁵ it now occupies a primetime slot, thanks in part to its influence on prominent LLM proponents like Judge Newsom, whose views we discuss in the following section.

Overall, Arbel and Hoffman's proposal is *narrow* in the questions it asks LLMs, *broad* in the evidence it asks them to consider, and *narrow* in the outputs it examines. They invite judges to turn to generative AI to help with discrete interpretive tasks—*e.g.*, does the term "other affiliates" as used in a contract refer only to entities already in existence?—and they are capacious in the text they provide LLMs to help answer these questions—frequently, entire contracts.⁵⁶ Indeed, they present the use of LLMs as a new form of contextualism: a way for contract interpretation to consider a wide array of evidence bearing on a contract's meaning while preserving the values of predictability, restraint, and judicial economy typically advanced for its alternative, textualism.⁵⁷ They are enthusiastic about giving LLMs access to the complete text of affidavits and other documents beyond the four corners of a contract.⁵⁸ Overall, however, they use LLMs as meaning meters rather than legal writers; rather than relying on the apparent meanings of a chatbot's generated responses, Arbel and Hoffman look closely at the underlying AI's internal calculations.

⁵³ See *supra* note 43; see also *Generative Interpretation*, *supra* note 2, at 495 ("[W]e have provided examples that showcase how large language models might power a stronger, cheaper, more robust form of textualism."); *id.* at 505 ("[O]n average, these models predict with great accuracy linguistic distinctions that humans make."); *id.* at 509 ("[G]enerative interpretation promises an accessible, relatively predictable, tool that will help lawyers and judges interpret contracts"). Compare *id.* at 485 ("Generative interpretation . . . offers courts a better sense of the relevant probabilities if the parties were intending to use English in its most public and common sense."), with *id.* at 485 n.165 ("[T]he probabilities shouldn't be interpreted literally."). We think the best way to interpret these claims and caveats together is as an endorsement of the models' *accuracy*, but a disavowal of the models' claimed *precision*. See *id.* at 496 n.210 ("[Y]ou should be skeptical of model's [sic] expressed confidence; the direction of change with every new piece of evidence, not its quantification, is informative."). For an explanation of the distinction between accuracy and precision, see, *e.g.*, *Accuracy and Precision*, WIKIPEDIA (last visited Mar. 17, 2025), https://en.wikipedia.org/w/index.php?title=Accuracy_and_precision&oldid=1281017807 [<https://perma.cc/79CF-BTUP>].

⁵⁴ *Generative Interpretation*, *supra* note 2, at 458.

⁵⁵ After the authors assert in the present tense that "generative interpretation is good enough," *id.*, they later present the question as, "whether [generative interpretation] is good enough, *if not today then soon*, for resource-deprived courts to adopt in ordinary cases," *id.* at 502 (emphasis added and original emphasis omitted).

⁵⁶ See *id.* at 487, 489–90.

⁵⁷ See *id.* at 510–13.

⁵⁸ See *id.* at 513.

2. Judge Kevin Newsom

The next major proponent of generative interpretation is Judge Kevin Newsom of the United States Court of Appeals for the Eleventh Circuit. In 2024, he wrote his concurrence in *Snell*, in which he characterized himself as “th[inking] the unthinkable” and “sa[ying] the unsayable.”⁵⁹ “LLM[s] like ChatGPT,” he mused, “might have something useful to say about the common, everyday meaning of the words and phrases used in legal texts.”⁶⁰

The relevant issue in *Snell* was whether an insurance policy that covered liability arising out of the insured’s “perform[ance of] landscaping”⁶¹ applied to the allegedly negligent installation of a trampoline.⁶² While the majority disposed of the case without ruling on the meaning of “landscaping,”⁶³ Judge Newsom argued that “[t]hose, like me, who believe that ‘ordinary meaning’ is the foundational rule for the evaluation of legal texts should consider—consider—whether and how AI-powered large language models like OpenAI’s ChatGPT, Google’s Gemini, and Anthropic’s Claude might—*might*—inform the interpretive analysis.”⁶⁴

On “a lark,” Judge Newsom directed one of his clerks to query ChatGPT, “What is the ordinary meaning of ‘landscaping’?”⁶⁵ ChatGPT produced a two-sentence definition that “aligned with [his] priors.”⁶⁶ He then asked, “Is installing an in-ground trampoline ‘landscaping’?”⁶⁷ Both ChatGPT and Google Bard answered in the affirmative.⁶⁸ The coherence of ChatGPT’s response impressed him, and he found that it “squared with [his] own impression” about how ordinary speakers might use the term.⁶⁹

A few months later, Judge Newsom reaffirmed his enthusiasm for LLMs in another concurrence, this time in a case, *United States v. Deleon*,⁷⁰ interpreting a provision of the United States Sentencing Guidelines.⁷¹ The Guidelines impose a sentencing enhancement for an armed robbery “if any person was physically restrained to facilitate commission of the offense or to facilitate escape.”⁷² The defendant had “walked into a store, pointed a gun at the cashier while demanding money from the register, received the money,

⁵⁹ *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1221 (11th Cir. 2024) (Newsom, J., concurring).

⁶⁰ *Id.* at 1234.

⁶¹ *Id.* at 1213 (majority opinion).

⁶² *Id.* at 1211–12.

⁶³ *See id.* at 1216–17.

⁶⁴ *Id.* at 1221–22 (Newsom, J., concurring).

⁶⁵ *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1224–25 (11th Cir. 2024) (Newsom, J., concurring).

⁶⁶ *Id.* at 1225.

⁶⁷ *Id.*

⁶⁸ *Id.*

⁶⁹ *Id.*

⁷⁰ 116 F.4th 1260 (11th Cir. 2024).

⁷¹ *See id.* at 1261.

⁷² *Id.* at 1261 (quoting U.S. SENT’G GUIDELINES MANUAL § 2B3.1(b)(4)(B)).

and . . . left . . . within about one minute.”⁷³ The court held the enhancement applicable, and Judge Newsom wrote separately to suggest that LLMs might “inform the interpretive analysis” of the Guideline at issue.⁷⁴ Several times he asked ChatGPT: “What is the ordinary meaning of ‘physically restrained’?”⁷⁵ While he received a different answer each time, he found it informative that despite “subtle, marginal divergences,” the models’ responses “coalesce[d], substantively, around a common core—there was an objectively verifiable throughline.”⁷⁶ He wrote, “For our purposes, what matters is that the LLMs consistently defined the phrase ‘physically restrained’ to require the application of tangible force, either through direct bodily contact or some other device or instrument. And that, again, squares comfortably with the results obtained through the traditional, dictionary-driven breaking-and-repiecing method.”⁷⁷

This exercise gave Judge Newsom “hope that the models have something significant to offer the interpretive enterprise.”⁷⁸ “LLMs are trained on actual individuals’ uses of language in the real world,” he reasoned, and “the LLMs’ responses to [his] repeated queries reliably revealed . . . a common core.”⁷⁹ Judge Newsom concluded his *Deleon* concurrence by observing, “I continue to believe—perhaps more so with each interaction—that LLMs have something to contribute to the ordinary-meaning endeavor.”⁸⁰

In both *Snell* and *Deleon*, Judge Newsom favorably cited *Generative Interpretation* to support LLMs’ authority concerning the ordinary meaning of legal language.⁸¹ Like Arbel and Hoffman, Judge Newsom believes in asking *narrow* questions of LLMs, using them to answer discrete questions about the meanings of particular terms. But in other ways, his approach differs from theirs. For one thing, he takes a *narrow* view of the sources the LLM should specifically consider. Instead of giving the LLM an entire contract to interpret, Newsom uses single-sentence prompts, e.g., “*Is installing an in-ground trampoline ‘landscaping’?*”⁸² For another, he takes a slightly *broader* view of the LLM’s outputs than Arbel and Hoffman’s prenuptial-petition example above. Instead of considering token probabilities as they do, Newsom reads the LLM outputs as natural-language text, much like he would a dictionary definition.⁸³ In *Deleon*, he borrows from Arbel and Hoffman in repeating his

⁷³ *Id.*

⁷⁴ *Id.* at 1270 (Newsom, J., concurring) (quoting *Snell*, 102 F.4th at 1211 (Newsom, J., concurring)).

⁷⁵ *Id.* at 1272.

⁷⁶ *United States v. Deleon*, 116 F.4th 1260, 1275 (11th Cir. 2024) (Newsom, J., concurring).

⁷⁷ *Id.*

⁷⁸ *Id.* at 1276–77.

⁷⁹ *Id.* at 1277.

⁸⁰ *Id.*

⁸¹ See *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1226–27 n.7 (Newsom, J., concurring); *Deleon*, 116 F.4th at 1266 n.1 (Newsom, J., concurring).

⁸² *Snell*, 102 F.4th at 1235 (Newsom, J., concurring).

⁸³ See *id.* at 1228–30 (comparing LLMs and dictionaries).

query 30 times (10 for each of 3 models) to assess their “confidence,” but he continues to treat the outputs as English text.⁸⁴

3. Judge Joshua Deahl

Some months later, Judge Joshua Deahl of the District of Columbia Court of Appeals followed Judge Newsom’s lead. In *Ross v. United States*,⁸⁵ Nina Ross was charged with animal cruelty for leaving a dog in her car on a hot day.⁸⁶ In overturning her conviction, the majority held that there was insufficient evidence to prove beyond a reasonable doubt that leaving a dog in a car for over an hour on a hot day would harm the dog.⁸⁷ The majority cited *Jordan v. United States*,⁸⁸ a similar case where the court had reversed a conviction for leaving a dog outside on a cold day.⁸⁹ In his dissent, Judge Deahl distinguished *Jordan* on the ground that it is common knowledge that leaving a dog in a car on a hot day “created a ‘plain and strong likelihood’ that [the] dog would be harmed,” whereas the risks of leaving a dog outside on a cold day were less clear-cut.⁹⁰

In support of this argument, Judge Deahl and his clerks asked ChatGPT, “Is it harmful to leave a dog in a car, with the windows down a few inches, for an hour and twenty minutes when it’s 98 degrees outside?”⁹¹ They received an “unequivocal ‘yes.’” They then asked: “Is it harmful to leave a German shepherd outside in 25 degree temperature for five hours?” and received answers that “boil[ed] down to ‘it depends.’”⁹² Carefully caveating that ChatGPT is “definitely not” a good proxy for common knowledge, Judge Deahl nevertheless avowed that those chats supported his claims because inclusion in ChatGPT’s responses is “at least a point in favor of something being common knowledge . . . ; and a far stronger point against it [being common knowledge] if it’s not.”⁹³

Citing to Judge Newsom, Deahl used the text produced by chatbots as a proxy for common knowledge.⁹⁴ Still, there is something different going on in his approach. Judge Deahl didn’t ask ChatGPT about the meaning of a term (what does “harmful” mean?); he asked it about how the world works

⁸⁴ See *Deleon*, 116 F.4th at 1273–74 (Newsom, J., concurring).

⁸⁵ 331 A.3d 220 (D.C. 2025).

⁸⁶ See *id.* at 222.

⁸⁷ See *id.* at 224–27.

⁸⁸ 269 A.2d 848 (D.C. 1970).

⁸⁹ See *Ross*, 331 A.2d at 224 (citing *Jordan*, 269 A.2d at 849) (“In the absence of testimony . . . that the shelter or protection from the weather supplied this dog on this occasion would cause the dog to suffer, the evidence was insufficient to sustain the conviction.”).

⁹⁰ *Id.* at 235–36 (Deahl, J., dissenting).

⁹¹ *Id.* at 236.

⁹² *Id.* at 237.

⁹³ *Id.* at 236 n.4.

⁹⁴ See *id.* (calling *Snell* “a thoughtful and engaging discussion . . . which I largely agree with”).

(is it harmful to leave a dog in a car on a hot day?). That is a different kind of inquiry. It is still *narrow* in the question being asked, *narrow* in the sources presented to the LLM, and somewhat *broad* in considering natural-language outputs—but it is not really an interpretive question at all.

C. LLMs as Adjudicators

Some commentators go beyond mere LLM *interpretation* to LLM *adjudication*. Instead of prompting the LLM to carry out a specific interpretive task (“What is the ordinary meaning of ‘landscaping?’”) they prompt the LLM with a question about the case as a whole (“You are a federal appellate judge. Please read the attached briefs. Please write a three paragraph decision. In the third paragraph, explain who wins and why.”).⁹⁵ This is a *broad* interpretive task, much broader than anything attempted by Arbel and Hoffman, Judge Newsom, or Judge Deahl. It almost necessarily involves considering a *broad* range of sources—briefs, precedents, statutes, and other sources a human adjudicator would consider. And it equally almost necessarily involves reading the output *broadly* as a natural-language text—the kind of work product that a human adjudicator would produce.

Numerous commentators have considered whether AIs in general, and recent LLMs in particular, could do the work of judges. What sets apart the two we discuss in this section—Adam Unikowsky and Jack Kieffaber—is their full-throated enthusiasm. While others ask whether LLMs are ready for judicial use, these authors have already arrived at their answer: an emphatic “yes!”

1. Adam Unikowsky

Adam Unikowsky—a partner at Jenner & Block who focuses on Supreme Court and appellate advocacy⁹⁶—has used a variety of methods to test LLM adjudication. Although the specifics vary, from the jump he has focused on prompts that generate extensive natural-language outputs: his first, for example, was “Write two paragraphs about the [*Snell*] concurrence. First: what does the concurrence argue? Second, do you agree with the concurrence’s argument?”⁹⁷ From there, he progressed to uploading the entire briefing in a case to Claude, asking it to write a brief opinion, and then asking probing follow-up questions.⁹⁸

Unikowsky was deeply impressed with the results, asserting that the Claude chatbot “is fully capable of acting as a Supreme Court Justice

⁹⁵ See *In AI We Trust II*, *supra* note 5.

⁹⁶ Adam Unikowsky, JENNER & BLOCK LLP, <https://www.jenner.com/en/people/adam-g-unikowsky> [<https://perma.cc/3LST-BY3X>].

⁹⁷ *In AI We Trust II*, *supra* note 5.

⁹⁸ See *id.*

right now. When used as a law clerk, Claude is easily as insightful and accurate as human clerks, while towering over humans in efficiency.”⁹⁹ He found that Claude “consistently decide[d] cases correctly,” by which he meant that it articulated the holding the Supreme Court in fact reached.¹⁰⁰ On the occasions when Claude articulated conclusions that differed from the Supreme Court’s, its output struck Unikowsky as “reasonable.”¹⁰¹ He has also asked it to critique the methodology in an expert report (much as judges do when discounting an expert’s conclusions),¹⁰² to generate persuasive opinions in support of a particular outcome (“You are a Supreme Court Justice who believes racial-gerrymandering claims should be non-justiciable. Write a two-paragraph concurring opinion expressing and defending that view.”),¹⁰³ and to develop novel doctrines to resolve a case (“Please come up with an extremely creative, brilliant, out-of-the-box alternative to the ‘primary purpose’ test that would improve Confrontation Clause law.”).¹⁰⁴ In each case, he found the results to be intelligent, creative, and persuasive.

In additional blog posts Unikowsky has suggested that AIs could write effective appellate briefs,¹⁰⁵ and that AI should be used to automate the entire veterans’ appeals process.¹⁰⁶ His proposed process is worth quoting in detail, because it gives a good sense of the ambition of his vision:

This can be trivially implemented today. Just upload the relevant portion of the VA benefits manual into an AI context window, upload the veteran’s documentary evidence, and ask AI to apply law to fact. If you’re concerned that AI won’t be sufficiently generous, you can keep track of the percentage of veterans that obtain benefits and then tweak the prompts until the percentage is at a satisfactory level. If you’re squeamish about having AI make decisions that affect people’s lives, then you can have the AI complete the initial layer of review (equivalent to what the VA does today) and then give the veteran the right to appeal to a human judge. If we just snap our fingers, we can make all tradeoffs go away.¹⁰⁷

For Unikowsky, LLMs are ready to perform the work that human judges do; all that remains is to resolve the details of integrating the technology into the courts.

⁹⁹ *Id.*

¹⁰⁰ *Id.*

¹⁰¹ *Id.*

¹⁰² *See id.*

¹⁰³ *In AI We Trust II*, *supra* note 5.

¹⁰⁴ *A Brief History*, *supra* note 5.

¹⁰⁵ *See Automating Criminal Appeals*, *supra* note 5.

¹⁰⁶ *See Ignore the Future*, *supra* note 5.

¹⁰⁷ *Id.*

2. Legal Tech Entrepreneurs

Another enthusiastic promoter of LLM adjudication is Jack Kieffaber, a former clerk to Judge Newsom involved with AI arbitration startup Fortuna Arbitration.¹⁰⁸ Kieffaber has promulgated several papers arguing strongly in favor of replacing judges with AIs like the one Fortuna sells. The place to start is actually with his second article, *We Built Judge.ai. And You Should Buy It*, written with Kimo Gandall and Kenny McLaren.¹⁰⁹ The authors built an LLM-based online system, Arbitrus.ai, that receives evidence, motions, and briefing from two parties and then issues rulings.¹¹⁰ To calibrate and verify the system, they used another LLM to generate synthetic briefs in 100 hypothetical disputes, and then had Arbitrus.ai issue opinions resolving those disputes.¹¹¹ Humans checked that cited authorities existed, were responsive to the issues in the case, and fully resolved the dispute.¹¹² Overall, they claim that the system usually fully resolved each case without hallucinations or ungrounded assertions.¹¹³

As the name suggests, Arbitrus.ai is designed and marketed as an arbitration system. Although there are some statutory and doctrinal complications about the applicability of the Federal Arbitration Act to AI arbitrations, party consent can bless a wide range of procedures in arbitration that would be unacceptable in adjudication.¹¹⁴ Gandall, Kieffaber, and McLaren argue that “automation will beat all comers” in the market for arbitrators.¹¹⁵ This is a predictive claim about parties’ satisfaction with AI arbitrations, or at least their willingness to accept form contracts selecting it, that does not necessarily depend on the quality of those arbitrations.

Kieffaber’s principal normative claims in favor of AI adjudication come, instead, in his earlier article, *Predictability, AI, And Judicial Futurism*.¹¹⁶

¹⁰⁸ See Kimo Gandall, Jack Kieffaber & Kenny McLaren, *We Built Judge.ai. And You Should Buy It*, SSRN (Mar. 27, 2025), <https://papers.ssrn.com/abstract=5115184> [<https://perma.cc/d5jw-8y7k>] (“Mr. Gandall, Mr. Kieffaber, and Mr. McLaren all have financial interests in Fortuna-Insights, Inc., the company that makes Arbitrus.ai.”).

¹⁰⁹ See Gandall et al., *supra* note 30. Gandall and McLaren are CEO and CTO, respectively, of the startup company behind Arbitrus.ai. See ARBITRUS.AI, <https://www.arbitrus.ai> [<https://perma.cc/P2HT-6CXZ>].

¹¹⁰ See Gandall et al., *supra* note 30 at 30–34.

¹¹¹ See *id.* at 48–49.

¹¹² See *id.* at 50–56.

¹¹³ See *id.* at 56–57.

¹¹⁴ See, e.g., David Horton, *Forced Robot Arbitration*, 109 CORN. L. REV. 679, 679 (2024) (arguing that “AI adjudication is more likely to flourish in . . . forced arbitration”); Michael J. Broyde & Yiyang Mei, *Don’t Kill the Baby! The Case for AI in Arbitration*, 21 N.Y.U. J.L. & BUS. 119, 119 (2024) (“This article examines the integration of AI into arbitration, arguing that the Federal Arbitration Act (FAA) allows parties to contractually choose AI-driven arbitration, despite traditional reservations.”).

¹¹⁵ Gandall et al., *supra* note 30, at 83.

¹¹⁶ Jack Kieffaber, *Predictability, AI, And Judicial Futurism: Why Robots Will Run The Law And Textualists Will Like It*, 48 HARV. J.L. PUB. POL. (forthcoming) (manuscript at 78), <https://papers.ssrn.com/abstract=4966334> [<https://perma.cc/8KSE-6ZV8>].

He anticipates a future in which “Judge.AI” can do what Arbitrus.AI does, but on a much grander scale, providing both binding opinions on the “back end” and also advisory guidance for citizens on the “front end.”¹¹⁷ He assumes—big if true—that “Judge.AI is a perfectly neutral arbiter and interprets words with perfect mathematical accuracy.”¹¹⁸ He then argues that “Judge.AI is optimal under a textualist framework and, indeed, is the logical end result of the textualist project.”¹¹⁹

This view might be taken either as a criticism of textualism, or an invitation for textualists to embrace generative AI. Kieffaber comes closer to the latter, much like his former boss,¹²⁰ Judge Newsom, who justified his explorations in *Snell* and *Deleon* by pointing to generative AI’s ability to advance the textualist project. Kieffaber argues that the principal goal—the only goal—of textualism is predictability,¹²¹ and asserts that an ideal AI can be more predictable than human judges.¹²²

D. Varied Justifications

To be sure, there are some notable differences among the group. We have flagged their divergences on broad versus narrow questions, sources, and outputs. There is also a notable theoretical division. Judge Newsom and Kieffaber are whole-hearted textualists; they regard the central judicial task as the extraction of linguistic meaning from a controlling textual authority. Kieffaber is more ambitious than Judge Newsom in how much of that task he would assign to an LLM, but they share a common conception of what a judge—or an LLM prompted to do a judge’s job—ought to do, with a common focus on clarity and predictability. They are intrigued by LLMs because they see them as, potentially, the truest textualists of them all.

None of the others are firmly opposed to textualism as such, but they are more broad-minded about what judges—and LLMs in judicial robes—should be doing. Unikowsky’s prompts include a mixture of textualist and

¹¹⁷ *Id.* at 6–7 (“Given the rate at which these large language models develop, it is starting to feel eerily as though the future is *now* — and, if not *now*, who is to say it won’t be *next month*?”); *see also id.* at 78 (“My Judge.AI hypothetical, at the end of the day, really isn’t a hypothetical at all. *It’s coming.*”).

¹¹⁸ *Id.* at 7.

¹¹⁹ *Id.* at 8.

¹²⁰ Kieffaber’s online biography indicates that he clerked for Judge Newsom roughly between Summer 2023 and late September 2024; *Snell* and *Deleon* were published on May 28, 2024, and September 5, 2024, respectively. *See* 2023 James Wilson Fellows, JAMES WILSON INSTITUTE, <https://www.jameswilsoninstitute.org/articles/2023-james-wilson-fellows> [https://perma.cc/JF5H-8SG8] (last visited June 13, 2025) (listing Kieffaber as an “incoming law clerk” for Judge Newsom); The James Wilson Institute, *Predictability, AI, and Judicial Futurism with Jack Kieffaber – Anchoring Truths* (2024), <https://www.anchoringtruths.org/2024/09/26/predictability-ai-and-judicial-futurism-with-jack-kieffaber/> [https://perma.cc/PJB2-WPTJ] (last visited June 13, 2025) (“[Kieffaber] has clerked for Judge Kevin Newsom on the 11th Circuit.”).

¹²¹ *See* Kieffaber, *supra* note 116, at 9–13.

¹²² *See id.* at 29–31.

non-textualist instructions, much as one would expect from a skilled appellate lawyer who must craft arguments that are persuasive both to textualist and non-textualist judges. Judge Deahl's one experiment with ChatGPT is decidedly non-textualist; the question he put to it is about dog safety (or what people commonly believe about dog safety), not about the meaning of a disputed term. And Arbel and Hoffman see generative interpretation as the synthesis that will finally bring together the seemingly opposed camps of textualism and contextualism; to them, it offers the certainty textualists seek while considering all the relevant evidence that contextualists want to consult.

