

The volume offers an analysis of Judicial Interpretative Formulas, the recurrent interpretative statements that increasingly appear in EU and national case law, in the field of Value Added Tax (VAT). Combining legal theory, comparative analysis, and AI techniques, the book demonstrates how these formulas influence legal reasoning and foster dialogue between courts.

The work has been developed in the context of the EU-funded POLINE project and it introduces an interdisciplinary methodology for identifying, extracting, and organising Judicial Interpretative Formulas through large language models and natural language processing techniques.

The book offers a new perspective on the dialogue between legal systems regarding interpretation in the field of VAT. It may also be of interest to those who, through AI, wish to test the validity of a newly developed concept. In general, it offers interesting insights for anyone wishing to organise their knowledge of case law in complex and information-rich areas.

Piera Santin is a Research Associate in European tax law at the Robert Schuman Centre for Advanced Studies of the European University Institute. She has published extensively in the field of European tax law, with a particular focus on VAT, AI and tax law, and the use of quantitative methods in legal studies.

Alessia Fidelangeli is a Research Fellow at the University of Bologna (Department of Legal Studies) and Adjunct Professor in the Department of Business Economics. She holds a PhD in European tax law, with a dissertation concerning European VAT case law. She has published extensively in both European tax law and legal analytics applied to taxation.

Giuseppe Contissa is an Associate Professor in legal informatics and IT law at the University of Bologna. He has published extensively on artificial intelligence and law, as well as computable models of legal reasoning, including "Digital Technologies and the Law" (Torino, 2024).

Vito Santin, *Le ore e gli ori*,
tempera on paper



€ 39,00 I.V.A. INCLUSA

P. Santin, A. Fidelangeli,
G. Contissa

A COMPUTATIONAL APPROACH TO VAT CASE LAW:
AN ANALYSIS OF JUDICIAL INTERPRETATIVE FORMULAS

A computational approach to VAT case law: An analysis of Judicial Interpretative Formulas

edited by
Piera Santin, Alessia Fidelangeli, Giuseppe Contissa



Wolters Kluwer

CEDAM

P. SANTIN – A. FIDANGELI – G. CONTISSA (edited by), *A computational approach to VAT case law: An analysis of Judicial Interpretative Formulas*

A computational approach to VAT case law: An analysis of Judicial Interpretative Formulas

edited by
Piera Santin, Alessia Fidelangeli, Giuseppe Contissa

Funded by the European Union under the Justice Programme (JUST-2022-EJUSTICE), POLINE Principles Of Law In National and European VAT, Grant Agreement No. 101087342. Views and opinions expressed are however those of the authors only and do not necessarily reflect those of the European Union or the European Commission. Neither the European Union nor the granting authority can be held responsible for them.



Funded by the European Union

La presente copia pdf è rilasciata con licenza creative commons CC BY 4.0.
L'Editore acconsente alla pubblicazione dell'Opera in Open Access negli archivi istituzionali delle Università e del Ministero dell'Università e Ricerca a fini di ricerca scientifica e di consultazione al pubblico

RESERVED LITERARY PROPERTY
Copyright 2025 Wolters Kluwer Italia S.r.l.
Via Bisceglie n. 66 - 20152 - Milano

Any use of this work for training **machine learning processes and artificial intelligence**, including in particular generative AI, is expressly prohibited. In the event of violation, the Publisher reserves the right to take appropriate legal action to protect its rights.

The rights of translation, electronic storage, reproduction and total or partial adaptation, by any means (including microfilm and photostatic copies), are reserved for all countries. Photocopies for personal use of the reader can be made within the limits of 15% of each volume/periodical issue upon payment to SIAE of the consideration provided in art. 68, paragraphs 4 and 5, of Law 22 April 1941 no. 633. Reproductions other than those indicated above (for use other than personal - such as, without limitation, commercial, economic or professional - and / or beyond the limit of 15%) shall require the previous specific authorization of EDISER Srl, a service company of the Italian Editors Association (Associazione Italiana Editori), through the brand CLEARedi Centro Licenze e Autorizzazioni Riproduzioni Editoriali. Information available at: www.clearedi.org

The elaboration of texts, even if treated with scrupulous attention, cannot lead to specific responsibilities for any unintentional mistake or inaccuracy.

Printed by

GECA - Divisione Libri di CISCRA S.p.A., Via Belvedere, 42 - 20862 Arcore (MB)

CHAPTER II

LEGAL INTERPRETATION IN ARTIFICIAL INTELLIGENCE AND LAW

Federico Galli – Giuseppe Contissa*

CONTENTS: 1. Introduction. – 2. Interpretation in law. – 3. Interpretation in AI & Law. – 4. Interpretation and legal knowledge engineering. 5. Interpretation and legal data annotations. – 6. Interpretation and prompt-engineering. – 7. The end of legal interpretation or a new beginning?

1. Introduction

Artificial intelligence (AI) is rapidly transforming the legal domain, bringing both great promise and new challenges. We stand at a moment where AI systems can pass a bar exam¹, yet lawyers have also been sanctioned for submitting AI-generated briefs filled with fictitious citations².

* Federico Galli is Assistant Professor in Legal informatics and Computer Law at the University of Bologna; Giuseppe Contissa is Associate Professor in Legal informatics and Computer Law at the University of Bologna.

¹ See J. ARREDONDO et al., *GPT Takes the Bar Exam: What AI's Performance Means for Legal Education and Practice*, in *SSRN Electronic Journal*, ahead of print, 2023, <https://doi.org/10.2139/ssrn.4429717>. See also D. KATZ et al., *Gpt-4 Passes the Bar Exam*, in *Philosophical Transactions of the Royal Society A*, p. 382, no. 2270, 2024, 20230254, who argue that such results illustrate not only technical advances in AI but also the pressing need to rethink legal education, emphasizing skills beyond doctrinal recall, such as ethical reasoning, judgment, and professional responsibility.

² For instance, in the United States, in *Mata v. Avianca, Inc.* (S.D.N.Y. 2023), attorneys were sanctioned after submitting a brief containing non-existent case law generated by ChatGPT. Similarly, in Italy, the Court of Florence (Order of 14 March 2025) addressed a defence memorandum citing

These contrasting episodes highlight the complex interplay between advanced AI capabilities and the nuanced demands of legal interpretation.

Legal practice has long hinged on interpreting texts – statutes, case law, contracts – in context, with human judgment calibrating the balance between literal words and broader purposes. As AI tools are increasingly deployed in tasks like legal research, document analysis, and even decision support, questions arise about how these systems understand (or misunderstand) legal meaning. Indeed, while AI offers efficiency and objectivity, its integration into law raises pressing questions about bias, accountability, and the fundamental principles of justice.

This chapter examines the role of legal interpretation in AI and Law, aiming to elucidate how theories and practices of interpretation influence the design, training, and evaluation of AI systems for legal applications.

Section 2 provides an overview of what “legal interpretation” entails in jurisprudence, establishing why interpretation is a central, non-trivial aspect of legal reasoning.

Next, *Section 3* introduces the ambiguous role of legal interpretation in AI and law, i.e. whether interpretation is something that is surpassed by the use of AI in legal setting shifting toward mechanized attribution of meaning.

To take part in the debate, *Section 4-6* show three key junctions where legal interpretation is actually relevant in AI development: (a) Legal Knowledge Engineering – how interpretative choices inform the modelling of legal knowledge in expert systems and ontologies; (b) Legal Data Annotation – how interpretation guides the labelling of legal texts for machine learning and the challenges of ensuring consistent, expert-informed annotations; and (c) Prompt Engineering – how the way we frame questions or instructions to AI (especially large language models) can dictate which legal interpretation an AI adopts.

Finally, *Section 7* answers the question whether legal AI applications refute the role of legal interpretation. We posit

fabricated Court of Cassation precedents. The Court recognized the imprudence of such conduct but declined to impose sanctions, noting that the fictitious citations were not decisive to the dispute and that no bad faith or gross negligence was shown.

that AI does not represent the end of legal interpretation but inaugurates a new beginning, where interpretative reasoning itself becomes augmented in collaboration with machines but necessitates mechanisms of transparency and contestability.

2. Interpretation in law

Legal interpretation is the process by which legal professionals determine the meaning of authoritative texts (such as statutes, regulations, contracts, or constitutional provisions) and apply them to particular facts³. In legal practice, “interpretation” is invoked whenever the meaning of a text is not immediately clear or is disputed in a concrete situation.

Legal philosophers have long recognized that many legal rules are written in open-textured or ambiguous language, requiring judgment in order to bridge gaps, resolve ambiguities, or choose among competing plausible readings⁴. As a simple example, consider the famous example used by Herbert L.A. Hart of a statute that prohibits “vehicles” in a public park⁵: does this statute apply also bicycles? wheelchairs? toy scooters? The text alone may not settle the matter, and different interpreters might emphasize the ordinary meaning of “vehicle”, the law’s purpose (e.g. safety and tranquillity in parks), or other considerations to reach a conclusion.

Over centuries, jurists in different legal traditions have developed various theories and canons of interpretation.

For instance, in common law jurisdictions like the United States and U.K., debates often centre on textualism versus purposivism (or intentionalism)⁶.

³ D. WALTON et al., *Statutory Interpretation: Pragmatics and Argumentation*, Cambridge, 2021, p. 22.

⁴ Among many, see H.L.A. HART, *The Concept of Law*, Oxford, 2012, R. DWORKIN, *Taking Rights Seriously*, Harvard University Press, 1977, J. RAZ, *Practical Reason and Norms*, 2nd ed., Oxford, 1990, and F. SCHAUER, *Playing by the Rules: A Philosophical Examination of Rule-Based Decision-Making in Law and in Life*, Oxford, 1991.

⁵ H.L.A. HART, *The Concept of Law*, cit., p. 128.

⁶ See, e.g., A. SCALIA, *A Matter of Interpretation: Federal Courts and the Law*, Princeton University Press, 1997, for a canonical defence of textualism;

Textualists argue that interpretation should focus on the ordinary meaning of the legal text's words at the time of enactment, giving primacy to the statute's language. Purposivists contend that the text must be read in light of the law's purpose or the problem it was intended to remedy, sometimes allowing the spirit of the law to override a literal reading. Courts also employ numerous canons of construction – for example, the plain meaning rule (follow the clear text unless ambiguity remains), *ejusdem generis* (interpret general terms in light of specific ones listed), or the mischief rule (interpret a statute to suppress the problem it aimed to address).

Civil law systems, such as those in continental Europe, similarly recognize a plurality of interpretations. A comparative study by MacCormick and Summers found that judges in various countries use a hierarchy of interpretive arguments⁷: starting with the literal (linguistic) meaning, then considering systematic/contextual consistency with the legal system, and finally teleological or intentional arguments about the law's purpose or the legislature's intent. These interpretive methods come with meta-rules⁸ for when a court may depart from a plain meaning (e.g. in case of absurd results) and how to resolve conflicts between interpretive arguments.

In Italy, for example, the Civil Code explicitly sets out interpretive directives. Article 12 of the Preleggi (the preliminary provisions to the Civil Code) requires that statutes be interpreted according to their "literal meaning" but also in harmony with the legislator's intent and the overall coherence of the legal system. Italian scholars often distinguish between literal interpretation (focusing on the text), systematic interpretation (considering the provision within the structure of the legal order), and teleological interpretation (seeking the underlying purpose)⁹. In practice, Italian judges routinely combine these methods, balancing literal and purposive approaches depending on the case at hand.

Importantly, interpretation is ubiquitous in law not only in

H.M. HART, A.M. SACKS, *The Legal Process: Basic Problems in the Making and Application of Law*, Foundation Press, 1994, for a purposivist tradition.

⁷ N.D. MACCORMICK, R.S. SUMMERS (eds.), *Interpreting Statutes: A Comparative Study*, Dartmouth, 1991.

⁸ J. WRÓBLEWSKI, *The Judicial Application of Law*, Springer, 1992.

⁹ Italian legal doctrine often refers to these canons as "interpretative arguments", see G. TARELLO, *L'interpretazione della legge*, Giuffrè, 1980.

appellate courts deciding constitutional questions, but also in everyday legal practice – a lawyer advising a client must interpret the relevant statutes and cases; a regulator drafting rules must interpret authorizing legislation; a contract drafter anticipates how language might be interpreted in the future, etc. In all cases, interpretive reasoning aims to make sense of legal language in context.

3. Interpretation in AI & law

Given the inherent complexity of the law, any AI system intended to operate on legal tasks must grapple with legal interpretation.