We will have much more to say in due course, but here at the outset, it is worth noting how profoundly these various advocates of LLM-based judging disagree on how to perform it and why it is justified. Judge Deahl thinks ChatGPT is valuable because it offers insight into what people know; Kieffaber denies that an ideal AI judge should consider anything besides linguistic meaning. Judge Newsom believes in prompting LLMs with disputed phrases in isolation; Arbel and Hoffman prefer to prompt them with every relevant document in the record they can. Unikowsky is perfectly happy to ask Claude to generate "off-the-wall" opinions to show off its creativity; Judge Newsom wants to restrict discretion as much as possible. And so on. That these LLM proponents have such fundamental differences of opinion suggests, perhaps, that none of them have rightly apprehended what generative interpretation is.

III. LLMs GENERATE TEXT BY PREDICTING IT

Before diving into our discussion of LLMs' use in adjudication and interpretation, we wish to emphasize an important point about how they function. We do not intend to rehash the history of LLMs or describe the technology in detail. At this point, there are excellent references available, and we assume that the reader has familiarity with the basics.¹²³

¹²³ See generally Memorandum from David Stein on AI Primer (July 12, 2024), <http://ai-memo.stein.fyi/> [<https://perma.cc/89AT-2449>] (providing an overview of the topic); Katherine Lee, A. Feder Cooper & James Grimmelmann, *Talkin' 'Bout AI Generation: Copyright and the Generative-AI Supply Chain*, 72 J. COPYRIGHT SOC'Y 251 (2025) (discussing how copyright law interacts with the chain of materials from truly novel human-made works to the outputs of AI algorithms). For readers looking for a deeper technical explanation of LLMs, there is a wealth of explanatory materials (of varying quality and accuracy) online and in the academic literature. Grant Sanderson's 2-hour video course on machine learning is by far the best mostly prerequisite-free resource we're aware of that covers the finer details of machine learning, neural networks, LLMs, and GPTs. 3BLUE1BROWN, *Playlist: Machine Learning*, (YOUTUBE, Oct. 5, 2017), https://www.youtube.com/watch?list=PLZHQObOWTQDNU6R1_67000Dx_ZCJB-3pi&v=aircAruvnKk [<https://perma.cc/SZB3-4EEN>]. The *Financial Times* has an accessible and accurate article on how LLMs generate text. See Madhumita Murgia, *Generative AI exists because of the transformer*, FIN. TIMES (Sep. 12, 2023), <https://ig.ft.com/generative-ai/> [<https://perma.cc/KUV9-KWTX>].

LLMs are at heart prediction machines. Given some text, an LLM predicts the next word.¹²⁴ The prediction process can be broken into two steps.

- Step 1: the LLM converts some input text (a “prompt”) into a sequence of numbers (a “vector”). Vectors represent the prompt in a computer-friendly format.
- Step 2: the LLM runs those vectors through a mathematical function that assigns a probability to each word in its dictionary. Those probabilities estimate the likelihood of each word directly following the prompt. For instance, given the prompt “*have a nice* _____,” a well-calibrated LLM might assign high probabilities to words like “*day*,” “*night*,” or “*stay*,” and almost zero probability to rare, grammatically incorrect, or semantically nonsensical words like “*thusly*.”¹²⁵

Generative AIs, including chatbots, use LLMs to generate text. A generative AI is a program that works by passing a prompt into an LLM to estimate next-word-likelihoods, using those likelihoods and a word selection algorithm (a “sampling strategy”) to select a word, then appending the selected word to the end of the prompt.¹²⁶ By repeating the process, a computer can generate large sequences of text.

Because the language models are where the magic of text generation happens, it is common to conflate the terms “generative AI” and “LLM.” LLMs are not manually programmed to predict words. Instead, their capabilities emerge during a process called “training.” Training involves repeatedly presenting text with omitted words to an LLM and programmatically adjusting the LLM’s configuration until it begins to correctly “guess” which word

¹²⁴ This is a simplification. LLMs can also “fill in” blanks in the middle of example text, or predict multiple words at a time. See generally Minkai Xu, Tomas Geffner, Karsten Kreis, Weili Nie, Yilun Xu, Jure Leskovec, Stefano Ermon & Arash Vahdat, *Energy-Based Diffusion Language Models for Text Generation*, 2025 PROC. INT’L CONF. ON LEARNING REPRESENTATIONS 1 (describing a technique that uses an LLM to predict many words simultaneously).

¹²⁵ For example, when presented with this example, OpenAI’s “babbage-02” model assigns the following probabilities: *day*: 47.7%; *weekend*: 9.8%; *evening*: 2.7%; *time*: 2.6%; *life*: 2.2%; . . . *picnic*: 0.0096%; . . . *thusly*: 0.000023%; . . .

¹²⁶ A simple sampling algorithm might pick the word assigned the highest probability. A more complex one might pick a word at random, weighted by the probabilities assigned by the LLM. In practice, word selection algorithms adjust the probabilities and then make a weighted random sample. The way those probabilities are adjusted is configured using “hyperparameters” with names like “temperature” and “top_p,” which you may have seen mentioned in related literature. See, e.g., *Fine-tuning Best Practices – OpenAI API*, OPENAI PLATFORM, <https://platform.openai.com/docs/guides/fine-tuning-best-practices#iterating-on-hyperparameters> [<https://perma.cc/9NW3-QA4P>] (last visited Oct. 20, 2025) (describing how to adjust sampling strategy hyperparameters when generating text); Xu et al., *supra* note 124. Note that because hyperparameters only affect the sampling process, they have no effect on how an LLM generates probabilities. See Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes & Yejin Choi, *The Curious Case of Neural Text Degeneration*, 2020 PROC. INT’L CONF. ON LEARNING REPRESENTATIONS 1, at 3–6 (describing the standard LLM text-generation workflow, where “models compute [the probability of each word] using the common left-to-right decomposition of the text probability,” which are then used to “generate the [output text] token-by-token using [one of the selection algorithms that accepts hyperparameters]”).

was omitted (that is, until it tends to assign high likelihood scores to missing words).¹²⁷ After training on billions of examples, some LLMs begin to make predictions that reflect surprising linguistic and contextual nuances—including grammar, semantic meaning, and style.¹²⁸

The LLM “learns” the patterns that appear in natural-language text: not just which words tend to follow which other ones (“nice day” is more common than “nice thusly”), but also which words are stylistically consistent with each other (few texts include both “detrimentally” and “butthead”), the ordering of textual passages (topic sentences precede supporting evidence), and other larger-scale and more abstract statistical patterns—which may or may not correspond to things readers identify as meaningful qualities of the text. Machine-learning researchers would say that an LLM is a model for the statistical distribution of natural-language text. Prompting one to generate text is a way of sampling from that distribution; if the model is a good one, the resulting texts should have the same statistical properties (word order, etc.) as the texts it was trained on.

In general, a reasonable, rough way to build instincts about how LLMs produce word likelihood estimations is to consider what factors would be relevant to filling in a scratched-out word in a document. For example, the fact that this sentence contains an em-dash—a generally uncommon punctuation mark beloved by law professors—might cause an LLM predicting the first word in the next sentence to assign significantly more likelihood to Latin words and prepositions. *Ergo*, some critics like Gary Marcus call LLMs “autocomplete on steroids,”¹²⁹ but this quality is also why some enthusiasts see them as ideal tools for interpretation.

It turns out that generating text by predicting it is strikingly effective at producing fluent text in a variety of genres. Unsurprisingly, some of those genres are typical work for lawyers: law school exams,¹³⁰ descriptive research

¹²⁷ See generally 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 Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever & Dario Amodei, *Language Models Are Few-Shot Learners* (July 22, 2020), <https://arxiv.org/pdf/2005.14165> [<https://perma.cc/VJB4-VXAL>] (discussing language model training).

¹²⁸ See generally Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Łukasz Kaiser & Illia Polosukhin, *Attention Is All You Need*, 31 CONF. NEURAL INFO. PROCESSING SYS. (2017) (discussing language model training and its outputs).

¹²⁹ Gary Marcus, *The Dark Rise of Large Language Models*, WIRED (Dec. 29, 2022), <https://www.wired.com/story/large-language-models-artificial-intelligence/> [<https://perma.cc/G78N-SEBL>].

¹³⁰ See *supra* note 9.

memos,¹³¹ and contract drafting,¹³² to name just a few. It is unsurprising that some researchers and judges have started to consider whether predictive text generation might also be a good fit for the work of judges. We now offer some reasons to question their optimism.

IV. IS GENERATIVE INTERPRETATION RELIABLE?

When a judge prompts an LLM, how do they know whether its response is based on insights about the meaning of language or on irrelevant cues in the prompt or training data? Intuitively, one should expect LLMs to use the semantic meaning of words to make objective choices—like dismissing nonsensical words—and to rely on other factors to predict the relative likelihood that semantically reasonable options comport with the rest of the text. But if LLMs' responses depend on factors other than the semantic meaning of particular fragments of text, then things are not so simple.

The problem is that LLMs' responses *do* depend on factors other than textual meaning. An LLM's output is determined by a complex generation algorithm, run on a model with a particular architecture and a specific set of weights, in response to a prompt that is both case-specific and expressed in natural language, subject to numerous configuration settings and implementation parameters, and typically starting from a randomly chosen seed value. All of these factors influence the output; that is simply the *definition* of what it means to carry out a generation using an LLM.

The reliability problem, then, is to show that in a given generative-interpretation protocol these other factors—model choice, prompt phrasing, configuration settings, etc.—do not significantly influence the aspects of the output being measured. If they do, then the protocol is measuring noise rather than signal.¹³³

In this section, we give reasons to believe that LLMs used for generative interpretation are indeed highly sensitive to irrelevant factors.¹³⁴ Even in

¹³¹ See, e.g., Daniel Schwarcz, Sam Manning, Patrick Barry, David R. Cleveland, JJ Prescott & Beverly Rich, *AI-Powered Lawyering: AI Reasoning Models, Retrieval Augmented Generation, and the Future of Legal Practice* (U. Mich. Pub. L. Working Paper 24-058), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5162111 [<https://perma.cc/X2LD-N8H3>].

¹³² See, e.g., *Navigating legal drafting: A how-to guide for law firms using AI-powered tools*, THOMSON REUTERS LEGAL BLOG (Oct. 2, 2024), <https://legal.thomsonreuters.com/blog/navigating-legal-drafting-a-how-to-guide-for-law-firms-using-ai-powered-tools/> [<https://perma.cc/NGM9-C6V7>].

¹³³ It might be tempting to avoid these issues by arbitrarily picking a model, settings, etc. That would make the protocol reliable in that it consistently generates similar results when repeated. But each arbitrary choice made to solve the reliability problem amplifies the epistemic problem, as we discuss *infra* Part IV.

¹³⁴ Work by other scholars reinforces our conclusions. See generally, e.g., Jonathan H. Choi, *Large Language Models Are Unreliable Judges* (manuscript at 21–30) (Aug. 11, 2025), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5188865 [<https://perma.cc/ZAD6-GU42>]; Waldon et al., *supra* note 16, at 39–47.

the LLM proponents’ carefully curated case studies, changes to minor implementation details have outcome-determinative effects, calling into question whether they are accurately measuring semantic meaning. We consider case studies from *Generative Interpretation* and show that in them, the results are sensitive to the choice of model, the prompting strategy, or implementation details.

Our choice to critique Arbel and Hoffman’s case studies should not be taken as singling out their methodology for special criticism. Quite the opposite: of the generative-interpretation proponents, they are by far the most careful and attentive to the reliability challenges. We use their case studies *because* they are admirably explicit and precise about their methods. Our point is that if even the most meticulous generative-interpretation experiments conducted to date face reliability challenges—as we show they do—then one should be even more skeptical of approaches that do not share Arbel and Hoffman’s attention to detail.

A. Famiglio: *Model Instability*

One extraneous factor to the semantic task is the choice of which model to use. A diamond’s mass should not vary based on whether it is weighed with a pan balance or a digital scale; words’ meanings should not vary based on whether they are interpreted by Claude or ChatGPT. But of course there are substantial differences in generative models; companies compete fiercely to differentiate them, and there are leaderboards and communities dedicated to comparing models and teasing apart their differences. Even within a model family, outputs can vary enormously from one version to the next—as though the definition of thousands of words in a dictionary changed each time it was reprinted.

Consider the *Famiglio* case study from *Generative Interpretation*.¹³⁵ The Famiglios had a prenuptial agreement that defined the marriage’s duration as the period from the wedding date to the date when either spouse filed for divorce.¹³⁶ The wife filed for divorce, withdrew her filing, then refiled in earnest several years later.¹³⁷ Millions hinged on whether the first or second filing date determined the marriage term.¹³⁸

Arbel and Hoffman present the contract clause to an LLM to demonstrate a generative-interpretation technique that, they argue, “offers courts a better sense of the relevant probabilities” of different interpretations of the prenup, assuming “the parties were intending to use English in its most public and common sense.”¹³⁹ This particular technique involves looking “under the

¹³⁵ *Generative Interpretation*, *supra* note 2, at 483.

¹³⁶ *Id.* (citing *Famiglio v. Famiglio*, 279 So.3d 736, 737–38 (Fla. Dist. Ct. App. 2019)).


¹³⁷ *Id.*

¹³⁸ *Id.*

¹³⁹ *Id.* at 485.

hood” of an LLM.¹⁴⁰ Instead of letting a chatbot pick words using an LLM’s word-likelihood estimates, they observe those estimates directly. They ask an LLM which filing date would control, observe the probabilities (reproduced in Figure 1), and conclude that the model favors the second date.¹⁴¹ At time of writing, the version of the language model used in *Generative Interpretation* is no longer available.¹⁴² A later version of the same model yields the conflicting result shown in Figure 2.


The second filing would determine the number of full years of marriage



Word	Probability
second	94.72%
date	4.44%
first	0.68%
number	0.13%
amount	0.01%

FIGURE 1: Reproduced from *Generative Interpretation*, at 484. This figure shows the probabilities produced by the davinci-003 model while responding to a question asking which date controls under the *Famiglio* facts.

The first filing would determine the number of full years of marriage



Word	Probability
first	52.55%
second	23.15%
date	20.51%
number	1.03%

FIGURE 2: The probabilities produced by the gpt-3.5-turbo-instruct model under otherwise identical conditions to those in Figure 1.

That is, while the model used by Arbel and Hoffman predicts that “second” is more likely than “first” by a factor of more than 100 to 1, a later version of the same OpenAI model predicts that “first” is more likely than “second” by a factor of more than 2 to 1. The fact that these two models’ predictions diverge so dramatically makes it harder to credit either of them. If they were predicting the frequency of “first” versus “second” based solely on

¹⁴⁰ *Id.* at 484.

¹⁴¹ *Generative Interpretation*, *supra* note 2, at 484–85.

¹⁴² *GPT-4 API General Availability and Deprecation of Older Models in the Completions API*, OPENAI (Apr. 24, 2024), <https://openai.com/index/gpt-4-api-general-availability/> [<https://perma.cc/K4FR-QJ2J>] (deprecating the “davinci-003” model used in *Generative Interpretation*, and recommending users transition to the “gpt-3.5-turbo-instruct” model used in this essay).

a deep understanding of language, we should expect consistency. But since they diverge, how should courts determine which model is more appropriate for a given question? At least with dictionaries, we can make inferences about editorial choices. LLMs are far more opaque. Perhaps models will converge in the future, but for now they do not always do so,¹⁴³ and that is a problem for those who want to perform generative interpretation.¹⁴⁴

B. Famiglio: Prompting Instability

Continue with the *Famiglio* example. There is a deeper problem here. What does the probability of 94.72% that the first model attached to “second” signify? It is tempting to say that it represents the model’s estimate of the probability that a typical English speaker *would believe that the second filing controls* (as the sentence overall asserts).¹⁴⁵ But that is an unwarranted leap. The probability represents only the model’s estimate of the probability that the word “second” *would fill the blank* in the sentence “The ____ filing would determine the number of full years of marriage.” This estimate, however, is exquisitely sensitive to the precise way the sentence is phrased.

Recall that an LLM predicts one word at a time, then text-generation software uses those predictions to select a word, adds that word to the text, and repeats the process. By exploring each possible word choice, we can map out a decision tree of possible generated texts.¹⁴⁶ Figure 3 shows the probabilities for different word sequences generated by GPT-3.5 in response to the *Famiglio* case study’s prompt. To read the chart, follow a path from left to right, noting the probability of each word given the previous words on that path. For example, there is a 52.5% probability that the first word is *The*. When the first word is *The*, there is a 20.5% chance that the second word is *date*. Following different paths reveals how prior word choices affect GPT’s predictions. Note that the relative likelihood of *first* and *second* shifts based on phrasing. For “*the ____ filing*,” the first filing is twice as likely, while for “*the date of the ____ filing*,” the second filing is five times more probable.¹⁴⁷

¹⁴³ Cf. *United States v. Deleon*, 116 F.4th 1260, 1272–76 (11th Cir. 2024) (Newsom, J., concurring) (querying ChatGPT, Claude, and Gemini with identical prompts, and discussing the variations in their answers).

¹⁴⁴ Cf. *Generative Interpretation*, *supra* note 2, at 501 (suggesting that “more sophisticated models tend to converge on meaning”).

¹⁴⁵ See *id.* at 485 (“Generative interpretation in this simple case thus offers courts a better sense of the relevant probabilities if the parties were intending to use English in its most public and common sense.”).

¹⁴⁶ For an interactive tool showing a similar probability tree for the conversation in Judge Newsom’s *Snell* concurrence, see David Stein, *GPT’s Possible Response to Judge Newsom’s Experiments in Snell v. United* (Sep. 28, 2025), <http://snell-gpt.stein.fyi/> [<https://perma.cc/QSD7-73LC>].

¹⁴⁷ Cf. *Generative Interpretation*, *supra* note 2, at 485 n.164.

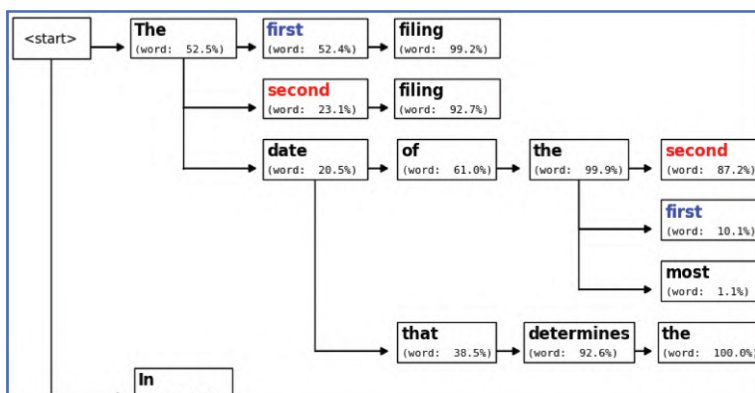


FIGURE 3: The decision tree gpt-3.5-turbo produces in response to the input used to produce Figure 2. Note that the relative probability of “first” and “second” flip based on the LLM’s phrasing choices.

One might try to fix this syntactic sensitivity by restricting the AI to one-word responses, perhaps by presenting the question as multiple-choice.¹⁴⁸ While that eliminates the phrasing decisions, it does so by arbitrarily committing to one specific phrasing. This behavior tells us nothing about the relative merits of the conflicting probability estimates we observed; it merely endorses one option while obscuring others.¹⁴⁹ In fact, presenting the *Famiglio* case as multiple-choice amplifies generative interpretation’s sensitivities. As shown in Figure 4, GPT-3.5 almost exclusively relies on the order of multiple-choice options.

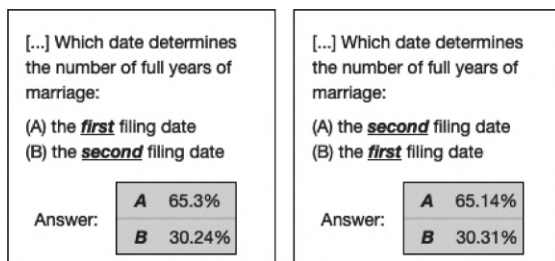


FIGURE 4: Phrasing the *Famiglio* case study as a multiple-choice question results in GPT ignoring substance and picking option “A.” Probabilities generated using gpt-3.5-turbo-instruct.

¹⁴⁸ See *id.* (proposing this solution).

¹⁴⁹ In fact, having the LLM “think out loud” when generating text is a common technique for improving the accuracy and quality of answers. See, e.g., Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le & Denny Zhou, *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models*, 36 CONF. NEURAL INFO. PROCESSING SYS. 1 (2022). At time of writing, clever application of chain-of-thought techniques is a meaningful differentiator for the generative AI techniques that perform best against major benchmarks. Daya Guo et al., *DeepSeek-R1: Incentivizing Reasoning Capability in LLMs via Reinforcement Learning*, ARXIV:2501.12948 [cs] (Jan. 2025), <https://arxiv.org/pdf/2501.12948> [<https://perma.cc/W3GM-D9VB>].

Similar problems haunt attempts to cross-validate between models.¹⁵⁰ If models disagree, which should we trust? If models agree, how do we know which factors they converged around? As Table 1 shows, OpenAI’s GPT-3.5 and Anthropic’s Claude-3 models both favor option “A,” regardless of what that option represents.¹⁵¹ Mere convergence does not imply convergence around information we consider relevant.¹⁵²

¹⁵⁰ But cf. *Generative Interpretation*, *supra* note 2, at 503–04 (“[A]s a best practice, judges would do well to cross-verify the answers that they get from one platform against another.”).

¹⁵¹ This is true of almost every version of both models. The table below shows the probabilities generated by every version of GPT-3.5 and Claude-3 available at time of writing. Note that the preference for option “A” is stable for most models, even when the contents of options “A” and “B” are swapped.

Model		Option		LLM Prediction	
		A	B	A	B
Flagship Models	gpt-3.5-turbo	first	second	91%	2%
	gpt-3.5-turbo	second	first	91%	1%
	claude-3-sonnet-20240229	first	second	96%	1%
	claude-3-sonnet-20240229	second	first	94%	5%
Other Models	gpt-3.5-turbo-0125	first	second	89%	4%
	gpt-3.5-turbo-0125	second	first	91%	3%
	gpt-3.5-turbo-0301	first	second	0%	92%
	gpt-3.5-turbo-0301	second	first	47%	23%
	gpt-3.5-turbo-0613	first	second	8%	4%
	gpt-3.5-turbo-0613	second	first	14%	5%
	gpt-3.5-turbo-1106	first	second	35%	26%
	gpt-3.5-turbo-1106	second	first	63%	3%
	claude-3-haiku-20240307	first	second	100%	0%
	claude-3-haiku-20240307	second	first	59%	41%
	claude-3-opus-20240229	first	second	9%	90%
	claude-3-opus-20240229	second	first	100%	0%

(Some predictions do not add to 100% because the model sometimes doesn’t commit to an option in the first few words.)

GPT-4 and GPT-4o are harder to measure because they often refuse to pick “A” or “B,” instead writing a few equivocating sentences. They appear to be less sensitive to the order in which options are presented. Instead, they converge on the importance of the parties’ names: both models are nearly twice as confident in their answers when the spouse filing for divorce is named “Jennie” and the spouse paying alimony is named “Mark.” This chart shows the probability that GPT-4 models immediately commit to the second filing date. That is, GPT-4 models equivocate about 80% of the time if Mark is paying Jennie, and only 50% of the time when Jennie is paying Mark.

model	Filing Spouse	Paying Spouse	P(“second”)
gpt-4-turbo	Jennie	Mark	18.30%
gpt-4-turbo	Mark	Jennie	49.00%
gpt-4o	Jennie	Mark	23.50%
gpt-4o	Mark	Jennie	47.30%

¹⁵² One might reasonably expect LLMs to converge more quickly around simple rules like basic syntax, or simple concepts like “people tend to pick the first option,” and more slowly around complex concepts like ways to resolve semantic ambiguity. There’s promising research exploring ways to probe the inner workings of LLMs, but those research efforts are years away from practical application. *See, e.g.*, Leo Gao, Tom Dupré la Tour, Henk Tillman, Gabriel Goh, Rajan Troll, Alec Radford, Ilya Sutskever, Jan Leike & Jeffrey Wu, *Scaling and evaluating sparse autoencoders*, ARXiv:2406.04093 [cs] (June 6, 2024), <https://arxiv.org/pdf/2406.04093> [<https://perma.cc/CU9Y-5QUG>]; Adly Templeton, Tom Conerly, Jonathan Marcus, Jack Lindsey,

Model	Option		LLM Prediction	
	A	B	A	B
gpt-3.5-turbo	first	second	91%	2%
gpt-3.5-turbo	second	first	91%	1%
claude-3-sonnet-20240229	first	second	96%	1%
claude-3-sonnet-20240229	second	first	94%	5%

TABLE 1: GPT and Claude models converge on the same, obviously arbitrary answer when presented with the *Famiglio* case study rephrased as a multiple-choice question.

Methods for determining which factors sway a model’s choices is an active area of technical research.¹⁵³ Still, in this case, GPT and Claude’s fixation on option order is not overly surprising. Semantic meaning is less useful for predicting which of two subjectively reasonable options is more likely to appear in text. In those situations, other factors (e.g., formatting) offer greater predictive power relative to semantic meaning and thus more heavily influence LLM predictions.

These examples illustrate how a multitude of factors—including many that users of generative interpretation would deem irrelevant—influence LLM predictions. Multiple-choice bias might feel obvious and easily mitigated. Other outcome-determinative choices are more subtle, like asking about “the first filing” (not “*the date of* the first filing”), or Judge Newsom asking about the “*ordinary* meaning” of landscaping.¹⁵⁴ Controlling for specific kinds of unwanted influences offers no assurances that LLMs will use relevant information. It merely introduces new, less obvious biases.¹⁵⁵

C. Snell: *Instability through Randomness*

We have seen that generative interpretation can be extremely sensitive to the choice of model and to the phrasing of the prompt. But these are still only two of the factors that affect a model’s outputs. Another is the randomness

Trenton Bricken, Brian Chen, Adam Pearce, Craig Citro, Emmanuel Ameisen, Andy Jones, Hoagy Cunningham, Nicholas L Turner, Callum McDougall, Monte MacDiarmid, Alex Tamkin, Esin Durmus, Tristan Hume, Francesco Mosconi, C. Daniel Freeman, Theodore R. Sumers, Edward Rees, Joshua Batson, Adam Jermyn, Shan Carter, Chris Olah & Tom Henighan, *Scaling Monosemanticity: Extracting Interpretable Features from Claude 3 Sonnet*, ANTHROPIC (May 21, 2024), <https://transformer-circuits.pub/2024/scaling-monosemanticity/> [<https://perma.cc/H3RK-6GFE>].

¹⁵³ See *supra* note 134.

¹⁵⁴ For additional discussion and data on the effect of phrasing on Judge Newsom’s *Snell* experiments, see *infra* note 220 and accompanying text.

¹⁵⁵ See Justin Curl, Peter Henderson, Kart Kandula & Faiz Surani, *Judges Shouldn’t Rely on AI for the Ordinary Meaning of Text*, LAWFARE (May 22, 2025), <https://www.lawfaremedia.org/article/judges-shouldn-t-rely-on-ai-for-the-ordinary-meaning-of-text> [<https://perma.cc/U3YJ-4Y4K>].

inherent to some LLM systems, including every major chatbot: asking the same question twice can yield different outputs.¹⁵⁶

Some of the consequences of this randomness can be surprising. Consider *Snell*, where Newsom asked ChatGPT whether installing an in-ground trampoline is “landscaping.” He included the following transcript of his conversation with ChatGPT:

Question: What is the ordinary meaning of “landscaping”?

Response[GPT]: “Landscaping” refers to the process of altering the visible features of an area of land, typically a yard, garden or outdoor space, for aesthetic or practical purposes. This can include activities such as planting trees, shrubs, flowers, or grass, as well as installing paths, fences, water features, and other elements to enhance the appearance and functionality of the outdoor space.

Question: Is installing an in-ground trampoline “landscaping”?

Response[GPT]: Yes, installing an in-ground trampoline can be considered a part of landscaping. Landscaping involves altering the visible features of an outdoor area for aesthetic or practical purposes, and adding an in-ground trampoline would modify the appearance and function of the space. It’s a deliberate change to the outdoor environment, often aimed at enhancing the overall landscape and usability of the area.¹⁵⁷

We posed the same questions to ChatGPT-3.5—and got back the opposite answer:

Question: What is the ordinary meaning of “landscaping”?

[GPT-3.5]: The ordinary meaning of “landscaping” refers to the act or process of designing, shaping, or beautifying the outdoor area around a building or property. This may include planting trees, shrubs, and flowers, as well as creating paths, lawns, and other features to enhance the aesthetic appeal of the space. Landscaping also encompasses maintenance and upkeep of these elements to keep the outdoor area looking visually appealing.

Question: Is installing an in-ground trampoline “landscaping”?

[GPT-3.5]: While installing an in-ground trampoline may involve shaping and altering the ground surface, **it is not typically considered landscaping** because it does not involve the traditional elements of design and beautification of an outdoor space.

¹⁵⁶ There is some well-documented non-determinism in commercial LLMs. See Seherman Chann, *Non-Determinism in GPT-4 is Caused by Sparse MoE*, 152334H (Aug. 5, 2023), <https://152334H.github.io/blog/non-determinism-in-gpt-4/> [<https://perma.cc/UE3N-3JAT>]. The magnitude of the non-determinism is fairly small.

¹⁵⁷ *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1234 (11th Cir. 2024) (Newsom, J., concurring) (emphasis added).

What explains this discrepancy? It turns out that ChatGPT’s answer to Newsom’s second question was heavily influenced by how it phrased its answer to the first. For example, if ChatGPT uses the phrase “other features” at the end of its list of examples, it is three times more likely to include trampoline installation in the definition of “landscaping” than to exclude it. Conversely, if ChatGPT uses the phrase “ordinary meaning” when answering the first question, it will tend to *exclude* trampoline installation from the definition.

These wording choices are determined through a randomized process, essentially reducing ChatGPT’s determination of whether trampoline installation falls within the meaning of the word “landscaping” to a roll of the dice.¹⁵⁸ The following table lists a few more examples:

If GPT-3.5’s response to Question 1 contains...	the “yes” to “no” likelihood ratio in GPT’s responses to Question 2 will...
The phrase “ordinary meaning”	drop by 29.8%
A list of examples ending with “or other elements” (versus “or other features”)	increase by 71.4%
A list of examples beginning with “such as” (versus “including” or “for example”)	increase by 32.6%
A list of examples containing two man-made features (e.g., fences, paths, retaining walls, etc.)	increase by 40.3%

TABLE 2: The effect of phrasing choices on GPT’s relative likelihood to say that an in-ground trampoline falls within the ordinary meaning of the word “landscaping” when presented with the questions in Judge Newsom’s *Snell* concurrence.

It is possible to explore GPT-3.5’s entire decision tree to build intuitions about what influences Judge Newsom’s experiments. As shown below, the probability of GPT saying that a trampoline is not landscaping drops precipitously if it uses the word “typically,” and increases substantially if it starts its sentence with “the” and then puts the word “landscaping” in quotes.

¹⁵⁸ Ari Holtzman, Jan Buys, Li Du, Maxwell Forbes & Yejin Choi, *The Curious Case of Neural Text Degeneration*, 2020 PROC. INT’L CONF. ON LEARNING REPRESENTATIONS (describing the word-selection technique used by most modern text generation AIs, and explaining that always picking the most likely word option leads to flat-sounding text).

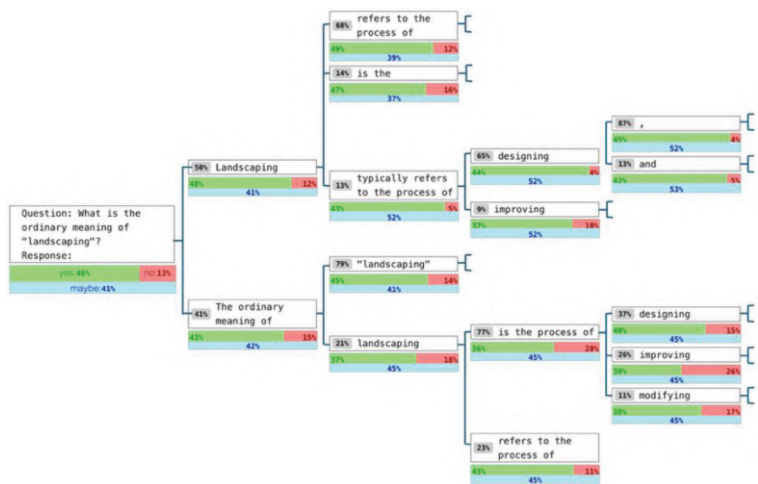


FIGURE 5: Part of the tree of possible responses to Judge Newsom’s first question in *Snell*. The bars next to each word represent the probability that GPT’s answer to the second question will be yes (green), no (red), or a non-committal response (blue). An interactive version of this tree exploration tool is available at snell-gpt.stein.fyi.

What should we make of this path dependence? One lesson, surely correct, is the one that Newsom, Arbel and Hoffman, and other pioneers in generative interpretation have drawn: it is important to conduct experiments multiple times.¹⁵⁹ Random variation can be quantified and managed statistically; indeed, statistical methods are essential to empirical validation of generative interpretation.

Another and subtler lesson is that LLMs can influence *themselves* in un-intuitive and extraneous ways. ChatGPT’s answers to Newsom’s first questions appear highly similar; the random variation between them does not substantially affect the interpretive takeaways. But the differences in those answers, insignificant though they seem, are enough to channel ChatGPT towards diametrically opposed answers to Newsom’s second question. If ChatGPT were responding solely based on the linguistic meaning of “landscaping,” the random variation in whether it used the phrase “typically” in answer to the first question should not influence its answer to the second. And yet it does.

¹⁵⁹ See, e.g., *United States v. Deleon*, 116 F.4th 1260, 1273–76 (11th Cir. 2024) (Newsom, J., concurring) (running each query ten times, and discussing the dependence of outputs on random variations). But see Gandall et al., *supra* note 30, at 32 (assuming “perfect predictability”).

D. Katrina Canal Breaches: *Implementation Instability*

Even this does not exhaust the factors that influence a generative AI's word choices. The generation process also depends on a large number of configurable parameters and other implementation details.

Consider another case study from *Generative Interpretation*, one that revisits a series of Fifth Circuit cases involving an insurance contract. In relevant part, it asks whether man-made disasters can fall within the meaning of the word “flood” as used in one of the contract's clauses.¹⁶⁰ To answer the question, Arbel and Hoffman analyze the “vectors” LLMs use to represent text.

A quick refresher: vectors are sequences of numbers that LLMs use to represent words and phrases. It is often helpful to think of them as coordinates. LLMs tend to place similar text in close proximity: the “London” and “Berlin” vectors lie closer to each other than to the “Umbrella” vector.¹⁶¹ LLMs also tend to arrange vectors so direction is meaningful: a line running from “London” to “England” might run roughly parallel to the line connecting “Berlin” to “Germany.”¹⁶² Those properties make it possible to mathematically estimate similarities and relationships between words and phrases.¹⁶³

Applying this insight to the flood question, Arbel and Hoffman estimate the similarity between the insurance clause and a list of terms related to natural and unnatural causes of floods. They include a few unrelated terms as quality checks.¹⁶⁴ Their result is reproduced in Figure 6.

¹⁶⁰ See *Generative Interpretation*, *supra* note 2, at 453–55, and accompanying footnotes (citing *In re Katrina Canal Breaches Consolidated Litig.*, 466 F. Supp. 2d 729, 747–63 (E.D. La. 2006); and Willy E. Rice, *The Court of Appeals for the Fifth Circuit: A Review of 2007–2008 Insurance Decisions*, 41 TEX. TECH. L. REV. 1013, 1039 (2009).

¹⁶¹ See Tomáš Mikolov, Ilya Sutskever, Kai Chen, Greg Corrado & Jeffrey Dean, *Distributed Representations of Words and Phrases and Their Compositionality* (Oct. 16, 2013), <https://arxiv.org/pdf/1310.4546> [<https://perma.cc/M4D7-J4TS>].

¹⁶² *Id.*

¹⁶³ See Tomáš Mikolov, Wen-tau Yih & Geoffrey Zweig, *Linguistic Regularities in Continuous Space Word Representations*, 2013 PROC. CONF. N. AM. CHAPTER ASS'N FOR COMPUTATIONAL LINGUISTICS: HUM. LANGUAGE TECHS. 746, 746.

¹⁶⁴ *Generative Interpretation*, *supra* note 2, at 457.

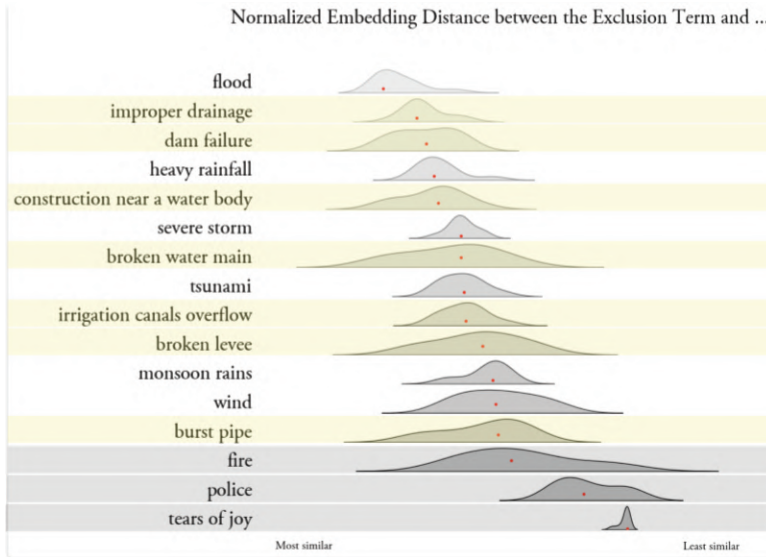


FIGURE 6: Reproduced from *Generative Interpretation*, at p. 457. This chart shows how different LLMs measure the distance between an insurance clause and various terms. We added highlights to demark *Generative Interpretation*’s “man-made” exclusion terms (yellow) and “quality check” terms (grey).

Generative Interpretation explains that the further the average “similarity” (denoted by a red dot) is to the right, the less semantically related that term is to the insurance clause. Noting that natural and unnatural causes are interleaved, they conclude that whether a flood was man-made or natural is not a major factor in determining its relationship to the insurance clause.¹⁶⁵ They present this insight as “objective, cheap support for the court’s judgment that floods can be unnaturally caused.”¹⁶⁶ With caveats, they suggest this kind of chart contains information that would be useful to courts, perhaps even superior to dictionary definitions:

[T]he model doesn’t provide (nor could it) a scientific answer to the question of whether certain words are sufficiently close to make the plain meaning of *flood* unambiguous. That choice is ultimately a normative one which judges must make. But there is a bit of difference between an informed conclusion based on a statistical analysis of billions of texts and a judgment by a few dictionary editors. And there is an ocean of difference between the baroque and expensive textualism the court used and code that is cheap,

¹⁶⁵ *Id.*

¹⁶⁶ *Id.*

replicable, quick, and most importantly, extremely straightforward to use. Simply put, generative interpretation is good enough for many cases that currently employ more expensive, and arguably less certain, methodologies.¹⁶⁷

But to generate that chart, Arbel and Hoffman make a series of implementation choices—and changes to these technical choices change the outcome of the experiment. Their measurements of similarity in ten LLMs’ vector representations of terms¹⁶⁸ are calculated using a measure called “cosine distance.”¹⁶⁹ But using cosine distance requires implementors to choose a reference frame¹⁷⁰ and, under certain circumstances, normalize

¹⁶⁷ *Id.* at 458. We did indeed find the code provided in tandem with the article cheap (~\$50), quick (a few hours), and relatively straightforward to use.

¹⁶⁸ *Id.* at 456 n.22.

¹⁶⁹ *Id.* (describing and justifying the use of the cosine distance metric). Cosine distance measures how far apart two points appear when viewed from a third reference point. For example, measuring the distance between stars by measuring how far apart they appear in sky. See G. Salton, A. Wong & C.S. Yang, *A Vector Space Model for Automatic Indexing*, 18 ASS’N. COMPUTING MACHINERY 613 (1975). This is a standard measure of similarity. It is popular because in high-dimensional spaces, *direction* is often more meaningful than absolute distance. (To see why, for example, consider a line, circle, and sphere with the same diameter. The longest possible line that can fit within each shape is the same, but the number of directions available increases with each new dimension. The amount of information conveyed by direction grows with each new dimension, while the number of possible distances remains the same).

¹⁷⁰ In high-dimensional spaces, the perspective from which you measure direction matters. For example, the vectors for “London” and “Berlin” will tend to fall in roughly the same direction from most vantage points (e.g., the relationship between “swimming pool” and “London” is roughly the same as the one between “swimming pool” and “Berlin”). Differences start to show up when measured from a relevant perspective (e.g., the relationship between the words “England” and “London” are very different than the relationship between the words “England” and “Berlin”). The two reasonable measurement choices for this case study are (a) to use some arbitrary neutral perspective and measure the distance between each term and the insurance clause, or (b) to use the perspective of the insurance clause and measure the similarity between terms. *Generative Interpretation* chooses the default, neutral perspective. Cf. Mikolov, Yih & Zweig, *supra* note 154, at 749 (illustrating how the direction between two vectors can be meaningful, and describing how to apply semantically meaningful linear transformations—what I refer to as “change of perspective”—to embedding spaces).

the results.¹⁷¹ These are choices so seemingly minor they don't even show up in the footnotes—they are subtleties buried in the implementing code.¹⁷²

And yet these choices make all the difference. We reproduced the “flood” case study four times. Each iteration used the same cosine-distance method described in *Generative Interpretation*, with the same models, terms, embeddings, data, and measurement techniques. The only differences are that we tried replacing Arbel and Hoffman's choices of reference frame and normalization with reasonable alternatives. In one trial (top left in Figure 7) we kept both of their choices; in a second (top right) we replaced their reference frame with one from the perspective of the insurance clause; in a third (bottom left) we kept their reference frame but used a simpler and more natural normalization algorithm; and in a fourth (bottom right) we made both changes.

The four trials reached four completely different results. Switching from Arbel and Hoffman's normalization technique to a simpler one weakened their results; switching from their choice of reference frame to another logical one caused the results to fail their quality check; and doing both at once completely flipped the result—seemingly showing that natural and unnatural causes are *not* significantly interleaved.

Our point is not that Arbel and Hoffman are making *unreasonable* choices; it is that they are making *unacknowledged* choices. Are the inferences enabled by these charts really “informed conclusion[s] based on a statistical analysis,”¹⁷³ or do they mistake implementation artifacts for semantic meaning; noise for signal?

¹⁷¹ *Generative Interpretation* applies min-max linear unit normalization without outlier correction to the measurements from each model. Because they intentionally include outliers as quality checks, the results are warped and hard to interpret. For example, the difference between *flood* and *fire* varies significantly less across models than the difference between *flood* and *police*. But *fire*'s curve is more spread out in the chart reproduced in Figure 6 because there is significant variance between *fire* and the “quality check” term *tears of joy*. Replacing the term *tears of joy* with a more fire-related term (e.g., *burn ward*) would cause “police” to appear closer to the origin than “fire.” Even the authors seem to get tripped up by this choice. *Id.* at 504 (“the policy exceptions were closer to ‘fire’ than to . . . ‘police’”). A more straightforward approach might measure everything relative to the distance between the reference clause and “flood.” See, e.g., Prashant V. Kamat, *Absolute, Arbitrary, Relative, or Normalized Scale? How to Get the Scale Right*, 4 ACS ENERGY LETT. 2005 (2019) (describing relative scale and discussing the difference between relative and normalized scales).

¹⁷² The specific code changes are available at <https://gist.github.com/davidbstein/27360bd c86bc5adc7a3050e44368ec15> [<https://perma.cc/7SVJ-556Z>]. The top-left chart in Figure 7 is reproduced from *Generative Interpretation*, *supra* note 2, at 457. Based on our reading of the article, the top bottom-right chart seems to track most closely to the context-dependent interpretive approach described in the article.

¹⁷³ *Generative Interpretation*, *supra* note 2, at 458.

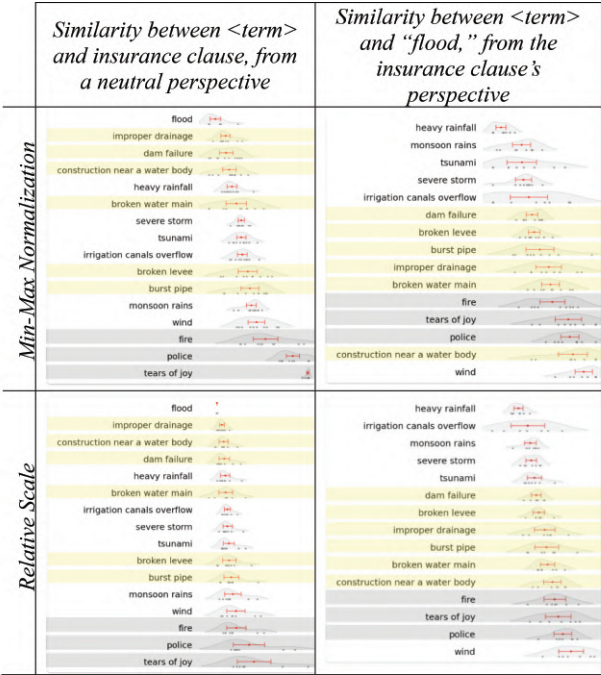


FIGURE 7: The data and method in *Generative Interpretation’s* Hurricane Katrina case study plotted using two different reference frames and normalization approaches. Highlights to demark *Generative Interpretation’s* “man-made” exclusion terms (yellow) and “quality check” terms (grey). Depending on implementation choice, the chart either: contradicts the Fifth Circuit (top left), supports the Fifth Circuit (bottom right), indicates an inconclusive measurement (bottom left), or fails *Generative Interpretation’s* “quality check” (top right).

Imagine a judge trying to use embeddings as interpretive tools or contracting parties pre-committing to an interpretation technique, as *Generative Interpretation* suggests they ought to do.¹⁷⁴ How should they navigate these kinds of implementation choices, possibly without expert guidance? How should they resolve disagreements about those choices? How can they determine which implementation choice is better?¹⁷⁵ They need a principled framework to pick one implementation over another, but *Generative Interpretation* does not supply one. None of the LLM proponents does.

¹⁷⁴ *Generative Interpretation*, *supra* note 2, at 501.

¹⁷⁵ Another issue: it is not clear—though it seems incredibly likely—that embedding vectors produced by LLMs have the same semantic meaning as embedding vectors used by other kinds of AI. Harald Steck, Chaitanya Ekanadham & Nathan Kallus, *Is Cosine-Similarity of Embeddings Really About Similarity?*, 2024 ACM WEB 887 (“cosine-similarity can yield arbitrary and therefore meaningless ‘similarities.’”). Compare with *Generative Interpretation*, *supra* note 2, 458 n.24 (citing to papers about embedding distances generated using non-LLM technologies).