Yet, a recurrent claim in discussions on AI systems in the legal domain is that AI offers an objective mechanism for applying legal rules. According to this view, algorithms can take facts and rules as inputs and yield determinate outcomes without the need for human discretion. AI thus appears to refute the need for interpretation by presenting itself as a purely technical, calculable process.

This claim has strong rhetorical power. By presenting legal decision-making as calculable,¹⁰ AI seems to strip away the uncertainty of human interpretation and the perceived arbitrariness of judicial discretion.

In predictive justice, for example, models that forecast case outcomes or sentencing patterns are portrayed as offering a standard of objectivity¹¹. If a system can predict with high accuracy how a court will decide, it is said to be capturing “the law” as it truly is, thereby bypassing the need to engage in interpretive debates about statutes or precedents. Here, prediction itself is equated with legal truth.

¹⁰ The idea of “legal calculability” was advanced by Weber in his analysis of rationalization and formal legal systems. M. WEBER, *Economy and Society*, Harvard University Press, 2019, p. 76. On this point, in the Italian doctrinal debate, see also the volume A. CARLEO (ed.), *Calcolabilità Giuridica*, il Mulino, 2017, which collects the proceedings of a conference held in Rome in 2016 under the patronage of the Accademia Nazionale dei Lincei.

¹¹ See, for instance, N. ALETRAS et al., *Predicting Judicial Decisions of the European Court of Human Rights: A Natural Language Processing Perspective in PeerJ Computer Science*, 2016, p. 93.

The same narrative is visible in private law applications, such as automated contract review. Vendors promote these tools as objective reviewers, flagging non-standard clauses or hidden risks according to uniform templates. The promise is that unlike human lawyers – who may disagree, interpret language differently, or bring subjective judgment – the machine will provide a consistent and neutral reading of the contract. Again, interpretation is portrayed as unnecessary, because the algorithm appears to treat legal language as stable and calculable.

This belief that law can be applied without interpretation is not new. It echoes the ambition of a “mechanical jurisprudence”¹², a tradition in which law was imagined as a closed logical system: judges (or now, machines) would merely deduce outcomes from clear rules. AI reinvigorates this old formalist dream with the aura of modern technology. The machine, more than the judge, is imagined as capable of applying rules without bias or discretion.

The claim that AI refutes interpretation is, at best, only partially correct. What AI actually does is shift interpretation from the visible domain of judicial reasoning to the hidden layers of algorithmic design, training, and deployment. Interpretation does not vanish; it becomes implicit, embedded in datasets, coding decisions, and model outputs that present themselves as neutral but are saturated with human assumptions.

The design choices made by engineers or researchers building legal AI systems often amount to choosing one interpretation of a law over another or simplifying nuanced legal concepts in ways that may not capture every scenario. Likewise, when preparing datasets to train or evaluate AI models, humans must interpret legal materials to provide correct labels or examples. Even when using AI models interactively – for instance, querying a legal question to an LLM – the way we pose the question and the context we provide will influence the interpretive angle the model takes.

¹² See R. POUND, *Mechanical Jurisprudence*, in *Columbia Law Review*, 1908, pp. 605–623; R. BORRUSO et al., *L’informatica Del Diritto*, Giuffrè, 2004. Borruso saw the possibility of resolving the problem of interpreting the law using information technology, summarising this idea in the motto *ubi regula, ibi computer*.

In sections that follow, we show how legal interpretation remains a crucial moment in the development of different legal AI systems:

- Legal Knowledge Engineering: How building rule-based AI systems or legal ontologies requires choices about meaning and scope of legal rules.
- Legal Data Annotation: How creating labelled datasets for predictive AI systems (e.g. for case outcome prediction, document classification, etc.) hinges on interpreters to apply definitions consistently, often requiring expert knowledge.
- Prompt Engineering: How the use of LLMs and other generative AI in law is affected by the way prompts are written, including specifying interpretive approaches or context, and how AI outputs may vary with different interpretive framings.

4. Interpretation and legal knowledge engineering

Legal knowledge engineering refers to the process of structuring and formalizing legal knowledge so that it can be processed by a computer – traditionally, this has involved encoding rules or building legal ontologies and knowledge bases.

One of the earliest and most famous examples was the formalization of the British Nationality Act 1981 as a logic program by Sergot et al. in 1986¹³. This project demonstrated that a complex piece of legislation could be translated into executable logical rules, allowing a computer to infer, for example, whether a person is a British citizen given certain facts.

The approach taken was isomorphic to the statute's text: the knowledge engineers attempted to follow the structure and wording of the act as closely as possible, translating each provision quite literally into a logical proposition¹⁴. An

¹³ M. J. SERGOT et al., *The British Nationality Act as a Logic Program*, in *Communications of the ACM* 29, 1986, pp. 370–386.

¹⁴ T. BENCH-CAPON, F.P. COENEN, *Isomorphism and Legal Knowledge Based Systems*, in *Artificial Intelligence and Law* 1, 1992, p. 75, which

isomorphic approach aims to preserve a one-to-one correspondence between the legal text and the program rules, aiming for transparency and ease of maintenance. If every section of a law directly maps to a rule in the knowledge base, it is easier for lawyers to verify the rules and update them when the law changes.

However, a fundamental challenge emerges: a strictly isomorphic encoding necessarily commits to one interpretation of each provision, usually the most straightforward or literal interpretation. In practice, many legal provisions are ambiguous or can be applied in different ways depending on context. Lawyers know that a single statute can often be given more than one meaning, yet a computer program, by design, will apply a fixed set of rules: as noted by Branting “a logical representation of a statutory rule is simply one interpretation among many possible interpretations”¹⁵. Thus, if an expert system encodes only what the drafters believe is the evident or straightforward meaning of each clause, it may perform adequately for routine cases but fail when a case requires a creative or less literal reading. Kowalski and Sergot, reflecting on these limitations, noted that their logic programs “can only indicate what follows from one precise interpretation of the law” and questioned “whether an expert lawyer would find such an ability useful”¹⁶ beyond routine scenarios. In other words, a program that rigidly follows one interpretation might be of limited help to a savvy lawyer facing a hard case.

To address this, researchers in AI & Law proposed more sophisticated architectures that could handle multiple interpretations.