E. Methodological Instability

There is a larger issue here. As computer scientist Andrew Tanenbaum quipped, “The nice thing about standards is that there are so many to choose from. And if you do not like any of them, just wait a year or two.”¹⁷⁶ The same is true of generative interpretation. It is not just that there are unsettled details that affect the results of a commonly used interpretive method. It is that the *interpretive method itself* is unsettled. “Generative interpretation” does not refer to a specific approach to using LLMs to answer a particular type of judicial question. Rather, it is a family name for a disparate collection of methods, many of which have in common that they make use of LLMs. Indeed, these methods diverge dramatically in terms of how to set up the query and how to interpret the results.

As we saw in Part I, generative interpretation’s proponents disagree on:

- Whether to ask discrete interpretive questions (Arbel and Hoffman, Judge Newsom, Judge Deahl) or to resolve entire disputes (Unikowsky, Kieffaber).
- Whether to prompt the LLM with a short question in isolation (Judge Newsom, Judge Deahl) or to give it as many materials pertaining to a dispute as possible (Arbel and Hoffman, Unikowsky, Kieffaber).
- Whether to interpret the LLM by using technical tools like word probabilities and cosine distances (Arbel and Hoffman) or by reading its outputs as natural language (Judge Newsom, Judge Deahl, Unikowsky, Kieffaber).

These are not small differences. There is no strong reason to expect that they will all yield similar results, and no strong reason to think one of them is obviously better than another. Indeed, even within this taxonomy, each box contains multitudes.

Consider, for example, the four different ways in which *Generative Interpretation*’s case studies produce numerical estimates of interpretive meaning.

Method 1: Character Values as Probabilities

As their first example, Arbel and Hoffman prompt ChatGPT with text from an insurance policy and instruct the chatbot,

please state your prediction—with the associated numerical level of confidence in parentheses—on the likely expectations of most policyholders under these terms for the following proposition[]:
1. The policy will provide compensation for losses resulting from a substantiated third-party burglary.¹⁷⁷

¹⁷⁶ ANDREW S. TANENBAUM & DAVID J. WETHERALL, *COMPUTER NETWORKS* 702 (5th ed. 2011).

¹⁷⁷ *Generative Interpretation*, *supra* note 2, at 475 n.129.

ChatGPT replies, “1. Likely Expectation (90%): The policy will compensate for third-party burglary.”¹⁷⁸ Arbel and Hoffman assert that this response constitutes “Chat GPT-4 [telling] us that it was 90% likely that the policy would pay in response to a ‘substantiated third-party burglary.’”¹⁷⁹ As they explain, “90%” is the string of text that ChatGPT predicted would appear at that particular place in its response.¹⁸⁰ It reflects an AI model’s determination about the string of *characters* that was most likely to be produced in a written response to that question. It does not reflect a direct, *numerical* assessment of the probability “that the policy would pay in response to a ‘substantiated third-party burglary.’”¹⁸¹

Method 2: Floating-Point Values as Probabilities

For its next example, *Generative Interpretation* invokes a number with an entirely different significance. This is the *Famiglio* case study discussed above, in which the authors extract a probability from the internals of GPT’s word-selection algorithm. Given the choice between the words “first” and “second,” Arbel and Hoffman find the AI will select the word “second” roughly 94.72% of the time (and our attempt to replicate this result with a newer LLM yielded a 54.16% chance of “first”).¹⁸² From this example, Arbel and Hoffman conclude, “Generative interpretation . . . thus offers courts a better sense of the relevant probabilities if the parties were intending to use English in its most public and common sense.”¹⁸³ We take them to be arguing that the different “probability” values that the model assigns to “second” versus “first” constitute “relevant probabilities” for ascertaining the meaning of the prenuptial agreement.¹⁸⁴

The numbers that Arbel and Hoffman are evaluating in this prenup example differ dramatically from the numbers that they evaluated in the earlier insurance-policy example. In the insurance-policy example, they evaluated the string “90%” that the model’s text-prediction engine produced as a response to their written query. In the *Famiglio* case study, they evaluate

¹⁷⁸ *Id.*

¹⁷⁹ *Id.* at 476.

¹⁸⁰ *Id.* at 482.

¹⁸¹ *Id.*

¹⁸² *Id.* at 484.

¹⁸³ *Generative Interpretation*, *supra* note 2, at 482.

¹⁸⁴ Arbel and Hoffman also caution that “the probabilities shouldn’t be interpreted literally. The model could, for example, continue the sentence with ‘The first filing would not control.’” *Id.* at 485 n.165. We interpret this warning to mean that the relative probabilities of “second” and “first” in their example do not necessarily reflect the respective probabilities that the model would have generated complete sentences that endorse each respective reading. This recognition seems appropriate, but acknowledging it undermines the probative value of those numbers.

the model's internal estimate of the probability that the text "second" should follow the text "The" in its response. *These numbers are the results of completely different processes.* We looked "under the hood" at the responses ChatGPT gave to the authors' first question about the meaning of the insurance policy and found that the model assigned an approximate probability of .5416 to "90%," a probability of .013 to "0%," and a probability of .001 to "yes."¹⁸⁵

There may or may not be some basis on which to ascribe legally relevant meaning to either the text "90%" from the first example or the number from the second example. How should courts understand and differentiate between the "90%" from the first example and the .9472 from the second example (or the .5416 from our replication attempt)?¹⁸⁶ Both can be notated as numerical percentages, but they are produced by vastly different processes.

Method 3: A Different Method for Characters as Probabilities

The third method returns to textual percentages. The authors pass the entire text of the promissory note from *Trident Center v. Connecticut General Life Ins. Co.*¹⁸⁷ to an LLM with the instructions: "Rate the following proposition on a 1-100 scale, where 0 is wrong and 100 is correct...[:] The language is REASONABLY susceptible to being read as providing the borrower the right to early prepayment." They repeat this process 100 times for each of 3 popular chatbots and report the results in a chart (reproduced in Figure 8(a)).¹⁸⁸ They conclude "the models roughly agree on average that prepayment is not allowed, with a mean score of ~41."¹⁸⁹

¹⁸⁵ Generated using the OpenAI chat completions API, using GPT-4 with recommended settings. Code available at <https://gist.github.com/davidbstein/f95d908344205ef583d1f57392ace76f> [<https://perma.cc/L9WL-34AY>].

¹⁸⁶ For example, how should the inherent biases in number-as-word selection inform interpretation of those numbers? If you ask GPT-3.5 to fill in the blank on "Rating = __ / 100" with no additional context, about half the time it will pick one of: 50 (p=.1488), 0 (p=.1220), 10 (p=.0667), 90 (p=.0558), or 60 (p=.0542). Computed using the OpenAI "playground" using the gpt-3.5-turbo-instruct model on default settings.

¹⁸⁷ 847 F.2d 564 (9th Cir. 1988).

¹⁸⁸ Taken from the source code accompanying *Generative Interpretation*, *supra* note 2, available at <https://github.com/yonathanarbel/generativeinterpretation/blob/67dc982ec69aa383d3c2422dda55294af68158a7/Scripts/Trident.ipynb> [<https://perma.cc/A8A4-JYWC>]. Input Box 19, line 9.

¹⁸⁹ *Generative Interpretation*, *supra* note 2, at 488.

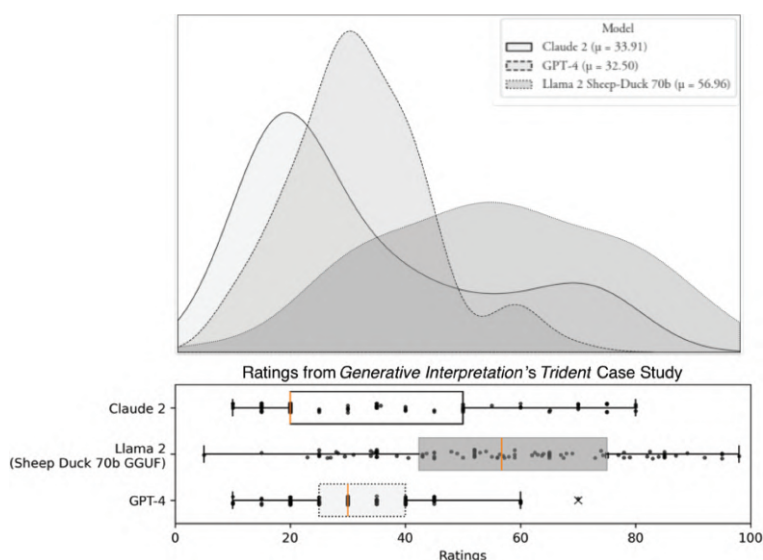


FIGURE 8: Results from the *Generative Interpretation's* Trident case study. (a) Above: the estimated “density” of responses, as presented in the original paper. (b) Below: the same data, displayed as a box-and-whisker chart. We include the second chart to make the data points (black dots) and averages (red lines) explicit.¹⁹⁰

Method 4: From “Yes/No” to Numerical Distributions

In a fourth example, Arbel and Hoffman present a contract to several LLMs and solicit yes/no answers about the meaning of a clause. They pose 20 linguistic “variations of the same legal question,” a “yes/no question[] where yes indicates agreement with the judge’s interpretation,” and graph the results.¹⁹¹ Here, the authors are not instructing the model to respond with its confidence, nor are they looking under the hood at the probability of the model

¹⁹⁰ *Id.* (asking models to rate whether prepayment is allowed on a scale of 0-100, where 0 means repayment is definitely not allowed, and 100 means repayment is definitely allowed). The curves shown in *Generative Interpretation* are computed using a method called kernel density estimation. As a part of that method, the authors needed to choose between one of several bandwidth “estimators,” the function used to estimate density based on a collection of datapoints. The different methods result in significant differences in the resulting chart. See David M. Bashtannyk & Rob J. Hyndman, *Bandwidth Selection for Kernel Conditional Density Estimation*, 36 COMPUTATIONAL STAT. & DATA ANALYSIS 279, 279 (2001) (reviewing methods for picking a bandwidth function). See also Part III, *infra*, discussing minor implementation details that affect interpretive results.

¹⁹¹ *Generative Interpretation*, *supra* note 2, at 490.

picking a certain response. Instead, this number is the relative frequency with which a model provides a “yes” or “no” answer. From a technical perspective, this method is mathematically equivalent to measuring the number from method 2 for each of the 20 questions, then averaging the results and adding a large dose of random noise.¹⁹²

* * *

We do not have a definitive opinion about which of these four methods is best.¹⁹³ Nor do we fault Arbel and Hoffman for exploring numerous different approaches. That is exactly what exploratory early work needs to do. Our point is that an embarrassment of riches can still be an embarrassment. Any one generative interpretation methodology raises difficult reliability problems on its own. Adding three more does not resolve those problems; it multiplies them. It may be that one or more of these methods can be made robust and legitimate. But a judge considering generative interpretation must still choose which ones to try, and that by itself is an entire can of worms.

V. IS GENERATIVE INTERPRETATION EPISTEMICALLY JUSTIFIED?

In this Part, we turn to the epistemic problem: should judges accept LLM outputs as providing authoritative guidance on legal issues? For generative interpretation, those issues involve the linguistic meaning of disputed texts. For generative adjudication, those issues involve the proper resolution of disputed legal issues. But in both cases, the epistemic question is the same. What makes *this* measure relevant to deciding *that* issue? To assert that an LLM output “accurately reflects real people’s everyday speech patterns,”¹⁹⁴

¹⁹² The choice between “yes” and “no” is technically equivalent to the choice between “first” or “second” in *Famiglio*: a binary selection between two possible words. We could have looked “under the hood” to get the underlying probabilities for each phrasing. Instead, for each phrasing, we use the result of a single weighted coin toss based on those probabilities. Even taking this number at face value, it is troubling: would we ascribe the same meaning to a human interlocutor who answered an identical question differently when polled multiple times in quick succession, or would we simply find their response incoherent? Cf. Thomas R. Lee & Jesse Egbert, *Artificial Meaning?*, 77 FLA. L. REV. (forthcoming) (manuscript at 43), <https://papers.ssrn.com/abstract=4973483> [<https://perma.cc/42P3-UHDW>] (observing that variation in LLM responses “is not indicative of the variability we would observe between different individuals in a larger population—just the variability we would see between responses from a single person who is repeatedly asked the same question.”).

¹⁹³ Though method 4 seems strictly worse than method 2.

¹⁹⁴ *United States v. Deleon*, 116 F.4th 1260, 1270 (11th Cir. 2024) (Newsom, J., concurring). To be clear, Newsom goes on to provide reasons, recognizing that this is a proposition that must be demonstrated, not simply assumed. As we detail in this Part, our point of divergence is that we think the reasons he gives do not warrant the conclusion that LLM outputs “accurately reflect[]” ordinary usage. *Id.*

without more, is to beg the question. *Why would* the outputs of a complicated algorithm necessarily reflect *anything*, let alone “everyday speech patterns?” Something more is needed to show that *this* algorithm, run in *this* way on *this* input, yields “correct” or “accurate” answers to the question being asked.

This is a general challenge for empirical methods; it goes by many names, including “construct validity.”¹⁹⁵ A method may be consistent and robust enough to measure *something*, but that *something* may not be the question the experimenter hopes to answer. A moment’s thought shows that establishing construct validity requires attending to the details of not just the method but also the question. The Schönhage-Strassen algorithm multiplies large numbers significantly faster than the algorithm taught in elementary school; it is an excellent source of knowledge about the products of integers, but it tells us nothing about ordinary English usage.

LLM proponents have offered, we think, four broad types of responses to the epistemic problem: they have tried to justify the use of LLMs based on *how they work*, based on *calibrating* their outputs against an accepted baseline, based on the rhetorical *persuasiveness* of their outputs, and based on a claim that they offer *predictability* as such. All four types are intuitively appealing, and three of them can be sound ways of establishing a method’s epistemic *bona fides*. But we think that in this case, for LLMs as a source of legal meaning, the case is unproven. LLM proponents’ claims depend on a mix of unstated assumptions, mistaken beliefs about how LLMs work, and misunderstandings of the judicial function.

A. *Relying on the Training Process?*

Start by considering the class of arguments that an LLM should be trusted because of what it is: a statistical model trained on a corpus of natural-language text. The problem here is that these are explanations of why an LLM *could* work at the task of capturing common linguistic usage, not an explanation of why a specific LLM *actually does* capture common linguistic usage in a specific setting.

1. *Deductive Arguments*

The strongest argument for a decisional legal AI system would be to validate the system *deductively*, using the tools of formal logic to craft an abstract model of law and formal verification to show that the AI system correctly implements that abstract model.¹⁹⁶ Indeed, there have been numerous previous attempts to bring software-based deductive rigor to legal reasoning,

¹⁹⁵ See Dag I. K. Sjøberg & Gunnar R. Bergersen, *Construct Validity in Software Engineering*, 49 IEEE TRANSACTIONS SOFTWARE ENG’G 1374, 1374 (2023).

¹⁹⁶ See generally James Grimmelmann, *Programming Languages and Law: A Research Agenda*, 2022 PROC. SYMPOSIUM COMPUT. SCI. & L. 1 (discussing formal approaches to law).

with mixed success.¹⁹⁷ To validate generative interpretation deductively would require demonstrating, *a priori*, that the internal workings of the LLMs *necessarily* lead it to estimate linguistic meaning to some degree of accuracy.

Deduction has the virtue of generating conclusions that are guaranteed to be true in specific instances. Of course, the veracity of a valid, deductive proof is as strong (and only as strong) as its assumptions. If those assumptions hold, the conclusion always follows. For example, say we want to know whether a two-sided coin is fair. If we assume that the coin is infinitesimally thin and perfectly symmetrical, that it is flipped with a fixed upward velocity v and an angular velocity drawn from a known distribution X at a fixed distance d above the ground, that it is unaffected by air resistance, and that it collides inelastically with the ground when it lands, then for suitable choices of v , X , and d , we can establish that this method for tossing coins is necessarily fair. The probability that it comes up heads on any given toss is equal to the probability that it comes up tails.

But proving LLMs' legal-interpretive authority isn't like positing the fairness of a coin, for two reasons. First, while "fairness" in coin-flipping has a straightforward definition—an independent, 50-50 expectation of heads or tails for any given flip—"accuracy" in legal interpretation has no equivalent formal specification. The second reason follows from the first. Because legal-interpretive accuracy is underspecified, we don't know what factual premises about LLMs would have to hold in order for them to be accurate. By contrast, because a coin's fairness is easily formalized, we *can* posit the factual properties that characterize a fair coin. Put another way, there are two things that stand in the way of using LLMs to reason deductively about law: the nature of law and the nature of LLMs.

Even assuming a deductive model of legal interpretation is reconcilable with law as a social enterprise, adopting such a model would require transformational changes to present-day attitudes. Indeed, even generative interpretation proponents seem to reject a deductive approach. In *Snell*, Judge Newsom wrote, "[w]ould the consideration of LLM outputs in interpreting legal texts inevitably put us on some dystopian path toward 'robo judges' algorithmically resolving human disputes? I don't think so. As Chief Justice Roberts recently observed, the law will always require 'gray area[]' decision-making that entails the 'application of human judgment.'"¹⁹⁸

Setting aside whether law can ever be a deductive enterprise, LLMs are extraordinarily poor vehicles for deductive validation.¹⁹⁹ The origins of

¹⁹⁷ See, e.g., L. Thorne McCarty, *Reflections on TAXMAN: An Experiment In Artificial Intelligence and Legal Reasoning*, 90 HARV. L. REV. 837 (1976). Note that McCarty, writing nearly 50 years ago, still correctly identified the ambiguities and open texture of law as key challenges to its formalization.

¹⁹⁸ *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1232 (11th Cir. 2024) (Newsom, J., concurring).

¹⁹⁹ Cf. E. Mark Gold, *Language Identification in the Limit*, 10 INFO. & CONTROL 447 (1967) (formally proving that it is impossible to learn something perfectly and completely through example-based processes like machine learning process).

LLMs trace back to a decision to dispense with logical rigor and precisely curated knowledge bases, and to rely instead on recognizing patterns in massive datasets.²⁰⁰ Today's LLMs did not develop because computer scientists deduced that a particular technical architecture would, as a matter of logical necessity, produce AI with remarkable text-generating abilities. Instead, they came about when engineers decided to run with the techniques that produced the best results, even when the reasons these techniques work so well were (and, in some cases, are) yet to be discovered.²⁰¹ As Arbel and Hoffman note in *Generative Interpretation*: "whatever [an LLM] tells you, it is really no explanation at all . . . working with LLMs admittedly requires a leap of faith, a realization that no better explanation is forthcoming."²⁰²

Kieffaber deals with these problems simply by assuming them away. One of the six premises of his "sci-fi-hypothetical" is that "Judge.AI is a perfectly neutral arbiter and interprets words with *perfect mathematical accuracy*."²⁰³ This assumption serves as an intuition pump for his thought experiment about the nature of textualism,²⁰⁴ but it tells us nothing about actually existing LLMs. Arbitrus.AI is an actual system, and while Kieffaber and his coauthors are enthusiastic about its abilities, they do not attempt to validate it deductively. They don't even detail their model architecture, training algorithm, training datasets, or prompts, let alone show that these particular design choices are logically guaranteed to produce correct outputs. To state the idea is to refute it. LLMs are not spherical cows.²⁰⁵

2. *Informal Arguments*

In practice, most LLM proponents offer a related but much weaker argument. As Unikowsky puts it:

To me, the proposition "AI is useful for determining the ordinary meaning of English words" should be approximately as controversial as "GPS is useful for determining directions." Cutting-edge LLMs have read the entire Internet, or almost all of it anyway. Of course they'll be useful in determining how words are ordinarily used.²⁰⁶

²⁰⁰ See Benjamin L.W. Sobel, *Artificial Intelligence's Fair Use Crisis*, 41 COLUM. J.L. & ARTS 45, 58 (2017).

²⁰¹ See generally Alon Halevy, Peter Norvig & Fernando Pereira, *The Unreasonable Effectiveness of Data*, 24 IEEE INTELLIGENT SYS. 8 (2009).

²⁰² *Generative Interpretation*, *supra* note 2, at 483.

²⁰³ Kieffaber, *supra* note 116, at 7 (emphasis added).

²⁰⁴ One might object that treating natural language as susceptible to "perfect mathematical" interpretation not only assumes away the debate between textualism and its critics but also distorts beyond recognition the texts that textualists understand themselves to be interpreting.

²⁰⁵ See David Kaiser, *The Sacred, Spherical Cows of Physics*, NAUTILUS (Apr. 25, 2014), <https://nautil.us/the-sacred-spherical-cows-of-physics-234898/> [<https://perma.cc/F3ZZ-FYWR>].

²⁰⁶ *In AI We Trust I*, *supra* note 5.

Similarly, Judge Newsom writes, “LLMs can be expected to offer meaningful insight into . . . ordinary meaning . . . because the internet data on which they train contain *so many* uses . . . , from *so many* different sources”²⁰⁷

These are not deductive proofs that the way LLMs are created means they *necessarily* provide accurate information on semantic meaning. Instead, they are arguments that the way LLMs are created gives us *good reason to think* that they provide information on semantic meaning.

Consider Unikowsky’s GPS analogy. The GPS system uses a set of satellites with precisely calibrated clocks that transmit synchronized radio signals; a receiver uses the timing information from signals from multiple satellites to calculate its position relative to them—and thus its position on the surface of the Earth.²⁰⁸ Given the design of the system—and the relevant orbital mechanics, wave physics, electronics, and so on—we have good reason to think that the calculation carried out by a particular receiver in a particular place at a particular moment is accurate. Similarly, given some reasonable assumptions about coins’ weight, dimensions, and initial flipping characteristics, we might have good reason to think that particular coin flips are likely to be close to fair.²⁰⁹

Once again, however, it is a much heavier lift for generative interpretation than for other methods, because LLMs are much more complicated than a coin or a GPS receiver, and generative interpretation asks LLMs to resolve far more complicated problems. The argument that LLMs work because they have been exposed to almost “the entire Internet” proves far too much. The Google Books corpus was assembled by scanning millions of hard-copy books.²¹⁰ The scanners that Google used have “read” millions of books; does it follow that “they’ll be useful in determining how words are ordinarily used?” With enough time, you could tap the entire corpus that trained ChatGPT in binary into the keys of a calculator. You could place a tortoise in front of a speaker broadcasting the entire corpus as audio. At the end of these (interminable) processes, the calculator will have “read” the same data that ChatGPT “read,” and the tortoise will have apprehended it, at least in the sense that they will have been presented with that information. Yet neither the calculator, nor the tortoise, nor the scanner will be a greater authority on ordinary meaning than it was before these indoctrination processes.

The same point holds regardless of whether the claim is phrased in terms of LLMs’ accuracy, their legitimacy, their adherence to democratic norms, or any other desideratum. There is nothing about “training on the entire Internet”

²⁰⁷ *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1227 (11th Cir. 2024) (Newsom, J., concurring).

²⁰⁸ See *How GPS Works*, GPS.Gov, https://www.gps.gov/sites/default/files/2025-08/Educational_Poster.pdf [<https://perma.cc/UVM4-DSF7>].

²⁰⁹ But see Persi Diaconis, Susan Holmes & Richard Montgomery, *Dynamical Bias in the Coin Toss*, 49 SIAM REV. 211, 211 (2007) (“We show that vigorously flipped coins tend to come up the same way they started.”).

²¹⁰ *Authors Guild v. Google, Inc.*, 804 F.3d 202, 208 (2d Cir. 2015).

that implies the output will satisfy any particular criteria. The calculator and the tortoise were trained on the entire Internet too; that doesn't make them democratically legitimate, either.

Mere exposure to large amounts of natural-language text does not automatically confer authority about linguistic meaning. Indeed, the road to ChatGPT is littered with the corpses of generative AIs that were emphatically terrible authorities on ordinary meaning.²¹¹ It was not that OpenAI's model-training methods were *a priori* better than all that had come before in a way that could have been expected to yield far more reliable results. Instead, a series of advances in model architecture, feasible model size, and large, high-quality training datasets came together to produce a model with demonstrably—but very surprisingly—better *a posteriori* performance.

In other words, “it trained on the entire Internet” is a plausible theory of why an LLM that usually generates fluent text does so, but does nothing to demonstrate that it will. What makes it reasonable to think that an LLM will usually generate fluent text is *that it has generated fluent text*. The demonstration is fundamentally empirical.

Now, to be sure, the fact that ChatGPT-3.5 generates fluent text also makes it more plausible that ChatGPT-4 will do so too—they are trained in similar ways on similar sources of data. Similarly, ChatGPT's fluency can transfer, in part, to Gemini, to Claude, and so on—it has become widely accepted knowledge that certain types of LLM architectures and training tend to work reasonably well. But notice that while we can appeal to “how Claude was trained” to validate Claude's fluency by drawing on ChatGPT's, this appeal still rests on the empirical demonstration of ChatGPT's fluency. The appeal to how LLMs are trained adds no new and additional justification; it just provides a way of organizing and marshalling the empirical evidence from various LLMs' usage. And so, it is to those empirical justifications that we now turn.

B. Relying on Empirical Measurements?

The second way one might validate generative interpretation is *empirically* demonstrating the integrity of a particular methodology for querying LLMs. Instead of deductively or informally arguing *a priori* that LLMs are necessarily or likely to be accurate, this approach demonstrates inductively on the basis of evidence that a particular process for using LLMs produces correct answers consistently, predictably, and robustly. If satisfied, this inquiry justifies trusting LLM outputs produced through that validated process. If an LLM consistently produces accurate information via a specified method in representative testing scenarios, then subsequent uses following the same

²¹¹ See generally JANELLE SHANE, YOU LOOK LIKE A THING AND I LOVE YOU: HOW ARTIFICIAL INTELLIGENCE WORKS AND WHY IT'S MAKING THE WORLD A WEIRDER PLACE (2019) (extensively describing AIs' humorously bad attempts to emulate human writing).

method in similar scenarios are probably trustworthy—and we may be able to quantify that probability.

Returning to our two-sided coin: if we toss a coin 400 times and it comes up heads 199 times and tails 201 times, we are now 95% confident the coin is fair within a 5% margin of error when tossed in the same way as the experiment. Additional experimentation could increase confidence and reduce the margin of error. Though empirical tests cannot provide the perfect confidence of a deductive proof, they can provide strong, statistical evidence that a certain process has a tested property.

1. *Internal Robustness*

But what would it actually take to show that a given LLM provides accurate linguistic interpretations? The coin has the advantage that the problem can be specified finitely and precisely: there is a single “flip” process with exactly two possible outcomes, and we seek to learn their respective probabilities. The problem of generative interpretation is significantly more open-ended, and the LLM is a significantly more complicated object of study.

Recall that an LLM is a prediction machine that uses whatever it can find in a textual prompt to make accurate predictions about what comes next. That might be the semantic meaning of the term it has been asked to gloss. Or it might be whether the prompt author asked the LLM to “explore” an issue or to “delve” into it.²¹² If authors who eschew contractions are more conservative and sesquipedalian, then a prompt that eschews contractions is more likely to elicit an output that not only eschews contractions, but also reads in a conservative and sesquipedalian register in other ways. While it is likely (though unproven) that insights about general language understanding are tucked away somewhere inside LLMs’ internal structures, that information is intertwined with myriad other factors that can help the predictive task.²¹³ Currently, there is no easy way to discern which information an LLM uses to answer questions. It might use the information courts care about. It might use some unrepresentative subset of that information. It might rely on some inscrutable blend of factors that courts would consider irrelevant. These are null hypotheses: possibilities that can only be ruled out by conducting tests to rule them out.

In short, we are completely dependent on running actual studies to generate the kind of empirical evidence needed to validate LLMs’ use for

²¹² See, e.g., Tom S. Juzek & Zina B. Ward, *Why Does ChatGPT “Delve” So Much? Exploring the Sources of Lexical Overrepresentation in Large Language Models*, 31 PROC. INTERN. CONF. COMPUTATIONAL LINGUISTICS 6397 (2025).

²¹³ For an accessible discussion exploring what LLMs “know,” see Kenneth Li, *Do Large Language Models learn world models or just surface statistics?*, THE GRADIENT (Jan. 21, 2023), <https://thegradient.pub/othello/> [<https://perma.cc/Q4JB-NK4V>], and the underlying paper, Kenneth Li, Aspen K. Hopkins, David Bau, Fernanda Viégas, Hanspeter Pfister & Martin Wattenberg, *Emergent World Representations: Exploring a Sequence Model Trained on a Synthetic Task*, 2023 PROC. INT’L CONF. ON LEARNING REPRESENTATIONS 1.

generative interpretation. The LLM proponents have started on this important task, some more rigorously (Arbel and Hoffman) and some more informally (Judge Newsom, Judge Deahl, and Unikowsky). And this is exactly right—a journey of a thousand miles begins with a single step.

But as we showed in the previous Part, the empirical evidence to date falls well short of showing that generative interpretation really is justified. In particular, we documented that existing methods are not robust. Switching models can flip the direction of a result; small changes to queries produce large variations in the responses.²¹⁴ If protocols *A* and *B* to measure phenomenon *X* are identical except in some small respect that has no apparent connection with *X*, but *A* and *B* produce very different results, it counts as evidence that *neither A nor B* is actually measuring *X* well. Waldon et al.'s results are similar to ours; they “show that LLMs’ metalinguistic judgments are highly sensitive to subtle prompting variations” and that “LLMs can be easily ‘gami-fied’ to reflect a user’s preconceived biases, even when the lawyer or judge adheres to what appears to be a rigorous and neutral methodology.”²¹⁵

This is a problem common to all empirical disciplines, and the obvious way to respond to it is to improve methodological robustness. If *A* and *B* produce highly correlated measurements despite their differences, and so do variations *C*, *D*, *E*, and so on, then we are more justified in believing that all the members of this family really are measuring the same thing. And some of this work is already occurring. From *Snell* to *Deleon*, Judge Newsom went from two queries to each of two models to ten queries to each of three models and more systematically analyzed them. In their work on LLM interpretation, Christoph Engel and Richard McAdams analyzed distributions of responses rather than a model’s single “best” response.²¹⁶ As the sophistication and scale of generative-interpretation experiments continues to increase, its practitioners may be able to find robust protocols that are not sensitive to small differences in experimental setups.

This is not to say that a protocol must be robust against all possible variations. Some can be excluded for principled reasons. Say a coin flip gets fairer—in a way that can be empirically measured with strong statistical significance—as the coin is tossed from a greater height. If that’s the case, the protocol in which a coin is tossed from two inches above a table may produce different results than the protocol in which it is tossed from two feet, but that will be because the two-inch protocol is demonstrably worse. From political polling to laser interferometry, empiricists in numerous fields have accumulated bodies of standards and best practices for conducting their studies using

²¹⁴ See Choi, *supra* note 134, at 20–30 (discussing sensitivity to prompts); *id.* at 30–36 (discussing wide variability between LLMs).

²¹⁵ Waldon et al., *supra* note 16, at 39.

²¹⁶ See Christoph Engel & Richard H. McAdams, *Asking GPT for the Ordinary Meaning of Statutory Terms*, 2024 U. ILL. J.L. TECH. & POL’Y 235, 235 (2024).

specific methods in ways that are broadly accepted over their alternatives.²¹⁷ Here too the LLM proponents are tentatively beginning to move beyond simply brainstorming different possible methods to testing them comparatively, and excluding the ones that fail robustness checks. Engel and McAdams, for example, have an extended discussion of different types of design choices for generative-interpretation studies, including tentative thoughts on how to make many of those choices.²¹⁸

It is worth emphasizing how daunting this task is for generative interpretation. LLMs are sensitive to their configuration and prompts *by design*. Consumer-grade LLMs aim to respond in ways their users prefer; a chatbot that can conform to its user’s preferences by picking up on subtle cues in their inputs will provide a better service.²¹⁹ Consider *Snell*: if Newsom’s clerk had asked, “whats landscaping mean?”, he would have seen vastly different results to those he received in response to “What is the ordinary meaning of ‘landscaping’?”²²⁰ Professional- and business-grade LLMs are used inside tools that need to conform to organizations’ preferred tone, values, opinions, and topics. For example, an LLM-powered filing assistant on a court website should never acknowledge potential innuendos or double-entendres, whereas

²¹⁷ See, e.g., THE AM. SOC’Y OF MECH. ENG’RS, B89.1.8 - PERFORMANCE EVALUATION OF DISPLACEMENT-MEASURING LASER INTERFEROMETERS (2011), <https://www.asme.org/codes-standards/find-codes-standards/performance-evaluation-of-displacement-measuring-laser-interferometers/2011/pdf> [<https://perma.cc/ZZ6P-ZFCJ>]; AM. ASSOC. FOR PUB. OP. RSCH., BEST PRACTICES FOR SURVEY RESEARCH (Mar. 2022), <https://aapor.org/wp-content/uploads/2023/06/Survey-Best-Practices.pdf> [<https://perma.cc/CB2X-DAF2>].

²¹⁸ See Engel & McAdams, *supra* note 216, at 271–88.
²¹⁹ Text generation is tuned to align with human feedback. See, e.g., Stephen Casper, Xander Davies, Claudia Shi, Thomas Krendl Gilbert, Jérémy Scheurer, Javier Rando, Rachel Freedman, Tomasz Korbak, David Lindner, Pedro Freire, Tony Wang, Samuel Marks, Charbel-Raphaël Segerie, Micah Carroll, Andi Peng, Phillip Christoffersen, Mehul Damani, Stewart Slocum, Usman Anwar, Anand Siththaranjan, Max Nadeau, Eric J. Michaud, Jacob Pfau, Dmitrii Krashenninnikov, Xin Chen, Lauro Langosco, Peter Hase, Erdem Bıyık, Anca Dragan, David Krueger, Dorsa Sadigh & Dylan Hadfield-Menell, *Open Problems and Fundamental Limitations of Reinforcement Learning from Human Feedback* (Sep. 11, 2023), <https://arxiv.org/abs/2307.15217> [<https://perma.cc/Q8WM-ZV8W>] (describing Reinforcement Learning from Human Feedback and enumerating many of its limitations).

²²⁰ Using the GPT-3.5-turbo-instruct model and repeating the questions from *Snell* 1000 times yields the distribution of responses shown in the following table. Newsom’s phrasing resulted in GPT answering “yes” twice as often as “no.” The more colloquial phrasing of the question results in GPT saying “no” six times more often than “yes.”

Text of Q1 & Q2	GPT-3.5’s Response to Q2		
	Yes	No	Maybe
Q1: What is the ordinary meaning of “landscaping”?	12.4%	7.8%	79.8%
Q2: Is installing an in-ground trampoline “landscaping”?	(± 2.0)	(± 1.7)	(± 2.5)
Q1: whats landscaping mean?	6.6%	42.4%	51.0%
Q2: is installing a in-ground trampoline landscaping?	(± 1.5)	(± 3.0)	(± 3.1)

The code used to generate this data is available at David B. Stein, *davidbstein*, *snell-question-phrasing.py*, <https://gist.github.com/davidbstein/1a60483a71fdeb2da91f88efe24a55ad> [<https://perma.cc/QS7G-WBD3>].

X's Elon Musk has specifically promoted Grok by claiming that it "loves sarcasm" and would answer questions with "a little humour."²²¹ The demand for configurable general-purpose tools leads to LLM outputs that are highly sensitive to technical configuration.²²² In short, many LLM sensitivities are features, not bugs. The very nature of many LLMs pushes against their ability to be robust interpreters of meaning.

2. External Calibration

Merely passing robustness checks, however, is not sufficient to show that an empirical method is fit for purpose. It may be measuring *something* robustly, but that something might not be the object of interest. Robustness is a necessary condition, not a sufficient one.

In other words, to show empirically that an LLM truly measures linguistic meaning, its outputs must be compared to *an accepted authority on linguistic meaning*. This is the point of Engel and McAdams's study of ChatGPT as a source of ordinary meaning: they compare ChatGPT-3.5 Turbo's responses (under various prompting conditions) to the results of Kevin Tobia's 2020 survey of 2,800 English speakers about the meaning of statutory terms.²²³ Strikingly, three of the four LLM protocols they tested produced results that significantly diverged from Tobia's survey data.

You may or may not agree that surveys are an appropriate measure of the meaning of legal language. Our point is that while it need not be surveys, any attempt to calibrate LLMs empirically depends on having some external benchmark to calibrate against. Without one, no amount of experiments can provide a way to differentiate between "informed conclusion[s] based on a statistical analysis of billions of texts,"²²⁴ and wild extrapolation based on extraneous factors.

Consider Unikowsky's experiments with Supreme Court cases. He found, "Of the 37 merits cases decided so far this Term, Claude decided 27 in the same way the Supreme Court did."²²⁵ One take on this result would be

²²¹ Lucy Hooker, *Musk says his new AI chatbot has 'a little humour'*, BBC NEWS (Nov. 5, 2023), <https://www.bbc.com/news/business-67327060> [<https://perma.cc/ZHA6-5URT>] (quoting Elon Musk).

²²² Sensitivities are also great for people using chatbots to answer factual questions, draft documents, brainstorm new ideas, or build bespoke tools. Users can fiddle with configurations and prompts to achieve desired outcomes.

²²³ See Engel & McAdams, *supra* note 216, at 256–70 (citing Kevin P. Tobia, *Testing Ordinary Meaning*, 134 HARV. L. REV. 726 (2020)); see also Johannes Kruse, *The Ordinary Meaning Bot: Simulating Human Surveys with LLMs* 3 (Max Planck Inst. for Research on Collective Goods, Discussion Paper No. 2025/12, 2025), https://papers.ssrn.com/sol3/papers.cfm?abstract_id=5378203 [<https://perma.cc/SJ2K-PZJX>] (presenting an "approach [that] differs crucially from Engel and McAdams . . . : Rather than asking GPT for aggregate yes-rates, it seeks to approximate the individual response behavior that underlies those aggregates.").

²²⁴ *Generative Interpretation*, *supra* note 2, at 458.

²²⁵ *In AI We Trust II*, *supra* note 5.

that Claude has an accuracy rate of 73%—definitely not good enough for government work. But that is not Unikowsky’s response; instead, he says, “I frequently was more persuaded by Claude’s analysis than the Supreme Court’s.”²²⁶ That may well be (and we will discuss the argument from persuasiveness in the next section), but note what it does to the project of validating Claude’s accuracy. It means that the Supreme Court’s actual decisions are no longer functioning as a source of ground truth. It means that Adam Unikowsky trusts *his own judgment* more than he trusts *the Supreme Court’s opinions* as a benchmark of legal analysis. This too may be right,²²⁷ but notice what it does to the project of validating LLMs for judicial use. It means that you too need to trust Adam Unikowsky’s legal judgment more than the Supreme Court’s. It means we would need a societal consensus that Claude is *better at this than the Supreme Court*. There is no way to bootstrap into such a consensus by counting up how frequently Claude agrees with the Supreme Court. If you think Claude is better at judging than the Supreme Court, you need to argue that Claude is better at judging than the Supreme Court, and back it up with your own standard of what constitutes good judging and evidence that Claude meets that standard.

Other LLM proponents confront this problem, but none of them are able to dispose of it. Newsom relies primarily on intuition, finding that LLMs’ responses “squared with what I had assumed,”²²⁸ and “squared with my own impression—informed by my own experience writing, reading, speaking, and listening to American English in the real world.”²²⁹ What, one wonders, would Newsom have done with their responses if those responses had *not* squared with his own interpretations? For its part, *Generative Interpretation* defines “[a]ccuracy” as “thinking that we really got as close as we could to knowing what the parties would have said.” Arbel and Hoffman acknowledge that “there is no ground truth at hand—we can’t really know what the parties intended at contracting and have to make instead our best guess.”²³⁰

The trouble is that this mode of analysis invites circular reasoning. Arbel and Hoffman define accuracy in terms of “best guess[es]” while simultaneously implying that LLM-produced guesses constitute our “best guess[es].”²³¹ But if LLM-produced guesses are *definitionally* our “best guess[es],” and thus the benchmark for any measurement of accuracy, then *Generative Interpretation*’s case studies do not measure accuracy. Instead, they project authoritative meaning onto LLM outputs, irrespective of what those outputs signify.

²²⁶ *Id.*

²²⁷ As readers of his newsletter, we certainly agree that his arguments are often more convincing than those of the judges and lawyers he critiques!

²²⁸ *United States v. Deleon*, 116 F.4th 1260, 1272 (11th Cir. 2024) (Newsom, J., concurring).

²²⁹ *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1225 (11th Cir. 2024) (Newsom, J., concurring).

²³⁰ *Generative Interpretation*, *supra* note 2, at 459–60, 462.

²³¹ *See id.* at 459–60, 485.

The diversity of the techniques employed for generative interpretation is a warning of the breadth of empirical validation that generative interpretation will demand. Every distinct technique for querying an LLM will require its own distinct, inductive validation. Empirical support that an LLM's word-likelihood estimates reflect linguistic meaning does not necessarily establish that LLMs' probability estimates reflect linguistic meaning. Even if generative-interpretation proponents can empirically validate *one* method of querying LLMs—and, we emphasize, they have not yet done so, although early research gives some reason for optimism²³²—that validation would not establish the accuracy of *other* methods of querying an LLM. Evidence that some LLM output is accurate does not entail that output is accurate because it comes from an LLM.

3. Protocol Specification

Another way of phrasing the difficulty with answering the epistemic challenge empirically is as a problem of specification. The more sensitive generative interpretation is to implementation details, the more comprehensive *ex ante* commitments need to be to achieve predictable or consistent results. Posing a question to an LLM is a wildly underspecified task. As we have seen, outputs can depend on subtle differences in question formulation, model choice, interpretation method, and the configuration and implementation of any software used to operate the LLM.

Proponents of generative interpretation imagine parties agreeing on specific models and prompting methods, then resolving disputes by presenting them to an LLM using those agreed-upon methods.²³³ At least for the case studies explored here, every phrasing and formatting choice seems to influence outcomes. Constructing a sufficiently detailed pre-commitment for predictable generative interpretation requires agreeing on outcome-altering details. For all but the easiest of cases, many of those details are dispute-specific.²³⁴ Why not negotiate the anticipated disputes directly? Put another way, picking one experimental setup over another merely displaces the arbitrariness from the empirical realm (how to reconcile the different results of protocol *A* and protocol *B*?) to the epistemic one (what makes protocol *A* better or worse than protocol *B*?)

Again, to repeat, our claim is not that this *cannot be done* in a principled way, just that it *has not yet been done* in a principled way, and that it *will have to be done* in a principled way to justify generative interpretation empirically. The empirical woods are lonely, dark, and deep—and generative interpretation has miles to go.

²³² See Engel & McAdams, *supra* note 216, at 43.

²³³ See, e.g., *Generative Interpretation*, *supra* note 2, at 501.

²³⁴ But see *id.* at 510 (“[I]f courts follow our proposed best practices, this method is also predictable *ex ante*.”).

4. No Easy Fixes

Rather than waiting for more robust validation, it is tempting to play whack-a-mole with problems like the ones we spotted in Part III. When an LLM generates text that influences its later answers—like in *Snell* and *Famiglio*—some might naturally respond by constraining the model to one-word responses.²³⁵ But that merely shoves the improvident salience into the prompt,²³⁶ and some research suggests that LLMs produce lower-quality answers when forced to be succinct.²³⁷ So maybe we cross-reference between models instead.²³⁸ But model quality can vary; models can (and do) converge on arbitrary answers.²³⁹ So maybe we generate multiple versions of the same prompt.²⁴⁰ But sampling across low-quality prompts can nudge models towards the same irrelevant details.²⁴¹ And so on. Each problem exposes or introduces another, with some solutions possibly decreasing answer quality. Without a reliable way to measure performance, we cannot know which changes improve generative interpretation and which choices exacerbate or obscure problems.

It is also tempting to over-extrapolate from simplified explanations of complex systems. It may feel natural to assume that because LLMs are trained on a large amount of text, their answers are informed by that text’s semantic content. But that inference is unsupported; semantic information is one (potentially inconsequential) part of a more complex whole. If the argument that LLMs are trustworthy because of how they work does not go through on its own (as we argued in section IV.A), it also does not establish that LLMs are empirically trustworthy (as we discuss in this section). That claim must

²³⁵ See *Generative Interpretation*, *supra* note 2, at 485 n.164 (proposing this approach).

²³⁶ See Felipe Maia Polo, Ronald Xu, Lucas Weber, Mirian Silva, Onkar Bhardwaj, Leshem Choshen, Allysson Flavio Melo de Oliveira, Yuekai Sun & Mikhail Yurochkin, *Efficient Multi-Prompt Evaluation of LLMs*, 38 CONF. NEURAL INFO. PROCESSING SYS. 1, (2024) (collecting sources).

²³⁷ Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed H. Chi, Quoc V. Le & Denny Zhou, *Chain-of-Thought Prompting Elicits Reasoning in Large Language Models*, 36 CONF. NEURAL INFO. PROCESSING SYS. 1, 6 (2022).

²³⁸ See *Generative Interpretation*, *supra* note 2, at 503–04 (advising use of this approach as a “best practice”).

²³⁹ See *supra* note 151 and accompanying text.

²⁴⁰ See *Generative Interpretation*, *supra* note 2, at 490 (“tr[ying] something new” by generating multiple prompts); see also Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel & Pontus Stenetorp, *Fantastically Ordered Prompts and Where to Find Them: Overcoming Few-Shot Prompt Order Sensitivity*, 60 PROC. ASSOC. COMPUTATIONAL LINGUISTICS 8086, 8089 (2022) (generating multiple prompts automatically when querying a model); Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi & Graham Neubig, *Pre-Train, Prompt, and Predict: A Systematic Survey of Prompting Methods in Natural Language Processing*, 55 ACM COMPUTING SURVS., Jan. 2023, at 8 (collecting and reviewing literature on “prompt paraphrasing” and “prompt generation”).

²⁴¹ See Albert Webson & Ellie Pavlick, *Do Prompt-Based Models Really Understand the Meaning of Their Prompts?*, 2022 PROC. CONF. N. AM. CHAP. ASS’N FOR COMPUTATIONAL LINGUISTICS: HUM. LANGUAGE TECHS. 2300, 2307.

be made, and supported, on its own merits. An LLM might use great data in an unhelpful way. The brittleness and arbitrariness of early experiments in generative interpretation should caution us that the method is not yet “good enough” for use in litigation.²⁴² Indeed, numerous empirical studies have documented that LLMs consistently make basic errors in legal reasoning.²⁴³

C. *Relying on Persuasiveness?*

The third major way that generative-interpretation proponents attempt to justify their methodology is by arguing that the results literally speak for themselves: that LLM-generated text is sufficiently *persuasive* that it should be treated as authoritative. This type of argument is particularly slippery, because it comes in two closely related forms that make very different assumptions. One is that the persuasiveness of LLM-generated text is *evidence* that it resulted from a reliable process of linguistic interpretation or legal reasoning. The other is that persuasive LLM-generated text is authoritative purely *because it is persuasive*, and that it is irrelevant how it was generated. To clarify the difference, and to understand why both forms of argument are deceptive, it is necessary to speak first of abduction.