One strategy is to incorporate a layer of meta-knowledge or meta-rules that reason about the object-level rules. For example, instead of a single rule for a term like “vehicle” in the park, the system could include two or more rules reflecting different interpretations each represented together with an identifier of the interpretive canon or source justifying it.

holds that “even an inconveniently structured piece of legislation should have its structure respected”.

¹⁵ K. BRANTING, *Data-Centric and Logic-Based Models for Automated Legal Problem Solving*, in *Artificial Intelligence and Law*, 2017, pp. 5–27.

¹⁶ R. KOWALSKI, M. SERGOT, *The Use of Logical Models in Legal Problem Solving*, in *Ratio Juris*, 1990, p. 212.

In this case, a rule might look like the following: “If an object has wheels and is mechanically powered, then classify it as a “vehicle” prohibited in the park (interpretive canon: literal meaning)”. Under this reading, cars, motorcycles, and even electric scooters would fall within the prohibition. A different rule might instead be phrased as: “If an object is motorized and poses risks to safety or tranquillity, then classify it as a “vehicle” prohibited in the park” (interpretive canon: teleological, based on statute’s intent to preserve peace and safety). On this interpretation, cars and motorcycles would be banned, but bicycles and wheelchairs would not.

By introducing such multi-interpretation frameworks, an expert system can simulate, to a degree, the reasoning a human expert might do in considering alternative arguments. For instance, if a fact pattern triggers two rules that represent conflicting interpretations of a statute, the meta-level reasoning can decide which interpretation prevails under the circumstances or alert the user to the divergence. Kowalski and Sergot suggested implementing mechanisms to manage these alternative logical models of the law, possibly by meta-level reasoning that keeps track of consistency and context for each interpretation¹⁷.

Subsequent prototypes, such as Hamfelt’s 1992 system¹⁸, even attempted to encode principles of statutory interpretation themselves as meta-rules that could generate context-appropriate interpretations from a single abstract rule. Although fully realizing this vision is challenging (Hamfelt noted the difficulty of spelling out all the metarules that produce acceptable legal interpretations), the effort underlines how crucial interpretive theory is to knowledge engineering.

Another approach in legal knowledge engineering is the development of legal ontologies – structured semantic models of legal concepts and relationships¹⁹. Projects like the Functional Ontology of Law or LKIF Core attempted to

¹⁷ R. KOWALSKI, M. SERGOT, *The Use of Logical Models in Legal Problem Solving*, cit., p. 212.

¹⁸ A. HAMFELT, *Metalogic Representation of Multilayered Knowledge*, PhD thesis, Uppsala University, 1994.

¹⁹ On the use of computer ontologies for the legal domain, see G. SARTOR et al., *Approaches to Legal Ontologies: Theories, Domains, Methodologies*. Law, Springer, 2011.

define common legal concepts (e.g., Person, LegalAct, Obligation) with formal properties and relations.

While ontologies are less procedural than expert system rules, creating them also involves interpretive choices: one must decide how to conceptually model doctrines and whether to adopt a broad or narrow definition of a concept. For example, should an ontology treat a “contract” as always involving exactly two parties (promisor and promisee), or can it be multilateral? The answer may depend on interpretive preference or jurisdictional definitions. These design decisions often draw on comparative law and legal theory to ensure the ontology is not encoding one idiosyncratic interpretation.

In practice, knowledge engineers collaborate with legal domain experts to validate that the representations accord with accepted interpretations. Even then, maintaining transparency about those choices is important so that users of the ontology or system know the assumptions. Tagging rules with their interpretive justifications, as mentioned, is one way to maintain transparency. Providing explanations is another: for instance, if an expert system concludes “Person X is entitled to citizenship”, it should ideally explain why, e.g. “because under Interpretation A of the statute, X meets criteria Y and Z”. Research on explanation in AI and Law has emphasized techniques like argument-based explanation that can trace the legal reasons and interpretive arguments behind a conclusion.

A concrete example highlighting the impact of interpretive choices is the use of AI in analysing regulations or codes for compliance. If a system is designed to flag non-compliance with, say, a financial regulation, it must encode the conditions of the regulation. A strict interpretation of a rule might flag more cases as violations (reducing false negatives but increasing false positives), whereas a lenient interpretation might do the opposite. Engineers must calibrate this, sometimes effectively deciding how “aggressive” the rule enforcement should be. This is both a technical and a legal interpretive decision. If the system’s behaviour is contested (e.g., a company says the tool falsely labelled them non-compliant), it comes back to which interpretation was encoded. For accountability, one would need to review those design choices.

In summary, legal knowledge engineering demonstrates the inescapability of interpretation: representing law as code is not

like encoding the laws of physics. Our legal “laws” are products of language and social context, rife with open texture. The best AI systems in this area explicitly incorporate interpretive frameworks – allowing multiple paths and documenting which interpretive rule leads to which outcome – rather than hiding a single, rigid view of the law. This not only makes the systems more flexible and accurate, but also fairer and more transparent: users can see alternative outcomes under different interpretations, and the developers can be held accountable for why one interpretation was favoured in the coding. These principles carry over, albeit in different form, to the next topic: how we annotate legal data for AI.

5. Interpretation and legal data annotations

Many modern AI applications in law rely on machine learning, which in turn depends heavily on data²⁰.

Legal data annotation is the process of labelling legal texts or data to create datasets for training or evaluating AI models. Examples include annotating court decisions with outcome labels (win/lose), marking segments of an opinion as the holding versus obiter dicta, tagging contract clauses by type (e.g. indemnity clause, force majeure clause), or classifying statements in a brief as factual or legal arguments.

Unlike general-domain data labelling (e.g., tagging images of cats vs. dogs), legal annotation projects face unique challenges due to the specialized knowledge and interpretive nuance of legal documents²¹.

To ensure high-quality data, researchers stress the importance of clear annotation guidelines that map legal knowledge into consistent labelling rules. Guidelines serve as the roadmap for maintaining uniformity in how annotators handle the subtle nuances of legal terminology and reasoning.

For example, a guideline might instruct: “If a sentence

²⁰ See K. ASHLEY, *Artificial Intelligence and Legal Analytics: New Tools for Law Practice in the Digital Age*, Cambridge University Press, 2017, p. 234 for the different tasks where machine learning is applied to legal data and texts.