1. *Generative Interpretation as Abduction*

Consider the claim, “This text looks like it was the result of a skilled process of legal reasoning, so therefore it was produced by one.” The name for this form of reasoning is *abduction*, or sometimes “inference to the best explanation.”²⁴⁴ Unlike *deduction*, which reasons from premises to their logical entailments, or *induction*, which reasons to find generalizations consistent with existing data, abduction reasons by identifying plausible causes for observed phenomena. “All the balls in that urn are black; these balls came from that urn; therefore these balls are black” is deduction; “These balls are black; these balls came from that urn; therefore all balls in that urn are black” is induction; “These balls are black; all the balls in that urn are black; therefore these balls came from that urn” is abduction.

Though potentially less robust than inductive or deductive proof, abductive reasoning is often the only option in the law. For example, the preponderance of the evidence test asks for the most reasonable or coherent interpretation of the evidence. The legitimacy of legal dispositions is also something we may infer abductively. As later subparts explore more deeply, a judge’s decision is

²⁴² *Contra Generative Interpretation*, *supra* note 2, at 458.

²⁴³ See, e.g., Varun Magesh, Faiz Surani, Matthew Dahl, Mirac Suzgun, Christopher D. Manning & Daniel E. Ho, *Hallucination-Free? Assessing the Reliability of Leading AI Legal Research Tools*, 22 J. EMPIRICAL LEGAL STUD. 216, 217 (2025).

²⁴⁴ See generally Igor Douven, *Abduction*, STAN. ENCYC. PHIL. (June 18, 2025), <https://plato.stanford.edu/entries/abduction/> [<https://perma.cc/6HHU-CGKY>].

assumed to be the product of a good-faith application of legal reasoning because her written decision describes that reasoning process. When a judicial opinion sets forth a persuasive chain of legal reasoning that accurately refers to relevant law, we assume abductively that it was in fact realized through the analytical process it describes, and hence that the opinion is legitimate.

Abduction is not in itself a wrongheaded analytical device. We rely on abduction all the time. When a student writes an exam that makes cogent, well-informed arguments about the assigned materials, we treat that document as evidence that the student has mastered the topic. We presume that producing such an exam response requires a sophisticated understanding of the course materials, and we regard the exam as proof that its author has such an understanding. Of course, our abductive inference might be wrong. Maybe the exam response was produced by a cat walking across a keyboard, or a fountain pen leaking into a bluebook in exactly the pattern of handwritten text. But, based on our experience, we conclude that such situations are implausible enough that we award the student a good grade unless we have evidence to believe that her exam response was not her own composition.

Abduction, however, is limited by what we know and by the accuracy of what we think we know. In 2025, for example, it is probably educational malpractice for an instructor to be unaware that students have access to LLMs, because the wide availability of LLMs weakens the abductive inferences that educators customarily make after inspecting students' work. LLMs allow students to generate exactly the sorts of prose that, previously, would have required comprehension of the course materials. A cogent term paper no longer signifies what it signified in the pre-AI age. To ensure that written documents continue to serve as rough measurements of students' knowledge, instructors must redesign examinations to ensure that they measure what they purport to measure. At least two of us, for example, have forgone take-home examinations in favor of proctored, in-class assessments, to ensure that students' written work supports the same inferences about their mastery of course materials as it did before the advent of LLMs.

Just as they complicate the inferences that a teacher can make from a well-written essay, LLMs complicate the inferences that a jurist can make from a well-written legal argument. It used to be fair to infer that someone who authored a cogent legal opinion had considered the relevant issues. But as the processes for generating legal writing change, so do the permissible inferences that a legal work product can engender.

Abduction clarifies the nature of the two kinds of arguments from persuasiveness. A claim that LLM-generated text is authoritative because its persuasiveness reassures us that it was generated through a valid process of legal reasoning is an abductive argument. It uses the superficial qualities of LLM-generated text to make a claim about the legitimacy of the underlying reasoning. That claim, however, is defeasible, and can be rebutted by a showing that there are good reasons to doubt the bona fides of the process that generated the text. A claim that LLM-generated text is authoritative simply because it

is persuasive is not abductive. It avoids the rebuttal by treating the generation process as irrelevant. But, as we will see, it is a claim that is dramatically at odds with the theory of adjudication on which our legal system currently rests.

2. *Rhetoric ≠ Adjudication*

LLMs are rhetoric machines. Rhetoric is the art of persuasive argumentation,²⁴⁵ and it is unsurprising that LLMs trained on an enormous corpus of humans doing their best to make persuasive arguments are capable of replicating the textual features of those arguments that make them persuasive. Indeed, the training process (particularly the stages that involve human feedback) may select for persuasiveness, thereby amplifying the connections in an LLM's internal model of language that contribute to fluent, persuasive outputs.

Rhetoric is successful precisely to the extent it is persuasive.²⁴⁶ As anyone who reads their output can tell, LLMs excel at rhetorical writing. Today's LLMs reproduce the form of high-quality legal reasoning. As Unikowsky and others have observed, LLM output now replicates rational, eloquent argumentation that applies precedent to novel facts.²⁴⁷ An LLM can produce text that may be formally indistinguishable from—or even formally superior to—the reasoning described by an opinion written by a human judge.

There are many ways to elicit responses from LLMs that sound plausibly correct. But as anyone who has dealt with lawyers well knows, sounding correct is not necessarily the same as being correct. Indeed, there is ample evidence from the use of LLMs in law to illustrate the difference: ask any lawyer who has been sanctioned for submitting a brief containing LLM-generated “hallucinations” whether the fluency of LLM-generated text guarantees that the legal claims it makes are sound.²⁴⁸ Though tempting, conflating persuasiveness and coherence with authority or accuracy is a category error.²⁴⁹

²⁴⁵ See, e.g., SHARON CROWLEY & DEBRA HAWHEE, *ANCIENT RHETORICS FOR CONTEMPORARY STUDENTS* 9 (1994).

²⁴⁶ Note that we are not using “rhetoric” to refer to the concept of “constructive rhetoric,” as it sometimes does in writing on legal philosophy. See, e.g., James Boyd White, *Law as Rhetoric, Rhetoric as Law: The Arts of Cultural and Communal Life*, 52 U. CHI. L. REV. 684, 701 (1985) (using a “highly expanded” definition for “rhetoric”). That definition is useful in interrogating the cultural and social relationship between rhetorical legal artifacts—like statutes—and their authors and audience, but assumes a human author and is therefore ill-suited for this context. See *id.* at 689–90 (defining constructive rhetoric in reference to the lawyers creating it).

²⁴⁷ See *In AI We Trust II*, *supra* note 5.

²⁴⁸ See, e.g., *Mata v. Avianca, Inc.*, 678 F. Supp. 3d 443, 448–49 (S.D.N.Y. 2023) (imposing sanctions on attorneys who “submitted non-existent judicial opinions with fake quotes and citations created by the artificial intelligence tool ChatGPT”).

²⁴⁹ Cf. Michael Townsen Hicks, James Humphries & Joe Slater, *ChatGPT Is Bullshit*, 26 ETHICS INF. TECH., at 3 (June 2024) (“The problem here isn’t that large language models hallucinate, lie, or misrepresent the world in some way. It’s that they are not designed to represent the world at all; instead, they are designed to convey convincing lines of text.”).

The problem with rhetoric is that sometimes an argument can be persuasive for the wrong reasons.²⁵⁰ Rhetoric encompasses arguments that comply with all the formal conventions of good-faith legal reasoning. It also, however, encompasses writing that is persuasive for other reasons, like well-concealed gaps in logic and appeals to emotion. Rhetoric covers everything from “legal reasoning” to “casuistry,” and the rhetorician’s good or bad faith is irrelevant. LLMs make newly salient the ancient anxiety about rhetoric itself²⁵¹—the danger that rhetorical skill can persuade the audience, regardless of whether the argument itself is right or wrong.

By contrast, we expect more from adjudication. To be sure, a good adjudicator will present her decisions in a rhetorically effective format. But adjudications are not successful simply because they reach persuasive conclusions. Rather, adjudication requires the integrity of the process of legal reasoning that realized that conclusion.²⁵²

That is, adjudication is also subject to procedural criteria. At minimum, these criteria require that a decision-maker make a good-faith, rational effort to produce a legal conclusion. If a formally adequate legal determination is produced in bad faith, or without rational effort, it does not adhere to the process of legal reasoning. A decision-maker’s good faith establishes that she has not elaborated a formally adequate conclusion for improper, extralegal reasons. For example, a conflict of interest will taint a decision-maker’s conclusion of law. Even if the conclusion is set forth in a formally adequate decision, it is not a legitimate legal conclusion because an observer reasonably can suspect that the judge’s personal interests influenced her determination of the appropriate outcome.

Procedural constraints also require a decision-maker’s conclusion to be at least minimally rational. Paul the Octopus was a cephalopod whose handlers would serve him food in boxes that bore the logos of competing soccer teams.²⁵³ When Paul chose to eat from one box before another, his handlers deemed him to have “predicted” that the corresponding team would win.²⁵⁴ When this predictive method accurately foretold the winners of a number of international soccer matches, Paul became world famous.²⁵⁵ If a judge issued a formally adequate legal opinion, and it were later revealed that she had

²⁵⁰ Consider, by way of contrast, the use of AI for victim impact statements. *See, e.g.*, Juliana Kim, *Family shows AI video of slain victim as an impact statement—possibly a legal first*, NPR (May 12, 2025), <https://www.npr.org/2025/05/07/g-s1-64640/ai-impact-statement-murder-victim> [<https://perma.cc/JD7P-Y6F5>]. Even here, in a domain where pathos is as valued as logos, such uses raise the fear that AI will be persuasive for the wrong reasons.

²⁵¹ *See, e.g.*, PLATO, *GORGIAS* 27 (Donald. J. Zeyl, trans., Hackett Publishing 1987).

²⁵² *See* Chad M. Oldfather, *Writing, Cognition, and the Nature of the Judicial Function*, 96 GEO. L.J. 1283, 1318 (2008) (discussing the widely-shared belief that judicial “writing provides an important discipline on thought”).

²⁵³ *Paul the Octopus*, WIKIPEDIA, https://en.wikipedia.org/wiki/Paul_the_Octopus [<https://perma.cc/2765-B2VT>].

²⁵⁴ *Id.*

²⁵⁵ *Id.*

chosen the winning party because her pet octopus ate from a box labeled with that party's name, this procedural deficiency would undermine the decision's status as legal reasoning.

A formally adequate legal conclusion is legal reasoning's end product, but is not in itself proof that legal reasoning has taken place. This is because a formally adequate legal work product can be realized through processes that do not constitute legal reasoning. In some intellectual traditions, the integrity of the form is the integrity of the process. If, for example, a mathematician authored a formally accurate proof and later claimed that it came to him as a divine message, the proof would still be mathematically legitimate.²⁵⁶ If, on the other hand, a judge issued a formally adequate opinion and later revealed that she had reached the conclusion that a divine message had commanded her to adopt, such a revelation would undermine her claim to have engaged in legal reasoning.²⁵⁷ In an essay written before her judicial appointment, Justice Amy Coney Barrett concludes that Catholic doctrine forbids observant judges from enforcing the death penalty, and that a moral refusal to consider a death sentence warrants recusal from the sentencing phase of a capital case.²⁵⁸ Even if the judge can justify her conclusion in a formally adequate legal document, her recusal is warranted because she cannot reach that conclusion through a process recognized as legal reasoning.

The form and process of legal reasoning are socially contingent. Even a single legal tradition admits variations: modern American jurists disagree about the precise form that legal reasoning may take. Some, for example, "object to the use of legislative history on principle,"²⁵⁹ while others encourage it.²⁶⁰ A legislative-history skeptic may find an opinion by a legislative-history enthusiast formally inadequate on the ground that it employs a categorically impermissible mode of argument. The legislative-history skeptic may find this legal reasoning just as illegitimate as legal reasoning that invokes the judge's dreams or her subjective "value preferences" to justify a conclusion.²⁶¹

²⁵⁶ The esteemed mathematician Srinivasa Ramanujan claimed scholarly inspiration from a family goddess, whose consort visited him in a dream. See Roshni Chakrabarty, *Srinivasa Ramanujan: The mathematical genius who credited his 3900 formulae to visions from Goddess Mahalakshmi*, INDIA TODAY (Dec. 23, 2022), <https://www.indiatoday.in/education-today/gk-current-affairs/story/srinivasa-ramanujan-life-story-973662-2017-04-26> [https://perma.cc/95H6-SMVZ].

²⁵⁷ Cf. Lon L. Fuller, *The Forms and Limits of Adjudication*, 92 HARV. L. REV. 353, 390–91 (1978) (arguing that a decision rendered on "grounds . . . largely unrelated to what occurred at the hearing . . . fall[s] beyond the limits of adjudication").

²⁵⁸ Amy C. Barrett & John H. Garvey, *Catholic Judges in Capital Cases*, 81 MARQ. L. REV. 303, 305, 335, 339 (1998).

²⁵⁹ Antonin Scalia, *Common-Law Courts in a Civil-Law System: The Role of United States Federal Courts in Interpreting the Constitution and Laws*, in A MATTER OF INTERPRETATION 3, 31 (Amy Gutmann ed., 2018).

²⁶⁰ See generally Stephen Breyer, *On the Uses of Legislative History in Interpreting Statutes*, 65 S. CAL. L. REV. 845 (1992).

²⁶¹ Cf. Robert H. Bork, *Neutral Principles and Some First Amendment Problems*, 47 IND. L.J. 1, 9 (1971) (arguing that when the Court departs from the text and history of the Constitution,

But granular disagreements about particular formal qualities can reveal the interlocutors' broader agreement about the general form that legal argumentation should take. Some arguments for constitutional originalism, for example, posit that the methodology is, in consequentialist terms, the best approach to constitutional interpretation.²⁶² Consequentialist arguments for originalism do not assert that originalism is the only true form of legal reasoning.²⁶³ Rather, these arguments suggest a broader view of legal argumentation—as, say, consequentialist arguments rationally propounded—and employ that argumentative paradigm to advocate for originalism.

Similarly, that a document formally resembles legal reasoning does not make it adjudication, nor does this resemblance necessarily make the document authoritative. Whether something constitutes adjudication, and whether it is authoritative, depends on social facts like Hart's rule of recognition.²⁶⁴ The authority of an adjudication depends not just on formal resemblance to a reasoned judicial opinion, but on a set of social practices that imbue particular artifacts of legal reasoning with authority. Although the formal adequacy of written legal reasoning produced by humans may have been a decent proxy for presuming that it was realized in accordance with a bona fide adjudicative process, the same can't be said for formally adequate legal reasoning generated using an LLM. To presume that these AI-generated simulacra of rationality constitute the social process of adjudication is to conflate form and process.

It is almost trivially easy to make an LLM produce any kind of output desired (e.g., "Write a textualist opinion finding for the defendant while making as many references as possible to the films of Steven Spielberg"). Without some other kind of reassurance that the results reflect the kind of reasoning we associate with competent and good-faith judging—the kind that could be provided in theory by inductive or deductive approaches—saying that LLM outputs are valid interpretations because they look good boils down to a vibe-check. Ironically, it is precisely the sort of "visceral, gut-instinct decisionmaking" that proponents of generative interpretation hope to avoid.²⁶⁵

it is left solely with unprincipled and illegitimate "value preferences" to decide constitutional cases).

²⁶² See Cass R. Sunstein, *Originalism*, 93 NOTRE DAME L. REV. 1671, 1680–81 (2018); William Baude, *Is Originalism Our Law*, 115 COLUM. L. REV. 2349, 2351 (2015).

²⁶³ Of course, some proponents of constitutional originalism do posit that it is the only legitimate method of constitutional interpretation, and this view comes closer to asserting that non-originalist reasoning is in fact not legal reasoning. ANTONIN SCALIA & BRYAN A. GARNER, *READING LAW* 89 (2012) ("[Originalism] is the only objective standard of interpretation . . .").

²⁶⁴ H.L.A. HART, *THE CONCEPT OF LAW* 92–93 (1961) ("[A] 'rule of recognition' . . . will specify some feature or features possession of which by a suggested rule is taken as a conclusive affirmative indication that it is a rule of the group to be supported by the social pressure it exerts.").

²⁶⁵ *Snell v. United Specialty Ins. Co.*, 102 F.4th 1208, 1224 (11th Cir. 2024) (Newsom, J., concurring); see *Generative Interpretation*, *supra* note 2, at 459 (describing how their case studies applying generative interpretation "illuminate how transparent and objective interpretative methodologies have advantages over intuitive ones").

Ideally, some sort of assessment could separate the know-nothing scanners and tortoises from authoritative intelligences—some attribute of LLMs that, if observed, would support the conclusion that its interpretations are accurate and meaningful. Unikowsky acknowledges that before delegating legal decision-making to LLMs, we want to “mak[e] sure the AI works.”²⁶⁶ “[T]his,” he says, “shouldn’t be that hard to do”; validating AI’s decision-making would require only “download[ing] a few hundred briefs off of PACER, input[ting] them into your AI of choice, ask[ing] the AI to adjudicate the case, and compar[ing] the results to the actual judicial decisions that came down in those cases.”²⁶⁷

But does Unikowsky’s proposed test really measure accuracy? Unikowsky’s writings on LLMs conflate at least three different desiderata. First is LLMs’ ability to get the “right answer”—that is, to reach the same disposition as a court on particular briefing. Second is LLMs’ ability to sound reasonable. Third is LLMs’ ability to be reasonable—that is, to “reason.” These are distinct capabilities. Possessing either of the first two does not entail possessing any of the others. For example, “getting the right answer” alone may not prove much: we wouldn’t delegate adjudication to a dog, even if it demonstrated a robust and uncanny ability to “predict” case outcomes in conformity with litigated outcomes by choosing to eat from a particular labeled bowl. Simply sounding reasonable isn’t sufficient either. An LLM, like a skilled lawyer, can produce reasonable-sounding justifications for incompatible legal outcomes. Generating a reasonable-sounding justification does not in itself demonstrate the reasoning we expect from an adjudicator.

3. *Why Text Matters—And Doesn’t*

It is important to understand the work that text in particular does in legal reasoning—and why LLM-generated text does not automatically do that work. The rule of law is not merely a collection of written artifacts; rather, it is the social process that produces those artifacts. In the contemporary United States, those artifacts happen to look like legal opinions. The texts are the visible traces of the system in action, but they are not the system itself.

It’s unsurprising that a lawyer or a judge might regard written legal reasoning as the paramount artifact of the American legal system. For centuries, the legal profession has produced written judicial opinions. These authoritative opinions form part of a process that constitutes the rule of law. Lawyers, laypersons, and members of other branches of government could read these decisions and coordinate future action based on the reasoning they set forth. Their holdings bind private citizens and government officials alike.

²⁶⁶ *In AI We Trust I*, *supra* note 5.

²⁶⁷ *Id.*

Now, LLMs can produce artifacts that are formally indistinguishable from the artifacts produced by courts that apply the law authoritatively.

Notwithstanding the pride of place it enjoys today, written adjudication is not essential to an authoritative legal process. Juries, for example, are never called upon to explain their findings rationally, even though they serve a vital legitimating function in the criminal law.²⁶⁸ A typical verdict has the barest formal qualities: it is an unelaborated yes or no. It is not the form of a verdict that primarily inspires faith in the jury. Rather, it is our understanding of the jury's deliberative process that legitimates its authority.

Comparing juries to judges illustrates that law is process in addition to formal output. Unlike the work product of, say, an appellate court, the output that a jury produces is trivial. Providing a verdict is, of course, a fundamental part of the jury's job. But the point of the jury system isn't just to manufacture yes/no verdicts. Human juries are resource-intensive; rendering verdicts by a coin flip or a random binary generator would be far more economical, offering great efficiencies and improving access to the courts. Yet we routinely reject innovations that would make verdicts cheaper and more efficient to produce.²⁶⁹ We do so to protect the integrity of the process that realizes the verdict. In other words, there is very little about the form of a verdict that makes it authoritative and legal. What gives a verdict its legal authority is almost entirely our beliefs about a jury's process.

LLM proponents assert that AI is authoritative not simply because of how it sounds, but because of how it works. This reasoning suggests that we ought to trust the conclusions that AI reaches even when they are unaccompanied by output that resembles legal reasoning.²⁷⁰ And indeed, this is precisely what practitioners of generative interpretation do when they find that

²⁶⁸ See, e.g., *Apprendi v. New Jersey*, 530 U.S. 466, 497 (2000) (referring to “the jury tradition that is an indispensable part of our criminal justice system”); see also *id.* at 477 (quoting Justice Story and William Blackstone to corroborate the importance of “trial by jury” “to guard against a spirit of oppression and tyranny”).

²⁶⁹ See, e.g., *Granfinanciera, S.A. v. Nordberg*, 492 U.S. 33, 63 (1989) (“It may be that providing jury trials in some fraudulent conveyance actions . . . would impede swift resolution of bankruptcy proceedings But ‘these considerations are insufficient to overcome the clear command of the Seventh Amendment.’”); *Apprendi*, 530 U.S. at 498 (Scalia, J., concurring) (“[T]he jury-trial guarantee . . . has never been efficient; but it has always been free.”); *Parklane Hosiery Co. v. Shore*, 439 U.S. 322, 343 (1979) (Rehnquist, J., dissenting) (“[T]he concerns for the institution of jury trial that led to the passages of the Declaration of Independence and to the Seventh Amendment were not animated by a belief that use of juries would lead to more efficient judicial administration.”).

²⁷⁰ Note that when an LLM produces an output that describes a reasoning process, this description does not necessarily correspond to the process that the LLM actually followed to generate that output. See, e.g., Yanda Chen, Joe Benton, Ansh Radhakrishnan, Jonathan Uesato, Carson Denison, John Schulman, Arushi Somani, Peters Hase, Misha Wagner, Fabien Roger, Vlad Mikulik, Samuel R. Bowman, Jan Leike, Jared Kaplan & Ethan Perez, *Reasoning Models Don't Always Say What They Think* (May 8, 2025), <https://arxiv.org/abs/2505.05410> [<https://perma.cc/2DKD-U9ZT>] (finding that LLMs prompted to explain their reasoning via a “chain of thought” frequently drew conclusions based on “hints” that they did not disclose in their supporting reasoning).

individual words and numbers generated by LLMs are probative to their understanding of a legal instrument or a particular term (as Arbel and Hoffman do in their *Famiglio* and “flood” case studies discussed above, for example).²⁷¹ If our faith in AI should persist even when it does not present its conclusions in rational-sounding language, then shouldn’t that faith justify replacing not just lawyers and judges with AI, but juries, too?

Many readers may find the prospect of replacing juries with AI unpalatable, particularly in criminal cases. Many readers probably suspect that the public will not regard AI juries as having the same sort of legitimacy and authority as a human jury. The intuition that the public will reject their authority is a very good reason to resist AI juries. It also underscores that legality is a social fact; it does not follow inexorably from form. Just because AI can render something formally indistinguishable from a verdict does not mean that it can render a legally acceptable verdict. Correspondingly, just because LLMs can output text that resembles well-reasoned adjudication, does not mean that it can engage in the social process that constitutes acceptable adjudication.

To put the point another way, juries are not important because they render binary verdicts; countless mechanisms can output a “1” or a “0.” They are important because they instantiate a legitimate and authoritative process for reaching a binary verdict. If we cared only about the formal artifacts of juries—that is, only the verdicts they produce—then we could delegate jury duty to a coin flip. The thought of deciding criminal culpability by coin flip appalls us because we value juries for the process they instantiate, not simply the artifacts they produce. The same is true of adjudication—but focusing on the formal qualities of an opinion’s text obscures this point. The temptation is to treat formally adequate legal text as evidence of the process we call “legal reasoning” or “adjudication.” But formally adequate legal text is not necessarily adjudication—just like a bare yes or no resembles a verdict but is not necessarily one.