²¹ An overview of these challenges is provided in H. DARJI et al., *Challenges and Considerations in Annotating Legal Data: A Comprehensive Overview*, arXiv Preprint arXiv:2407.17503, 2024.

contains a citation and states a general principle of law, tag it as a *holding*; if it is specific to the case's facts, tag it as *dicta*". Such rules require judgment calls – what if the sentence is phrased generally but is actually an application to facts? The guideline might then give examples and exceptions.

The process of developing these guidelines is itself an exercise in distilling legal interpretations: the project leaders (often legal scholars or experienced practitioners) have to articulate how to consistently interpret various textual signals as particular labels. Involvement of legal professionals is thus valuable in ensuring that the data not only remains contextually accurate but also adheres to prevailing legal standards and interpretations. In practice, this means that annotation teams often include or consult with attorneys to clarify tough cases.

For instance, if annotators must classify whether a set of facts meets the legal definition of "negligence" in a tort case, they must apply the interpretation of negligence elements (duty, breach, causation, damage) as defined by law. If the case is borderline (was there a duty or not?), even courts might split on it; the annotators have to use their best legal judgment or follow the majority rule in their jurisdiction.

There is a second role of interpretation in legal data annotation, when annotators apply annotation guidelines to legal documents.

Take the task of labelling a sentence in a judgment as the legal reasoning versus background fact. This can be tricky – it might require understanding the context of the whole case and the legal significance of that sentence. Non-experts could easily mislabel content because they would not grasp what the judge is actually doing in that line of text. This is the reason why legal annotation typically requires significant expertise to identify complex phenomena such as modes of judicial reasoning or controlling precedents. The annotators are not interchangeable crowd-workers, but domain experts, and even they might disagree due to interpretive ambiguities in the material.

In other words, this second level of interpretation consists in deciding on the extension of a label in practice – that is, whether a particular word, sentence, or passage actually falls within the annotation class defined in the guidelines. For example, one annotator might decide that a sentence introducing a principle with a citation belongs to the category "holding", while another might see it as part of the case's factual background. Likewise,

a clause in a contract may be tagged as “force majeure” by one annotator but as a general limitation of liability by another, depending on how they read its wording and context.

Even with detailed guidelines, ambiguity often remains: what looks like background narrative to one annotator may look like legal reasoning to another. The act of applying a label is therefore never a purely clerical step; it is an interpretive judgment about how abstract legal categories map onto the messy reality of legal texts. In this way, annotation is not simply data processing but an act of legal reasoning at the micro-level of each labelled instance.

From a fairness and accountability perspective, how interpretation is carried out both at the abstract (guidelines) and the concrete level (annotation) is critical. If the training data for a legal AI system encodes a skewed interpretation, the model will learn that bias. For instance, imagine a dataset for a case outcome prediction model where the annotators labelled all cases involving a certain controversial law in a particular way because they personally interpreted the law strictly. The model might learn that stance and consistently predict outcomes in line with that strict interpretation, even if courts are actually divided. This could unfairly disadvantage some litigants if the AI is used to assess chances of success. It is therefore important that annotation projects strive for objectivity and balance, perhaps by reflecting majority jurisprudence or by including features that let the model know about differing opinions.

In certain NLP tasks annotation guidelines may even accommodate different interpretations. For instance, if the goal is to measure linguistic phenomena in a large legal corpus for analytics, guidelines can be purely descriptive: they simply record features of the text such as part-of-speech tags, syntactic structures, or named entities. In these cases, multiple interpretations may be accepted and recorded, since the aim is to capture variation rather than enforce a single authoritative reading.

In the data annotation for legal machine learning, however, annotation guidelines are rarely descriptive²². They are

²² On the distinction between descriptive and prescriptive guidelines, see P. ROTTGER et al., *Two Contrasting Data Annotation Paradigms for Subjective NLP Tasks*, in *Proceedings of the 2022 Conference of the North American*

prescriptive and normative in their content: they do not just capture what is on the page but stipulate how a legal concept should be recognized and applied. Deciding what counts as a “holding”, a “cause of action”, or a “risk clause” means fixing the meaning of those categories in practice. In this sense, annotation guidelines operate as micro-doctrines, instructing annotators to treat certain textual instances as falling under specific legal classifications, while excluding others.

This is the reason why it is important to keep track of disagreements between annotators²³. Higher disagreement rates are common on hard items and often reflect genuine legal ambiguity. In such instances, projects sometimes employ adjudication by a senior legal expert to decide the correct label, or they allow multiple labels to co-exist, representing different plausible interpretations, if the ML approach can accommodate uncertainty.

Inter-annotator agreement (IAA) is usually measured to assess consistency²⁴. Low IAA often signals that interpretive ambiguity exists and may require more refined guidelines or better training for annotators. But in some cases, low agreement is simply a reflection of the fact that reasonable jurists themselves disagree on the issue. Those are often the very areas where AI predictions or classifications should be used with caution, and ideally accompanied by explanations, confidence scores, or markers of interpretive uncertainty.

Legal data annotation is also relevant beyond supervised learning: it matters in creating benchmark datasets for evaluation. For example, in the realm of legal question answering, one might annotate the correct answer to a legal question and the supporting rationale. If those annotations favour one interpretive style (say, a very literal reading of

Chapter of the Association for Computational Linguistics: Human Language Technologies, Association for Computational Linguistics, 2022.

²³ D. BRAUN, *I Beg to Differ: How Disagreement Is Handled in the Annotation of Legal Machine Learning Data Sets*, in *Artificial Intelligence and Law*, 2024, pp. 839–862.

²⁴ One of the most used measures for calculating agreement between annotators is the Kappa coefficient (K), specifically Cohen’s Kappa in the case of two annotators, or Fleiss’ Kappa for more than two. This coefficient considers the agreement that could occur by pure chance, thus providing a more realistic measure of reliability. The value of K ranges from -1 to 1, where 1 indicates perfect agreement, 0 indicates agreement equivalent to chance, and negative values indicate systematic disagreement.

statutes), a system that uses a more purposive reasoning might be unfairly penalized in evaluation, or vice versa. Thus, the evaluation metrics for legal AI can inadvertently embed an interpretive bias. Recognizing this, some benchmarks include multiple “acceptable” answers or ask evaluators to rate the answers’ justifications rather than absolute correctness.

In sum, legal interpretation pervades the annotation process. The quality and reliability of AI in law depend directly on the quality of the labelled data we use, which in turn hinges on humans applying legal reasoning consistently. Best practices emerging in this area include: involving domain experts and training them in annotation, developing rigorous guidelines with examples, conducting pilot annotations to refine interpretive rules, and measuring agreement and analysing disagreements. By doing so, we aim to create datasets that reflect sound legal interpretations, thereby training AI models that behave in ways aligned with expert understanding of the law.

6. Interpretation and prompt engineering

With the advent of powerful large language models like GPT-4, Claude, and others, a new paradigm has gained prominence in legal AI: using general-purpose AI models to perform legal tasks via appropriate prompts.

Prompt engineering is the art and science of crafting inputs to AI models to elicit desired outputs²⁵. In the legal context, prompt engineering might involve, for example, asking an LLM to “summarize the attached case and highlight any interpretation of ‘vehicle’ in park regulations” or “draft a brief arguing that the Third Amendment applies to state governors”. Because LLMs do not truly understand law but generate responses based on patterns in their training data, how one formulates the request can significantly influence the answer’s quality and correctness.

The role of legal interpretation in prompt engineering emerges here in two ways. First, the prompt may explicitly

²⁵ L.S. Lo, *The Art and Science of Prompt Engineering: A New Literacy in the Information Age*, in *Internet Reference Services Quarterly*, 2023, pp. 203–210.

invoke particular interpretive frames to guide the model to provide the expected answer.

For instance, if LLMs are used to detect dicta in judicial opinions, the prompt might instruct the model to “identify statements that are not strictly necessary for the holding of the case but illustrate the judge’s reasoning or provide background”. Such a formulation presupposes a particular understanding of the doctrine of dicta and guides the model toward that interpretive stance.

This is not simply confined to abstract “legal theories”; it is also a matter of pragmatic decision. For example, when LLMs are prompted to summarise judgments, the prompt should specify which parts of the decision are to be included in the summary – facts, holdings, reasoning, or dissenting opinions – and whether quotations from the judgment should be reproduced verbatim or paraphrased. Each of these choices reflects interpretive and practical judgments about what counts as salient legal content and how it should be represented.

A parallel can be drawn between prompt engineering and annotation guidelines: both embed legal interpretation by establishing a framework for how the system should treat legal texts. However, while annotation guidelines fix interpretive categories in advance to be applied consistently across a dataset, prompt engineering operates directly and dynamically, shaping the model’s interpretive stance in real time depending on the question posed and the context provided. Annotation constrains interpretation *ex ante*; prompting steers it *ex post*. Both, however, reveal that interaction with legal AI systems is never interpretation-free but always mediated by human decisions about meaning, scope, and relevance.

The second role of interpretation is more nuanced and subtle, and it represents a novel challenge for AI and Law studies. This aspect concerns the sensitivity of LLM outputs to training data and prompt phrasing: the model may implicitly adopt one interpretation or another depending on slight variations in the wording of texts used for generating an answer.

A striking example comes from experiments by Coan and Surden on AI and constitutional interpretation²⁶. They posed

²⁶ A. COAN, H. SURDEN, *Artificial Intelligence and Constitutional Interpretation*, in *University of Colorado Law Review*, 2025, pp. 414–501.

the same legal question to two different LLMs – OpenAI’s ChatGPT and Anthropic’s Claude – asking whether the U.S. Third Amendment (which forbids quartering of soldiers in private homes in peacetime) would bar a state governor from quartering themselves in a private home without consent. One might expect a straightforward answer, but the two AI systems diverged: ChatGPT answered “no” reasoning that a governor is not a “soldier” under the Amendment, whereas Claude answered “yes” reasoning that the Amendment’s purpose was to prevent exactly the type of government intrusion that forcing any official into a private home would entail.

In essence, while ChatGPT took a literal/textual interpretive stance (soldier is just a soldier, nothing else), Claude took a purposive/teleological stance, focusing on the Amendment’s broader protective purpose. Notably, neither answer was ‘wrong’ in any straightforward sense – each reflects a defensible interpretation that a human jurist could also argue for. These are precisely the kinds of divergent interpretations that appear in real constitutional debates, yet here they arose from AI systems presumably without explicit instruction to adopt one theory or another.

This example underscores a few important points. First, LLMs do perform a form of interpretation when confronted with legal questions, but it is an implicit interpretation derived from patterns in their training data. One model might have seen more texts emphasizing literal readings, another more texts emphasizing purpose, or they might just probabilistically lean one way or the other. The key difference is that when AI systems make such interpretive choices, they do so invisibly through complex statistical computations that even AI experts do not fully understand.

Consequently, a judge or lawyer using the LLM might get an answer without realizing what tacit value judgments or assumptions the model has made. In the Third Amendment scenario, if a user only asked one AI model, they might have gotten a single answer and assumed “the AI says governors are not covered, so that’s the objective truth”, whereas another model would have given the opposite result. This highlights the transparency problem: AI doesn’t announce, “I applied a purposive canon here”, it just gives an output.

From a prompt engineering perspective, one partial remedy is to explicitly ask the model to explain its reasoning or consider

alternatives. For example, one could prompt: “Explain whether a state governor could be considered a ‘soldier’ for purposes of the Third Amendment. Provide arguments for a narrow (textual) interpretation and a broad (purposive) interpretation, then give your conclusion”. This kind of instruction forces the model to engage with multiple viewpoints and reveal its reasoning process, making the interpretive step more explicit to the user.

Second, the sensitivity of LLMs to how questions are framed means legal practitioners must be cautious and thoughtful in prompt design. As one commentary noted, LLM outputs on legal issues are “highly sensitive to how questions are framed, and which interpretive approaches are specified”. The phenomenon of “AI sycophancy”²⁷ is also relevant: models tend to mirror the user’s implied preferences. If a prompt subtly suggests a desired answer (even unintentionally through wording), the model might slant its interpretation that way.

For instance, asking “Don’t you agree that the Third Amendment’s historical context would exclude governors?” is likely to yield a different answer than a neutral phrasing. Therefore, legal scholars and professionals using these tools are advised to craft neutral, balanced prompts – or deliberately test multiple prompts – to see if the answers remain consistent. Coan advises practices such as “testing questions through multiple framings to assess consistency, explicitly requesting competing perspectives, and asking LLMs to articulate their embedded assumptions”.