Consider one of Unikowsky’s proposals. He observes that submitting cases nationwide to a single AI adjudicator could reduce the unpredictability that arises in a judicial system that vests power in thousands of individual judges.²⁷² However, he writes,

intellectual diversity across the judiciary is healthy, at least up to a point, and . . . it also seems bad to concentrate all judicial power in one (or a small number of) computer systems. But this problem can be dealt with by, e.g., adding a random ideological factor to the AI’s outputs, having the AI highlight when a particular dispute would come out differently across different philosophies, using multiple independent AIs, or through other mechanisms. The point is, if we use AI, we can add exactly as much unpredictability and ideology as

²⁷¹ See *Generative Interpretation*, *supra* note 2, at 455–58, 483–85.

²⁷² See *In AI We Trust I*, *supra* note 5.

we want via effective prompt engineering rather than having unpredictability thrust upon us by the constraint of individual judges' bandwidth.²⁷³

Why is it that adding a “random ideological factor”—a fancy version of a coin flip—would make AI adjudication more, rather than less, legitimate? Or rather, why does it seem to Unikowsky that taking a system that is arguably inferior to human judging (an AI adjudicator) and combining it with a system that is definitely inferior to human judging (random chance) will fix its deficiencies?

Unikowsky's proposal to inject arbitrary amounts of “unpredictability and ideology” misses what is important to the process of adjudication. Outputs do not have an “ideological factor”; judges do. Intellectual diversity in the judiciary is not desirable because it guarantees that judges will resolve cases differently. (It guarantees no such thing.) What intellectual diversity does guarantee is that judges will approach cases differently. And it is a diversity of approaches that is valuable, not a diversity of dispositions for its own sake.²⁷⁴ In fact, it tends to inspire confidence in the judiciary when judges of different ideological persuasions agree on the disposition of a case.²⁷⁵

Unikowsky's other suggestions come closer. “[H]aving the AI highlight when a particular dispute would come out differently across different philosophies” explicitly puts the diversity of judicial approaches into a single LLM's considerations; “using multiple independent AIs” attempts to have a diversity of LLMs emulate a diversity of judges.²⁷⁶ Explicitly spelling out these considerations, however, raises exactly the same questions we have been discussing—what kinds of validation would it take to make us confident that the process of LLM adjudication sufficiently captures the process values of human adjudication?

So far, this discussion has centered on the artifacts produced by legal reasoning. But generative-interpretation moderates do not propose replacing judges entirely. Instead, they suggest merely using LLMs as an alternative to common sense or a dictionary.

²⁷³ *Id.*

²⁷⁴ See generally David Orentlicher, *Politics and the Supreme Court: The Need for Ideological Balance*, 79 U. PITT. L. REV. 411 (2018) (arguing that an ideological balance of justices on the Supreme Court will result in better decisions); Deseriee A. Kennedy, *Judicial Review and Diversity*, 71 TENN. L. REV. 287 (2004) (arguing that increasing judicial diversity is necessary to bring additional perspectives into decisions and ensure fair review of cases and issues relevant to people of color).

²⁷⁵ See, e.g., Devin Dwyer, *Supreme Court Defies Critics with Wave of Unanimous Decisions*, ABC NEWS (June 29, 2021), <https://abcnews.go.com/Politics/supreme-court-defies-critics-wave-unanimous-decisions/story?id=78463255> [<https://perma.cc/6TRW-RXBK>] (quoting Kate Shaw: “I suspect the justices feel there is value in conveying to the American people that in a hyper-partisan moment, . . . the court remains a largely nonpartisan institution.”).

²⁷⁶ See *In AI We Trust I*, *supra* note 5.

There is a procedural oddity to judicial use of dictionaries: looking something up in a dictionary is in some respects a factual inquiry, yet one that judges perform *sua sponte*, outside of the ordinary adversarial process. The legitimacy of judicial use of external sources of fact during the interpretive process relies *even more heavily* on social norms and the judge's legal-reasoning process than does the legitimacy of using facts that have passed through the gauntlet of a trial. Judge Newsom complains that judges "rarely explain in any detail the process by which they selected one definition over others."²⁷⁷ Yet even without explanation, *the selection was a product of the judge's reasoning*.

The counterintuitive upshot is that interpretation is one of the *last* places where LLM use is justified within the adjudicatory process. When juries produce a "yes" or "no" response, they provide no explanation. The fact that the jury produced the answer is the source of its legitimacy. In short, when explanations run thin and discretion peaks, procedural and institutional sources of legitimacy acquire special import. Textual interpretation is an adjudicative task where judges' social role has special importance. LLM proponents suggest that computer programs might provide objective answers to subjective questions. Unilaterally abdicating the interpretive role in favor of an LLM (especially without substantial supporting evidence for those claims) actively subverts an important judicial responsibility.

4. *The Risk of Projection*

We make these points at length because of the danger that LLM users might unreflectingly project authority on to LLM outputs simply because the outputs sound right.²⁷⁸ Arbel and Hoffman mention the concern explicitly.²⁷⁹ And in both *Snell* and *Deleon*, Judge Newsom observed, "I definitely didn't want to fall into the trap of embracing ChatGPT's definition just because it aligned with my priors."²⁸⁰ Yet Judge Newsom wrote that his fears were "reassur[ed]" after he employed Arbel and Hoffman's proposed technique for measuring LLMs' "confidence" and found that he could produce similar answers by posing the same interpretive question multiple times and to multiple LLMs.²⁸¹

²⁷⁷ *Snell v. United Specialty Ins. Co.*, 102 F.4th at 1208, 1229 (11th Cir. 2024) (Newsom, J., concurring).

²⁷⁸ See Joe McIntyre, *Pareidolic Illusions of Meaning: ChatGPT, Pseudolaw and the Triumph of Form over Substance* 39–45 (Mar. 16, 2025), <https://papers.ssrn.com/abstract=5181165> [<https://perma.cc/FMS6-2JGP>].

²⁷⁹ See *Generative Interpretation*, *supra* note 2, at 476 ("[J]ust because the probabilities are reasonable doesn't mean they are accurate.").

²⁸⁰ *Snell*, 102 F.4th at 1225 (Newsom, J., concurring); see also *United States v. Deleon*, 116 F.4th 1260, 1272 (11th Cir. 2024) (Newsom, J., concurring).

²⁸¹ *Deleon*, 116 F.4th at 1274 (Newsom, J., concurring) ("[R]eassuringly, the 30 results I received—10 apiece from each of the three leading LLMs—largely echoed the initial response that I got from ChatGPT.") (citing *Generative Interpretation*, *supra* note 2 at 487–88); see also

Judge Newsom is right to identify the risk of projection, but the measures he adopts do not guard against it. His confidence checks do not demonstrate that he is not vesting LLMs with authority simply because their output matches his expectations. Rather, they demonstrate only that he is more comfortable vesting LLMs with authority when he thinks they are confidently and consistently producing output that matches his expectations. Indeed, he said as much: “If the model . . . returned essentially consistent answers—responses that coalesced around a common core—then one could probably say, with a higher degree of confidence, that the model was getting at the ‘correct’ response.”²⁸² This framing assumes that confidence and consistency corroborate correctness. They do not.²⁸³ The confidence and consistency of a model’s output bolster its probative value only when there is reason to believe that the model’s determinations accurately reflect facts about linguistic meaning in the first place. And this latter belief is something generative-interpretation proponents justify, at least in part, with appeals to the superficial reasonableness and plausibility of LLM output.²⁸⁴ So long as they do so, the danger remains that users will project authority onto LLMs not because LLMs are objectively “accurate,” but because they like what they hear.

D. Relying on Predictability?

The final general type of argument offered by LLM proponents is that they are more *predictable* than human judges.²⁸⁵ The importance of predictability to the rule of law is well established, and we need not dwell on why predictability is a worthwhile goal. Decades of legal scholarship highlight arbitrariness in current interpretative techniques.²⁸⁶ If anything, predictability looms even larger for LLM proponents than it does for other judges and scholars. Kieffaber, for example, argues that other goals for interpretation are “really predictability in disguise.”²⁸⁷ If LLMs provide consistent answers, they might improve judicial consistency and offer parties a reliable way to determine the meaning of terms without litigating. LLMs that sometimes produce arbitrary results might still improve upon the status quo if those results are predictable.

Snell, 102 F.4th at 1225, 1233 (Newsom, J., concurring) (proposing that users “try different prompts” and “query multiple models to ensure that the results are consistent”).

²⁸² *Deleon*, 116 F.4th at 1274 (Newsom, J., concurring).

²⁸³ See McIntyre, *supra* note 278, at 45–46.

²⁸⁴ See, e.g., *Snell*, 102 F.4th at 1228 (Newsom, J., concurring) (“And as anyone who has used them can attest, modern LLMs’ results are often sensible . . . [LLMs] are[] high-octane language-prediction machines capable of probabilistically mapping, among other things, *how ordinary people use words and phrases in context*.” (emphasis added)).

²⁸⁵ See *Generative Interpretation*, *supra* note 2, at 509.

²⁸⁶ See *id.* at 464–66 (summarizing and collecting relevant literature); see also *In AI We Trust II*, *supra* note 5.

²⁸⁷ Kieffaber, *supra* note 116, at 12.

LLM proponents offer, broadly, four reasons to think that LLMs are more predictable. First, there are claims that LLMs will be more predictable because they are algorithmic and automated, and thus not subject to human biases. But this is an unwarranted generalization about how computers work. Computations are typically predictable step-by-step: the results of each instruction a computer executes are formally specified and highly reliable. But the overall results of a computation need not be predictable in the slightest; indeed, because a computer can execute so many steps of an algorithm so quickly, it can produce outputs of astonishing complexity. This complexity is on full display for LLMs. Even compared with other types of programs, they are notoriously inscrutable. While scholars have begun to probe the internal representations employed by LLMs as they generate outputs, this work is still at a relatively early stage.²⁸⁸

Second, there are claims that LLMs will be more predictable because they are trained on huge corpora of human-written text.²⁸⁹ This argument fails for the same reason as claims that LLMs will be more accurate because they are trained on huge corpora of human-written text. The training process by itself guarantees nothing. Whether or not an LLM's outputs are predictable is an empirical question, and our empirical demonstrations in Part III show that frequently they are not.

Third, there are claims that the low cost and high accessibility of LLMs will give people access to them in advance of litigation, so that they can predict how courts or arbitrators would rule.²⁹⁰ This is not necessarily a claim based on accuracy. It is not necessary that the LLM be correct for it to be predictable; all that is required is that parties can obtain the same results *ex ante* by consulting the LLM as judges would obtain *ex post* by consulting it. The fly in the ointment is that for parties to predict outcomes they need to anticipate courts' prompting and configuration choices—many of which are dispute-specific. The more sensitive a method of generative interpretation is to those choices, the more precisely parties must predict those choices. If parties must negotiate or litigate dispute-specific choices to achieve predictable results, the efficiency benefits of generative interpretation shrink. And if those

²⁸⁸ See, e.g., Jack Lindsey, Wes Gurnee, Emmanuel Ameisen, Brian Chen, Adam Pearce, Nicholas L. Turner, Craig Citro, David Abrahams, Shan Carter, Basil Hosmer, Jonathan Marcus, Michael Sklar, Adly Templeton, Trenton Bricken, Callum McDougall, Hoagy Cunningham, Thomas Henighan, Adam Jermy, Andy Jones, Andrew Persic, Zhenyi Qi, T. Ben Thompson, Sam Zimmerman, Kelley Rivoire, Thomas Conerly, Chris Olah & Joshua Batson, *On the Biology of a Large Language Model*, ANTHROPIC (Mar. 27, 2025), <https://transformer-circuits.pub/2025/attribution-graphs/biology.html> [<https://perma.cc/5U4J-F2QP>].

²⁸⁹ *Generative Interpretation*, *supra* note 2, at 485 (referring to “the vast corpus on which [GPT-4] sits”); see also David Hoffman, *Consumers’ Unreasonable Textual Expectations*, 15 HARV. BUS. L. REV. 43, 50 (2025) (“The gist of our claim [in *Generative Interpretation*] is that large language models are particularly good at offering majoritarian readings of phrases read in the context of contracts, and can probabilistically estimate what meaning the ordinary reader of that document will take.”).

²⁹⁰ See, e.g., Gandall et al., *supra* note 30, at 4–5.

negotiations now require access to a new, expensive kind of additional technical expertise, generative interpretation might simply exacerbate challenges indigent parties face in gaining access to justice.

Finally, it is important to remember that predictability alone is not sufficient for legal legitimacy. An interpretive method that always reads text in the way that most favors the party with the longer name is cheap, accessible, *and* predictable. It would improve access and reduce costs. Yet we suspect most readers would agree that name length is not a legitimate factor in adjudicatory reasoning.²⁹¹

VI. REWORKING GENERATIVE INTERPRETATION

We aren't writing to announce unbridled AI pessimism. Our concern is that generative-interpretation proponents' influential proposals treat LLM output as probative of facts that they do not actually signify. Although these particular uses of LLMs are problematic—or rather, the problems they raise have not yet been surmounted—the good news is that there are many legal applications where LLMs can realize these benefits without running into the same problems.

A. *How AI-Generated Artifacts Can Acquire Normative Weight*

Broadly speaking, we see three ways in which LLMs may be fit for use in the legal process: for their purely rhetorical value, when they have been experimentally validated, or when they are subject to adversarial testing.

1. *Rhetorical Value*

LLM proponents often appeal to AI's strength as a rhetorician—and sometimes being rhetorically effective is exactly what one wants, and *all* that one wants. Unikowsky, for example, writes that after he presented Claude with the merits-stage briefs of adjudicated cases from the Supreme Court's most recent term, Claude decided all but ten of them “in the same way the Supreme Court did. In the other 10 . . . , I frequently was more persuaded by Claude's analysis than the Supreme Court's.”²⁹² If Claude is as persuasive a legal writer as Unikowsky suggests, we see no downsides to harnessing LLMs' rhetorical power *per se*. LLMs are legitimate assets to the legal profession at least insofar as they are harnessed for their rhetorical power. Lawyers could ask LLMs to generate persuasive arguments for their client's positions, then personally

²⁹¹ Cf. Brainerd Currie, *Conflict, Crisis, and Confusion in New York*, 1963 DUKE L.J. 1, 11–13 (1963) (arguing that if uniformity and predictability were all that mattered, a choice of law rule selecting Alaska law for all disputes, whether or not they had anything to do with Alaska, would be ideal).

²⁹² *In AI We Trust II*, *supra* note 5.

ensure that these outputs adhere to the other formal requirements of legal briefing. And lawyers and judges alike could ask LLMs to generate compelling counterarguments to their conclusions, then use those counterpoints to strengthen their briefs and opinions.²⁹³ Deployed for their rhetorical power, LLMs might be valuable tools to help judges and advocates strengthen their reasoning. Our argument here lines up with Waldon et al.'s recommendation for "dialectical legal AI," which similarly endorses LLM usage to "generate novel perspectives that challenge or expand a judge's theory of a disputed legal text," but not the "delegat[ion] [of] . . . authority to the model."²⁹⁴

To return to our fair-coin analogy one last time, consider the practice of tossing a coin to make a major decision. Deciding to follow the coin's "choice" is arbitrary.²⁹⁵ But tossing the coin and introspecting on whether the outcome feels right or wrong is a popular trick for getting insight into your internal compass. In that case, the coin toss's outcome is arbitrary, but the insights it generates are not.²⁹⁶ And, just as importantly, the coin's value as a tool to probe one's own thoughts does not depend on whether the coin actually is unbiased.

²⁹³ See, e.g., Andrew Coan & Harry Surden, *Artificial Intelligence and Constitutional Interpretation*, 96 U. COLO. L. REV. 413, 482 (2025) ("Another benefit of LLMs in legal decision-making is their ability to provide multiple perspectives when prompted appropriately.").

²⁹⁴ Waldon et. al, *supra* note 16, at 53–55.

²⁹⁵ A surprising number of judges have been sanctioned for decision-by-coin-flip. *E.g.*, In re Daniels, 340 So. 2d 301, 309 (La. 1976) (sanctions for pretending to decide guilt or innocence of criminal defendants by coin-flip); In re Application of Friess, 91 A.D.2d 554, 556 (N.Y. App. Div. 1982) (sanctions for determining sentence-length by coin flip); In re Brown, 662 N.W.2d 733, 733 (Mich. 2003) (sanctions for deciding which parent has child custody on Christmas by coin-flip); Jud. Inquiry & Rev. Comm'n v. Shull, 651 S.E.2d 648, 659–60 (Va. 2007) (sanctions for determining child visitation rights by coin-flip). Some legal theorists have considered the (narrow) circumstances under which randomness is preferable to human judgment. See, e.g., Bernard E. Harcourt, *Post-Modern Meditations on Punishment*, in CRIMINAL LAW CONVERSATIONS 163, 165, 167–68 (Paul H. Robinson, Stephen Garvey & Kimberly Kessler Ferzan eds., 2009) ("[T]here always came a moment when the empirical facts ran out . . . and yet the reasoning continued. There was always this moment when the moderns—those paragons of reason—took a leap of faith. . . . Where our social scientific theories run out, where our principles run dry, we should leave the decision-making to chance. We should no longer take that leap of faith, but turn instead to the coin toss."); JEPH LOEB & TIM SALE, BATMAN: THE LONG HALLOWEEN 330–31 (DC COMICS, 1998) (depicting comic book villain Two-Face arguing that coin flips are normatively preferable to courts or legislatures because random chance is impartial); Adam M. Samaha, *Randomization and Adjudication*, 51 WM. & MARY L. REV. 1, 34 (2009) ("In a subset of these cases, however small, randomization will be the theoretically superior option for reasons of practical indivisibility, equality norms, nagging uncertainty, incentive effects, and/or experimental value. That we have difficulty identifying this class of cases with precision is no reason to think it is an empty set."); see generally Ronen Perry & Tal Zarsky, "May the Odds Be Ever in Your Favor": *Lotteries in Law*, 66 ALA. L. REV. 1035 (2015) (examining the role of randomization in law).

²⁹⁶ See Mariela E. Jaffé, Leonie Reutner & Rainer Greifeneder, *Catalyzing Decisions: How a Coin Flip Strengthens Affective Reactions*, PLOS ONE, Aug. 2019, at 13 (explaining that "a decision aid [such as a coin flip] can function as a catalyst in the decision-making process"); Mariela E. Jaffé & Rainer Greifeneder, *Deciding Advantageously after Flipping a Coin*, 223 ACTA PSYCHOLOGICA, Mar. 2022, at 6 (finding that the coin-flip technique may improve decision-making under certain controlled conditions).

Of course, rhetorical use of AI isn't risk-free—but the risks are largely that rhetorical use might invite unfounded assumptions about LLMs' authority. For example, a lawyer or a judge may ask an LLM to produce competing arguments about an issue in a case. Any of the arguments that the LLM produces can help its interlocutor think through the issues and refine a brief or an opinion. The LLM user may, however, find that the LLM is able to argue for one result, A, more persuasively than for a competing result, B. If the LLM user takes this behavior as evidence that the LLM "believes" in the result it argued for more persuasively, and thus that A is correct and B is incorrect, then the user has stopped using the LLM as a rhetorical tool and has instead begun to defer to it, unfoundedly, as an epistemic authority. An LLM might craft a more persuasive argument for A than B for any number of arbitrary reasons, and the arguments that an LLM produces are not perforce the strongest ones. That an LLM argued for A more persuasively than B is not evidence that the strongest argument for A beats the strongest argument for B, and it would be an error to interpret as such.

But evaluating arguments is precisely what lawyers and judges know how to do best. So long as they approach LLMs as rhetoricians—not as oracles whose minds need to be read—lawyers and judges can use LLMs' powers of persuasion to test their own reasoning and make their writing more persuasive.

2. *Experimental Validation*

Limiting LLM use to purely rhetorical applications would obviate one of the most exciting potential benefits of AI: using automation to reduce judicial and administrative backlogs. By limiting our qualitative assumptions about LLM-produced text to those that are susceptible to objective benchmarks, we can expand the set of situations in which LLM use is justifiable. These applications of AI do not involve the splashy, novel issues of interpretation that make for exciting technical demonstrations. Instead of answering unfamiliar interpretive questions, LLMs can assist with the exact opposite: run-of-the-mill cases, where fact patterns rarely deviate from a handful of common situations, and voluminous precedent provides ready-made benchmarks for testing predictive accuracy.

In these cases, inductively validating the factual accuracy of LLM outputs is well within the technical capabilities of most social scientists. For example, an AI that can detect choice-of-venue clauses in a contract, or specific, common but invalid clauses in rental agreements from certain jurisdictions, could be tested against a curated list of similar agreements. A tool that performs well against that benchmark might be able to organize dockets and direct judicial attention in ways that increase efficiency without supplanting any part of the adjudicatory process in any individual case. Similarly, AIs could be tested on their ability to detect or correct procedural deficiencies in court and agency filings, allowing petitioners to amend those filings without an extra visit to the courthouse.

The common thread across these examples is *verifiability*. By using LLMs to perform rote tasks with verifiable results, we can establish ex ante the likelihood that those results are accurate and can monitor ex post whether the system remains accurate over time. With well-designed and continually retested benchmarks, LLMs could automate portions of the legal process while conforming to the same principles and standards used for other kinds of scientifically verified mechanization and measurement.²⁹⁷

Inductively validating LLMs is at least a tractable problem; this is a kind of science that computational legal scholars know how to do. But a comparison with immensely simpler and yet still highly controversial empirical interpretive projects—trademark surveys and corpus linguistics—shows that here the mountains are high and we are barely in the foothills.²⁹⁸ Effective legal benchmarks would need to provide a set of tests representative of the variety and complexity of interpretation tasks that face courts, and would need some well-defined “ground truth” to measure against.²⁹⁹

Today’s high-profile legal-reasoning benchmarks are far from representative.³⁰⁰ An LLM’s ability to pass the bar exam, for instance, is impressive but uninformative.³⁰¹ The bar exam is designed to test humans, and it assumes that test-takers are capable of lateral applications of knowledge. No LLM has been shown to be capable of that kind of lateral “thinking.”³⁰² Unless courts need an LLM to do *exactly* what a benchmark tests (e.g., answer bar exam questions from the same corpus, using the same format, writing style, level of complexity, and subject matter), that benchmark is not necessarily predictive of LLM accuracy. Though initial work on building better benchmarks is underway, those efforts currently focus on questions with well-defined

²⁹⁷ We do not mean to imply that all the questions surrounding the use of scientific and empirical evidence in courts are settled. Just that these situations—unlike blind faith in an LLM’s output—do not require courts to make any *new* descriptive or epistemic assumptions.

²⁹⁸ See generally, e.g., Barton Beebe & Jeanne C. Fromer, *Are We Running Out of Trademarks? An Empirical Study of Trademark Depletion and Congestion*, 131 HARV. L. REV. 945 (2018) (thoughtful empirical study of trademark depletion); Jens Frankenreiter & Michael A. Livermore, *Computational Methods in Legal Analysis*, 16 ANN. REV. L. & SOC. SCI. 39 (2020) (overview of computational textual methods).

²⁹⁹ See generally Michael J. Hasday, *Accuracy and the Robot Judge*, 25 J. APP. PRAC. & PROCESS 1 (2025) (arguing that AI judges could be benchmarked by comparing their decisions to decisions made by human judges).

³⁰⁰ See Florencia Marotta-Wurgler & David Stein, *Building a Long Text Privacy Policy Corpus with Multi-Class Labels*, 63 PROC. ASSOC. COMPUTATIONAL LINGUISTICS 8156, 8156 (2025).

³⁰¹ See Will Douglas Heaven, *AI hype is built on high test scores. Those tests are flawed*, MIT TECH. REV. (Aug. 30, 2023), <https://www.technologyreview.com/2023/08/30/1078670/large-language-models-arent-people-lets-stop-testing-them-like-they-were/> [<https://perma.cc/8756-7LDZ>] (discussing how to evaluate AI test results).