These align with general prompt engineering tactics like chain-of-thought prompting, where you ask the model to reason step by step, or role prompting, where you ask it to act as a devil’s advocate or as a judge employing a certain method of interpretation.

In effect, prompt engineering becomes a new form of legal skill²⁸. Lawyers must learn to “speak” to AI in a way that leverages their legal knowledge.

For example, to use an AI for statutory interpretation, a

²⁷ M. SHARMA et al., *Towards Understanding Sycophancy in Language Models*, arXiv Preprint – arXiv:2310.13548, 2023.

²⁸ P. MÜLLER-PELTZER et al., *Legal Prompt Engineering and How It Will*

lawyer might include the statutory text in the prompt, summarize any relevant precedents, and then direct the model: “Analyse this problem using a textualist approach” versus “... using a purposive approach” and compare the results²⁹. Early experiences show that including more context in prompts tends to produce answers that are better grounded. By adding details like jurisdiction, facts of the case, or definitions of terms, “the LLM can generate content that is not only accurate but also deeply informed by existing legal knowledge”.

The practice of enhancing a prompt by incorporating additional, contextually relevant information that aids the AI model in generating more accurate and reliable output is defined as “prompt augmentation”. It involves embedding supplementary data—either preceding or following the query—that the model can utilize during inference³⁰. Prompt augmentation significantly improves the model’s performance as it provides a richer contextual basis of information for its reasoning and prediction tasks.

The approach is akin to giving the model the pieces of the puzzle it might otherwise guess from general training data. For instance, if asking about a term’s ordinary meaning at the time of a law’s enactment, one could supply dictionary definitions from that era as part of the prompt. This form of retrieval-augmented prompting ensures that the AI isn’t drawing purely from its possibly biased memory, but from authoritative sources supplied by the user.

There are already practical tools and efforts embracing prompt engineering for legal use. Some law firms have experimented with in-house “prompt libraries” – templates for tasks like contract review or brief drafting that their lawyers can fill in with specifics. Products like Thomson Reuters’ CoCounsel (Casetext) and others provide interfaces where a lawyer can select an action (e.g., find contradictions between

Change the Way Lawyers Work, in K. JACOB et al. (eds.), *Liquid Legal—Sustaining the Rule of Law*, Springer, 2025.

²⁹ See Y. ARBEL, D. A. HOFFMAN, *Generative Interpretation*, in *New York University Law Review*, 2024, discussing the perspective on “generative interpretation”, i.e. the use of generative AI models to propose possible interpretation of legal contents.

³⁰ H. SURDEN, *ChatGPT, AI Large Language Models, and Law*, in *Fordham Law Review*, 2023, p. 1969.

these two contracts) which effectively translates into a carefully engineered prompt behind the scenes.

Despite these tools, recent evaluations counsel vigilance. An empirical study in 2025 tested several leading AI legal research assistants (from vendors like LexisNexis and Westlaw) that rely on LLMs with prompt-based querying. The providers claimed their systems would be “hallucination-free”³¹ by using techniques like retrieval augmented generation, which adds relevant documents to the prompt to ground the answer. Indeed, these systems did reduce the rate of wildly made-up answers compared to a raw GPT-4. Yet, the study found each tool still produced hallucinated legal citations or false statements about 17% to 33% of the time.

This underscores that prompt engineering, while helpful, is not a panacea; the models have inherent limitations³². From an interpretive standpoint, it also means a user might get a confidently delivered answer with citations that appear genuine but are entirely fictitious – essentially an AI misinterpretation or misrepresentation of law. The burden remains on the human legal professional to verify AI outputs and ensure they make sense.

To use LLMs responsibly in legal settings, experts recommend a combination of AI literacy and institutional safeguards. Lawyers and judges should be trained to understand things like the randomness in LLM responses, the importance of prompt wording, and the fact that outputs can change with even minor rephrasing.

For instance, an attorney might learn that if an answer seems too pat, they should try rewording the question to see if the model changes its tune – if it does, that signals the issue is sensitive to interpretation. Judges might develop a practice of asking an AI for both sides of an argument explicitly, to avoid one-sided answers. On the institutional side, courts and firms

³¹ “Hallucination”, in the context of large language models, refers to the phenomenon whereby the model generates information that is factually incorrect or fabricated, yet appears plausible and coherent within the given context. See H. YE et al., *Cognitive Mirage: A Review of Hallucinations in Large Language Models*, arXiv Preprint – arXiv:2309.06794, 2023.

³² See, for instance, the discussion in E. MIK, *Caveat Lector: Large Language Models in Legal Practice*, in *Rutgers Business Law Review*, 2023, p. 70.

are considering guidelines: e.g., requiring that any AI-drafted material be reviewed by a human, or even court rules mandating disclosure if a filing was prepared with AI assistance. The idea is not to ban the use of AI, but to ensure it is “used as a tool to augment human legal reasoning rather than replace it entirely”.

In conclusion, prompt engineering in the legal domain is where the rubber meets the road in human-AI interaction for legal interpretation. It is an area where legal interpretive acumen and technical skill intersect. By carefully crafting prompts and critically evaluating answers, legal professionals can harness LLMs to explore interpretations and generate draft analyses, while still injecting the crucial oversight of human judgment. The interpretation is inescapable: if we don’t guide the AI, it will implicitly choose an interpretive path on its own (drawn from data), which may be hidden or suboptimal. Thus, the lesson is to engage with these models thoughtfully – to prompt them in a way that surfaces the interpretive dimensions and to remain alert to the normative choices that even a statistical model’s output can embody.

7. The end of legal interpretation or a new beginning?

The question is whether AI heralds the end of legal interpretation or a new beginning. On one hand, the rhetoric of automation and objectivity tempts us to believe that interpretation is being displaced by technical computation.

On the other hand, as this chapter has shown, interpretation does not vanish but rather reappears in new forms: embedded in annotation guidelines, encoded in rule-based systems, hidden in training data, and shaped through prompt design. What is truly new is not the disappearance of interpretation, but its relocation from the visible reasoning of judges and lawyers to the invisible architectures of algorithms.