³⁰² See Marotta-Wurgler & Stein, *supra* note 300, at 8163–64 (“At least for the systems we tested, an LLM’s ability to answer one legal question appears to not be predictive of that LLM’s ability to answer other questions, even within extremely narrow domains . . .”).

right answers—like checking outputs for rhetorical and factual consistency.³⁰³ Subjective tasks—like the interpretation of ambiguous or contested terms—have no well-defined right answers to test against.³⁰⁴

3. Adversarial Testing

When a judge uses an LLM, they are recruiting an expert without the *Daubert* rule, without impeachment, without a jury, without the adversarial checks on which our legal system relies. By contrast, if a litigant introduces LLM outputs as evidence, they are forced to justify that use in a way that can survive evidentiary checks and adversarial challenge. In contrast to generative interpretation, where parties have no opportunity to challenge factual assumptions that judges make about the LLM-generated text, we take no similar issue with the use of LLM outputs by litigants. In those cases, the trial process *is* the verification process.³⁰⁵

B. Two Examples

Two other approaches to interpretation—trademark surveys and corpus linguistics—show what it takes to make empirical measurements of linguistic meaning normatively legitimate in adjudication. In each case, a field of expert specialists has grappled seriously with the problem of providing both descriptive and epistemic justifications for their methods. They have provided descriptive justifications by developing a consistent and relatively standardized set of methods; they have provided epistemic justifications by being precise about what their methods measure and modest in their ambitions. They show a possible road forward for generative interpretation, but that road is a long and hard one.

³⁰³ See generally Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi, Tom Zur, Varun Iyer & Zehua Li, *LEGALBENCH: A Collaboratively Built Benchmark for Measuring Legal Reasoning in Large Language Models* (Aug. 23, 2023) <https://arxiv.org/pdf/2308.11462> [<https://perma.cc/K7G5-WX8M>] (describing LEGALBENCH, presenting an empirical evaluation of 20 LLMs, and illustrating the types of research explorations LEGALBENCH enables).

³⁰⁴ Again, there is some nascent progress here; some recent benchmarks make the initial steps towards measuring an LLM's ability to navigate the grey areas in legal interpretation. See, e.g., Marotta-Wurgler & Stein, *supra* note 300, at 8164.

³⁰⁵ Cf. *Generative Interpretation*, *supra* note 2, at 506–09 (“[Judges will] want to be careful about parties’ manipulative [use of LLMs] . . .”).

1. Trademark Surveys

Survey evidence is broadly accepted and widely used in trademark and false advertising cases.³⁰⁶ Many core issues in trademark litigation involve questions of how consumers understand words or other marks, including whether a mark describes a product feature or the source of the goods, whether the defendant's use of a mark refers to their own goods or to the plaintiff's, and what a challenged statement conveys to consumers. These questions are all fundamentally empirical. While a judge's intuitions about language may play a role, the legal question is not whether "Fish-Fri" in the abstract uniquely refers to the plaintiff's coating mix or to any coating mix for frying fish, but whether *consumers understand* "Fish-Fri" narrowly to refer to the plaintiff's coating mix or broadly to refer to any coating mix. Indeed, even evidence that specific individual consumers have been confused about the source of the goods they are buying is only one factor in the tests used to ask whether consumers in general are likely to be confused.³⁰⁷ The interpretive community of consumers is sovereign in trademark law.³⁰⁸

To establish descriptive validity, judges, scholars, and experts have built up an extensive body of standardized best practices for trademark surveys.³⁰⁹ A typical survey design must consider how to identify consumers in the relevant market, how to present the marks to them, how and in what order to ask questions about their perceptions of those marks, what control questions to include to establish a relevant baseline, and what statistical analyses to apply to the raw data of participants' answers. All of these constrain the process so that it is more reliable and replicable. Even so, litigated trademark cases frequently feature the parties trying to impeach each other's surveys, requiring judges to probe the fine details of how questions are worded.³¹⁰

³⁰⁶ See 5 J. THOMAS MCCARTHY, MCCARTHY ON TRADEMARKS AND UNFAIR COMPETITION § 32:158 (5th ed. 2025) ("Survey Evidence is Routinely Received and Weighed by the Courts."); Shari Seidman Diamond & David J. Franklyn, *Trademark Surveys: An Undulating Path*, 92 TEX. L. REV. 2029, 2040 (2014) (noting the presence of over 1000 trademark surveys in reported cases between 1991 and 2012).

³⁰⁷ See, e.g., *JL Beverage Co., LLC v. Jim Beam Brands Co.*, 828 F.3d 1098, 1106 (9th Cir. 2016) (reciting an "eight-factor" test for likelihood of confusion, one factor of which is "evidence of actual confusion").

³⁰⁸ See Stanley E. Fish, *Interpreting the Variorum*, 2 CRITICAL INQUIRY 465, 483–84 (1976) (arguing that the meaning of text is created by interpretive communities).

³⁰⁹ See generally, e.g., JAMES T. BERGER & R. MARK HALLIGAN, TRADEMARK SURVEYS: A LITIGATOR'S GUIDE (2023); TRADEMARK AND DECEPTIVE ADVERTISING SURVEYS: LAW, SCIENCE, AND DESIGN (Shari Seidman Diamond & Jerre B. Swann eds., 2d ed. 2022); JACOB JACOBY, TRADEMARK SURVEYS, VOLUME 1: DESIGNING, IMPLEMENTING, AND EVALUATING SURVEYS (2015).

³¹⁰ See, e.g., *Novartis Consumer Health, Inc. v. Johnson & Johnson-Merck Consumer Pharmaceuticals Co.*, 290 F.3d 578, 590–95 (3d Cir. 2002) (considering whether survey questions were improperly leading). Scholars continue to propose refinements to survey methodology. See, e.g., Irina D. Manta, *In Search of Validity: A New Model for the Content and Procedural Treatment of Trademark Infringement Surveys*, 24 CARDOZO ARTS & ENT. L.J. 1027, 1065–70

As for epistemic validity, it is striking how modest judicial reliance on trademark surveys is, notwithstanding the immense effort and expense that go into them. Every circuit has a detailed multifactor test for consumer confusion, including factors that are not susceptible to surveying, such as the defendant's good or bad faith and the quality of their products.³¹¹ Judges admit but discount survey evidence when they have concerns about its methodology,³¹² and they weigh survey evidence against other sources of evidence on consumer understandings.³¹³ In short, the legal system treats all of the issues that trademark surveys speak to as discrete factual questions about the beliefs of actual consumers. Surveys are a powerful and potentially probative source of evidence about those discrete factual questions, but they are not regarded as conclusive even on those questions, let alone on an entire case. Scholars continue to offer powerful critiques that question whether trademark surveys actually measure what they purport to.³¹⁴

2. *Corpus Linguistics*

If trademark surveys are a mature empirical interpretive tool, corpus linguistics is in the process of attempting to become one. Corpus methods are well-established in linguistics, where they provide new sources of data on large-scale patterns of linguistic use. Like generative interpretation, corpus linguistics and law developed from the combined efforts of academics³¹⁵ and judges.³¹⁶ The basic idea is to take a large corpus of existing linguistic data—a collection of texts written by actual humans, such as the Corpus of

(2007) (proposing new quality standards for product stimuli, control stimuli, survey questions, and study environments, coupled with judicial pre-clearance of surveys).

³¹¹ See generally Barton Beebe, *An Empirical Study of the Multifactor Tests for Trademark Infringement*, 94 CALIF. L. REV. 1581 (2006) (discussing multifactor tests carried out by different circuits).

³¹² See, e.g., *Jordache Enters., Inc. v. Hogg Wyld, Ltd.*, 828 F.2d 1482, 1487–88 (10th Cir. 1987).

³¹³ Cf., e.g., Dan Sarel & Howard Marmorstein, *The Effect of Consumer Surveys and Actual Confusion Evidence in Trademark Litigation: An Empirical Assessment*, 99 TRADEMARK REP. 1416, 1430 (2009) (“[I]n the aggregate[,] the odds of winning improve only slightly with the presentation of actual confusion or survey evidence.”).

³¹⁴ See, e.g., Barton Beebe, Roy Germano, Christopher Jon Sprigman & Joel H. Steckel, *Consumer Uncertainty in Trademark Law: An Experimental Investigation*, 72 EMORY L.J. 489, 543 (2023); Sepehr Shahshahani & Maggie Wittlin, *The Missing Element in Trademark Infringement*, 110 IOWA L. REV. 1247, 1251 (2025).

³¹⁵ See generally Stephen C. Mouritsen, *The Dictionary Is Not a Fortress: Definitional Fallacies and a Corpus-Based Approach to Plain Meaning*, 2010 BYU L. REV. 1915; Stephen C. Mouritsen, *Contract Interpretation with Corpus Linguistics*, 94 WASH. L. REV. 1337 (2019); James C. Phillips, Daniel M. Ortner & Thomas R. Lee, *Corpus Linguistics and Original Public Meaning: A New Tool to Make Originalism More Empirical*, 126 YALE L.J.F. 21 (2016).

³¹⁶ See, e.g., *State v. Rasabout*, 2015 UT 72, 356 P.3d 1258, 1275–90 (Lee, A.C.J., concurring in part and concurring in the judgment); *Wilson v. Safelite Grp., Inc.*, 930 F.3d 429, 439 (6th Cir. 2019) (Thapar, J., concurring in part and concurring in the judgment); *Richards v. Cox*, 2019 UT 57, 450 P.3d 1074; *People v. Harris*, 885 N.W.2d 832 (Mich. 2016); *Facebook, Inc. v. Duguid*, 592 U.S. 395, 410–11 (2021) (Alito, J., concurring in the judgment).

Contemporary American Usage (1.1 billion words from 1990 to 2010)—and run an algorithm over it to obtain an answer to an interpretive question of interest.

Again, consider this empirical tool through the lenses of descriptive and epistemic validity. A linguistic corpus by itself says nothing; it is simply a collection of things people have written or said. To use a corpus to answer specific interpretive questions—e.g., does a person “discharge” a weapon once for each shot, or once for a group of shots fired in close succession?—one must find a way to interpret the corpus. In this example, from *State v. Rasabout*,³¹⁷ Associate Chief Justice Lee used a collocation analysis. He identified 86 cases in which “discharge” was used within 5 words of a term relating to a firearm; of them, 27 said or implied that a “discharge” was a single shot, and only 1 referred to a group of shots.³¹⁸

This is a specific empirical protocol, and like any empirical protocol, it can be challenged on grounds of unreliability. Scholars like Kevin Tobia have done significant empirical work questioning whether corpus-linguistic methods for legal interpretation are replicable in the scientific sense; Tobia found that they are “surprisingly unreliable.”³¹⁹ Prominent proponents have responded in defense of their methods, and the debate continues in the law-review pages.³²⁰ This back-and-forth is a normal part of the scholarly process; the point is simply that it is also a necessary part of establishing an empirical method as reliable in the sense required for use in adjudication.³²¹

Note, by way of comparison, how much more tightly constrained the task specification is in corpus linguistics than in generative interpretation. The algorithms used in contemporary legal corpus linguistics are straightforward: they can generally be specified in a single search query or a few lines of code, and they have broadly accepted common implementations. In contrast, generative methods are remarkably open-ended. One must write prompts

³¹⁷ 2015 UT 72, 356 P.3d 1258.

³¹⁸ *Id.* ¶¶ 88–91, 356 P.3d at 1281–82 (Lee, J., concurring in part and concurring in the judgment).

³¹⁹ Tobia, *supra* note 223, at 799.

³²⁰ See, e.g., Thomas R. Lee and Stephen C. Mouritsen, *The Corpus and the Critics*, 88 U. CHI. L. REV. 275, 279 (2021); Kevin Tobia, *The Corpus and the Courts*, U. CHI. L. REV. ONLINE (2021), <https://lawreview.uchicago.edu/online-archive/corpus-and-courts> [<https://perma.cc/58XE-X6DX>].

³²¹ Of course, it is possible that this process could cut in favor of generative interpretation compared with other empirical interpretive methods. Tobia’s surveys found that dictionaries and corpus linguistics did a worse job of tracking lay survey participants’ assessments of linguistic meaning; judges who did not use these tools came closer to the participants’ assessments. Tobia, *supra* note 223, at 753–77. It is entirely plausible that generative interpretation could outperform dictionaries and corpus linguistics, precisely because generative systems are engineered to emulate usage patterns in their training data. But this is a fundamentally empirical question, and note what could be measured. No experiment can determine whether a generative method yields correct results, because there is no accessible source of ground truth for legal meaning. (If there were, then legal interpretation would be unnecessary, and so would generative interpretation.) All that can be measured, even in principle, is the extent to which generative methods *track* other measurements of meaning, such as dictionaries, survey participants, and judges’ opinions.

(a natural-language task!), give them to a model whose internal functioning and semantic correlations are incredibly poorly understood, and then often interpret the outputs as well (another natural-language task!). Similarly, there are only a handful of corpora in wide use for corpus-linguistics purposes in law, while there are dozens of plausible models one could consult for generative interpretation, and more are being introduced constantly.³²² Indeed, two of legal corpus linguistics' most prominent advocates, Thomas Lee and Jesse Egbert, have criticized generative interpretation because it lacks “*transparency* in the methods used and *replicability* and *generalizability* in the findings that are generated.”³²³

Just as with trademark surveys, to pick among competing methods, one must have a theory of epistemic validity based on an underlying normative interpretive theory. Here, the theoretical debates among corpus linguists and their critics show why these theories are deeply contested; there is no simple linguistic truth that can be read off unproblematically from a corpus. Tobia points to the deep tension about what to do when corpus-linguistics methods conflict with other sources of meaning; the divergence could mean that corpus linguistics comes closer to the true collective linguistic meaning, that corpus linguistics falls further away from collective linguistic meaning, or that collective linguistic meaning itself is a construct that does not exist as a ground truth in the world that can be objectively ascertained.³²⁴ Scholars like Anya Bernstein and Carissa Byrne Hessick observe that there is necessarily a gap between the empirical facts that corpus linguistics can ascertain and the normative claim that its results should be treated as authoritative.³²⁵

3. *Implications*

Our point is not that these gaps cannot be overcome, either for corpus linguistics or generative interpretation. Trademark surveys have overcome them, within their domain. They provide a kind of interpretive evidence that is regularly used and is widely accepted as legitimate. Corpus linguistics is doing its level best to do the same. Still, it is being forced to articulate and explain why particular methods are empirically reliable and why the results of those methods should be regarded as legitimate interpretive data. It is slow, laborious, contentious work.

Our greatest fear is that generative interpretation will attempt to circumvent these challenges rather than take them seriously. Given the laboriousness

³²² Peter Henderson, Daniel E. Ho, Andrea Vallebuono & Cassandra Handan-Nader, *Corpus Enigmas and Contradictory Linguistics: Tensions Between Empirical Semantic Meaning and Judicial Interpretation*, 25 MINN. J.L. SCI. & TECH. 127, 153 & n.95 (2024).

³²³ Lee & Egbert, *supra* note 192, at 4.

³²⁴ See Tobia, *supra* note 223, at 735–36.

³²⁵ Anya Bernstein, *Legal Corpus Linguistics and the Half Empirical Attitude*, 106 CORNELL L. REV. 1397, 1400 (2021); Carissa Byrne Hessick, *Corpus Linguistics and the Criminal Law*, 2017 BYU L. REV. 1503, 1509 (2018).

of corpus-linguistic methods, generative interpretation can seem like a remarkable shortcut: rather than specifying a detailed algorithmic query that will provide circumstantial evidence of meaning, why not just ask a model the interpretive question directly and go with what the output says? But that approach has, to quote Bertrand Russell, all “the advantages of theft over honest toil.”³²⁶ For generative interpretation to fulfill the hopes of its advocates, it will have to confront all of the tedious empirical challenges and all of the contentious normative issues that trademark surveys and corpus linguistics have had to deal with. That will be no mean feat.

Consider how generative interpretation would look if it followed a trademark-survey model. Parties to a case involving an interpretive question would commission their own experts to conduct generative-interpretation studies. The experts would specify detailed experimental methods: the models and systems they use, the prompts they provide to those models, the outputs they examine, how they interpret those outputs, the number and nature of variations on the inputs they try, and the statistical tests they apply to summarize the outputs across numerous prompts. Every choice they made in doing so would be adversarially probed by opposing counsel and experts: Why Claude 4 Opus and not ChatGPT-o3? Why did the prompt use “plaintiff” and “defendant” instead of the parties’ names? As we have shown above in Part III, many of these choices can dramatically affect the results of generative interpretation. They will have to be pinned down in a principled way, and that will require hard conversations about what makes one form of generative interpretation more reliable than another. As we have shown in Part IV, there is no way to resolve these debates without a normative theory of what it is that generative interpretation is supposed to be measuring. Trademark surveys have just such a theory, but it is narrower and less ambitious than what generative interpretation currently aims for. Corpus linguistics shows what the process of working out such a theory looks like—a drawn-out debate that simultaneously gets deep into the technical weeds and ascends to the highest peaks of jurisprudence.

Or compare Lisa Larrimore Ouelette’s proposal to use search results as evidence of a trademark’s strength and of the degree of overlap between two marks’ likelihood of confusion.³²⁷ Her suggestion combines trademark surveys’ domain (consumer understandings of marks), corpus linguistics’s quantitative methods (count occurrences with a given word sense), and generative interpretation’s use of Internet datasets and reliance on complex black-box machine-learning systems.³²⁸ But the modesty of Ouelette’s proposal leaps off

³²⁶ BERTRAND RUSSELL, INTRODUCTION TO MATHEMATICAL PHILOSOPHY 71 (2d ed. 1920).

³²⁷ See generally Lisa Larrimore Ouelette, *The Google Shortcut to Trademark Law*, 102 CALIF. L. REV. 351 (2014).

³²⁸ For additional hybrids, see generally Jake Linford & Kyra Nelson, *Trademark Fame and Corpus Linguistics*, 45 COLUM. J.L. & ARTS 171 (2022) (arguing that corpus-linguistic methods have advantages over surveys for assessing trademark fame); Shivam Adarsh, Elliott Ash, Stefan

the page: she discusses in extensive detail the ways in which search results might be unstable or unrepresentative of linguistic usage; she individually examines every case in her dataset where a court and Google disagree; she emphasizes the limited nature of her empirical results; and her bottom line is only that search results are probative evidence that courts should not categorically exclude.³²⁹ This is the kind of scholarly foundation on which a persuasive argument for generative interpretation will ultimately need to rest.

VII. CONCLUSION

The legal philosopher Ronald Dworkin famously deployed the figure of Judge Hercules, an interpreter of unlimited competence and intelligence, with perfect access to all of the relevant authorities and all the time and diligence needed to study and harmonize them. Judge Hercules, Dworkin argued, was the ideal towards which human judging aspires.³³⁰

Generative AI proponents see AIs as modern miracles, capable of tackling any cognitive task with superhuman skill. So perhaps it should come as no surprise that a version of this boosterism has taken hold in the legal academy. In their various ways, the thinkers we have been calling “LLM proponents” describe ChatGPT, Claude, and the like in terms that recall Dworkin’s Judge Hercules. Like him, LLMs have mastered a corpus of text far larger than any human ever could. Like him, they can produce an answer to a complex problem far faster than any human ever could. And like him, they can knit together disparate and conflicting sources into a seamless justification as persuasively as any human could.

To all of this we say: *perhaps*. Even granting that LLMs are capable of Herculean feats, a deep and fundamental problem remains. *How can we be confident that Judge Hercules really has decided the case before him competently and in good faith?* Anyone could claim to be Judge Hercules, to have mastered the world’s accumulated statutes and precedents. And even if Judge Hercules is who he claims to be, the fact that he purports to present the best

Bechtold, Barton Beebe & Jeanne Fromer, *Automating Abercrombie: Machine-Learning Trademark Distinctiveness*, 21 J. EMPIRICAL LEGAL STUD. 826 (2024) (using machine learning to predict trademark distinctiveness).

³²⁹ Ouelette, *supra* note 327, at 369, 384–95, 406–07.

³³⁰ See generally RONALD DWORKIN, *LAW’S EMPIRE* (1986); cf. Kieffaber, *supra* note 116, at 31–32 (arguing that an ideal AI would be “Scalia’s Hercules, with a complete understanding of *language alone* rather than the various capricious elements that inform a Dworkinian positivist”). But of course, as we have been arguing at length, there is no reason to assume that any particular LLM’s outputs reflect “language alone.” For one thing, a typical LLM’s training process includes extensive pre-training on heavily curated datasets, fine-tuning for specific use cases, and reinforcement learning from human feedback to optimize its responses against human preferences. See generally Lee et al., *supra* note 123 (discussing model development and the numerous choice points it offers to shape a system’s behavior). To claim that this process learns “language” and nothing else is to hold up a randomly selected kitchen tool and assume that it is a can opener.

harmonization of the law is no proof that he is telling the truth about it. It does no good to point out that he can also write iambic pentameter or functional JavaScript code, when the question that needs answering is whether each interpretation he produces really is faithful to the legal authorities. Perhaps he forgot, or misunderstood, or hallucinated the nuances of the ERISA regulations on which the case depends. Perhaps he is a textualist in insurance cases but consults legislative history in admiralty cases. Perhaps he systematically favors corporate plaintiffs. Perhaps he knows how most people would understand a phrase, but ignores that knowledge when a simpler interpretation is available. Perhaps he finds arguments more persuasive when they SOMETIMES BREAK INTO ALL CAPS. Or perhaps the superhuman judge is also a superhuman liar.

LLMs are like humans in some ways and unlike them in others, and whether an LLM should or should not be treated as equivalent to a human in a particular context depends on the LLM, the human, and the context.³³¹ A merely human judge's written opinion is a kind of proof of work. It is meant to persuade the reader of the legal conclusion it presents, but it is also meant to persuade the reader that the conclusion was arrived at competently and in good faith. A judge can fall short of the judicial ideal not just by being *too bad* at legal writing—writing an opinion that fails to persuade because it lays bare the threadbare reasoning that resulted in it—but also by being *too good* at it—writing an opinion that is too persuasive because it carries the reader along too briskly to notice the better road not taken. True, the test of persuasiveness has never been sufficient to ensure complete honesty and competence in judging. But LLMs threaten to explode it altogether.³³²

Something more is required to make generative interpretation legitimate, some assurance that *this particular* output from among *all possible* outputs was generated fairly, objectively, on the basis of the evidence, and in accordance with the relevant authorities. To date, LLM proponents have not paid sufficient attention to what that something more might be. We hope that they will.

³³¹ Cf. Robert Brauneis, *Copyright and the Training of Human Authors and Generative Machines*, 48 COLUM. J.L. & ARTS 1, 22–56 (2024) (exploring in detail whether LLMs should be regarded as fundamentally human-like, fundamentally inhuman, or somewhere in between for purposes of copyright infringement).

³³² Cf. Eric A. Posner & Shivam Saran, *Judge AI: Assessing Large Language Models in Judicial Decision-Making* 28 (Univ. Chi. Coase-Sandor Inst. for L. & Econ. Rsch., Working Paper No. 25-03, 2025) (“If the goal is to produce AI judges that operate like human judges, success would be achieved only if the AI judges decide cases in a realist way while using formalist reasoning—meaning that they do not explain how they actually decide the cases. It is hard to imagine such AI judges being acceptable in a democracy or any well-ordered political system.”); Re, *supra* note 23, at 1585 (“The existence of legal norms and elites has always depended on there being a scarcity of persuasive resources and arguments. . . . By undermining these constraints, a surfeit of persuasiveness threatens the effectiveness of legal norms.”).