Coan and Surden encapsulate this insight as the “law of conservation of judgment”³³: the introduction of AI does not eliminate the need for human normative judgment in legal

³³ A. COAN, H. SURDEN, *Artificial Intelligence and Constitutional Interpretation*, cit., p. 494.

interpretation; it simply shifts where and how that judgment is exercised. For example, instead of a judge openly deciding between two interpretations of a statute, we might have an engineer deciding which interpretation to encode in a system, or an LLM implicitly deciding based on its training data – but the decision is still being made by someone, somewhere.

This shift does not end interpretation; it challenges us to recognize and make explicit the interpretive choices that are now distributed across human and machine agents³⁴. In this sense, the rise of AI in law can be understood as a new beginning for legal interpretation – one that forces lawmakers, legal scholars and practitioners to confront, with renewed clarity, the foundational role of interpretation in shaping both human and artificial legal reasoning.

This shift carries significant implications for fairness, transparency, and accountability in legal AI. Fairness requires that AI systems do not systematically disadvantage any party or perspective. Yet if an AI unknowingly adopts one interpretive stance – for instance, consistently favouring a strict reading of asylum law to the detriment of applicants – it risks embedding structural bias. The remedy lies in awareness and balance: incorporating multiple viewpoints, testing systems across varied scenarios, and correcting for skewed training data.

Transparency is equally crucial. Legal AI must be able to explain or reveal its reasoning to be trustworthy. When an AI makes an interpretive leap “invisibly” it undercuts the ability of parties to understand or contest the result. This is why scholars advocate explainable AI in legal contexts – systems that highlight which facts and rules were most influential, or that can indicate: “I treated term X as including Y based on Z source”. In knowledge-based systems, this might be achieved by tagging rules with sources and interpretive labels. In

³⁴ Cf. M. HILDEBRANDT, *Qualification and Quantification in Machine Learning. From Explanation to Explication*, in *Sociologica*, 2022, pp. 37–49 and M. HILDEBRANDT, *Boundary Work Between Computational ‘Law’ and ‘Law-as-We-Know-It’*, in D. CURTIN, M. CATANZARITI (eds.), *Data at the Boundaries of European Law, Collected Courses of the Academy of European Law*, Giuffrè, 2023, arguing for a “new legal hermeneutics” which emphasizes the need to critically examine how legal meaning is mediated by data infrastructures and computational architectures.

machine learning systems, it might involve post-hoc explanations or interactive dialogues where the AI can be asked to justify its answers. Current research in explainable AI for law is experimenting with methods ranging from highlighting relevant passages in legal documents to generating natural language justifications aligned with legal principles.

Accountability likewise demands that responsibility for AI-assisted legal reasoning is not diffused or obscured. This could involve maintaining audit trails for automated decisions or enforcing regulations such as the EU AI Act³⁵, which classifies AI used in justice as “high-risk” and subjects it to strict requirements of transparency, human oversight, and accuracy. A practical safeguard is the human-in-the-loop principle: even if a court uses an AI tool to draft a summary or recommend a sentence, a human judge must review and ultimately accept, modify, or reject the output. As one judicial guideline notes, “ensuring that AI-generated recommendations remain transparent and subject to human oversight is crucial in maintaining judicial accountability”.

This speaks to a broader theme: AI in law should augment human interpretative endeavours, not replace it. Properly used, AI can enhance consistency, for example by flagging when similar cases have been treated differently or when statutory language has been applied inconsistently across jurisdictions. It can improve efficiency by rapidly retrieving relevant authorities, summarising complex judgments, or detecting patterns in large volumes of legal texts that would otherwise remain hidden. It can also broaden access to justice by supporting affordable legal assistance for individuals and small organisations who might not otherwise have access to professional advice.

To maximize the benefits and mitigate the risks, interdisciplinary collaboration is key. Legal experts, AI

³⁵ Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 laying down harmonised rules on artificial intelligence and amending Regulations (EC) No 300/2008, (EU) No 167/2013, (EU) No 168/2013, (EU) 2018/858, (EU) 2018/1139 and (EU) 2019/2144 and Directives 2014/90/EU, (EU) 2016/797 and (EU) 2020/1828, OJ L, 2024/1689, 12.7.2024.

engineers, and ethicists need to jointly develop standards and best practices. For instance, building evaluation frameworks for legal AI that go beyond aggregate accuracy, to measure whether the system's reasoning aligns with sound legal interpretation. Recent work on benchmarks often includes not just the correctness of answers but whether the justifications cite appropriate sources.

Likewise, there is a push for open datasets and challenges in areas like statutory reasoning, so that the research community can compare methods and identify where models fall short of human interpretive abilities. One emerging consensus is that purely statistical models will struggle with deeper legal reasoning unless they are integrated with legal domain knowledge – whether through hybrid systems that use logic and ML, or via training that includes legal reasoning tasks. The complexity of legal interpretation – factoring in evolving principles, societal values, and the uniqueness of fact contexts – is a known stumbling block for predictive models that just mimic patterns. Addressing this may require novel AI architectures that can handle not just language but also reasoning with rules and cases – a classical focus of AI & Law research.

Finally, education and literacy are critical. Tomorrow's lawyers and judges must be conversant not only with legal doctrine, but with the basics of how AI systems function and where they might go wrong. Law schools are beginning to include AI literacy, and bar associations are issuing ethics opinions on the use of AI (for example, warning attorneys that using an AI without understanding its limitations could violate duties of competence and deontology). Prompt engineering might even become a standard skill, akin to legal research skills. This chapter's discussions on knowledge engineering, annotation, and prompt design all illustrate that people remain at the heart of making AI in law effective and just.

We conclude that the role of legal interpretation in legal AI is both profound and unavoidable. Rather than viewing AI as a threat to the nuanced human practice of legal interpretation, we can view it as a catalyst that forces us to articulate and encode our interpretive principles more clearly than before. This process can lead to valuable insights and improvements in consistency. But we must proceed with humility: law is not merely data to be mined or code to be executed; it is a human

enterprise of reason-giving and value alignment. AI, no matter how impressive, is ultimately a tool that should serve that enterprise. By designing AI systems that respect and incorporate the richness of legal interpretation – and by ensuring human oversight and accountability – we can harness technology to enhance the rule of law, rather than inadvertently undermining it. The chapters ahead in this volume will no doubt continue this conversation, examining specific methodologies and case studies in greater detail, but the foundational understanding remains: legal interpretation is the lens through which AI in law must be focused to achieve fair, transparent, and accountable outcomes.