

Informative Patents? Predicting Invalidity Decisions with the Text of Claims

JAMES HICKS*

Patent litigation presents a puzzle. On the one hand, lawyers and academics spend a huge amount of time and attention on the intricacies of patent drafting and the proper construction of claims. On the other, empirical studies of patent litigation consistently find that the intrinsic characteristics of a patent have little relationship with its success or failure in invalidity proceedings. The conventional wisdom is that by the time they reach a decision on the merits, validity decisions—like much of patent litigation—are a largely unpredictable affair. This widespread sense that patent rights are uncertain has long been a concern to patent commentators and scholars. It has significant costs to the patent system, making investments in R&D more risky, and litigation both more expensive and difficult to insure.

In this Article, I use computational legal analysis to demonstrate that the conventional wisdom is incomplete. In fact, the content of patent claims is a surprisingly powerful predictor of invalidity decisions: the best performing model can correctly predict the outcome of approximately 70% of federal court decisions, using only the claim text of the disputed patent. These findings suggest that computational analysis can help to shed light on patent quality, and they demonstrate where predictive analytics might be usefully applied in patent litigation.

* Academic Fellow, Berkeley Center for Law and Business; james.hicks@berkeley.edu. I am grateful to John Allison, Mark Lemley, and David Schwartz both for sharing their hand-coded litigation outcomes and for providing thoughtful feedback on the project. For valuable comments and conversations, I thank Adam Badawi, Kristelia García, Uri Hacoheh, Ryan Hübert, Sonia Katyal, Aniket Kesari, Justin McCrary, Peter Menell, Tejas Narechania, Julian Nyarko, Kevin Quinn, Rachel Stern, Molly Van Houweling, Sam Zyontz, participants at the 2020 IP Scholars Conference, the 2021 Berkeley LEAP workshop, and a CPIP retreat. This work was supported by a grant from Berkeley's Law, Economics, and Politics Center.

Introduction	3
I. Prior Studies of Patent Litigation.....	7
II. Computational Analysis in Law	12
III. Predicting Invalidity.....	16
A. Data sources	16
B. Preparing the text for analysis.....	18
C. Methodological approach.....	20
D. Results.....	23
E. The challenge of interpretability	28
IV. Implications.....	31
A. New approaches to research on patent quality	31
B. What does this mean for patent litigation?.....	34
Appendix	38
A. Receiver operating characteristic (ROC) curves.....	38
B. Separation graphs.....	39

INTRODUCTION

Does the text of a patent convey important information in litigation? This seems, on the face of it, to be a strange question. Deciphering and delineating the scope of the claimed invention is perhaps *the* central question of patent law. The preliminary stage of litigation, during which a patent’s written claims are construed by the court—so-called *Markman* hearings—occupy a central place in patent commentary.¹ Landmark patent law cases involve the rules for the proper interpretation of claim language.² As Mark Lemley puts it: when drafting a patent claim, “every word matters.”³ And yet, when it comes to empirical studies of patent litigation, scholars find something surprising: with remarkable regularity, the observable, intrinsic features of the patent have no apparent relationship with the disposition of the crucial question of validity.

How can we explain this apparent paradox? It could be that the process of selection into (and progression through) litigation results in the adjudication of a set of patents whose validity is highly uncertain to the litigants—coin flips, essentially.⁴ Alternatively, it could be that patents have become highly technical scientific documents, largely impenetrable to juries (or even non-specialist district court judges) such that the actual text of the claims is, in practice, less important than other contextual factors of the litigation.⁵ In this Article, I suggest and test a different explanation: that the content of patent claims *is* consequential, but our empirical tools have yet to identify that relationship.

¹ See, e.g., J. Jonas Anderson & Peter Menell, *Informal Deference: A Historical, Empirical, and Normative Analysis of Patent Claim Construction*, 108 NW. U. L. REV. 1 (2014) (summarizing *inter alia* the history of, and scholarly debate, on claim construction).

² See *Phillips v. AWH Corp.*, 415 F.3d 1303 (Fed. Cir. 2005) (en banc) (holding that words in a claim must be given their ordinary meaning in the context of the whole patent, as understood by a person skilled in the relevant art).

³ Mark A. Lemley, *Without Preamble*, 100 B.U. L. REV. 357, 364 (2020).

⁴ This is an implication of the well-known Priest-Klein hypothesis. See George L. Priest & Benjamin Klein, *The Selection of Disputes for Litigation*, 13 J. LEGAL STUD. 1 (1984). Keith Hylton makes this argument directly in relation to patent litigation, observing that patent law has relatively little “systemic informational uncertainty,” and consequently litigants’ converging expectations of success (or failure) will tend to induce settlement. Keith N. Hylton, *Patent Uncertainty: Toward a Framework with Applications*, 96 B.U. L. REV. 1117, 1125–30 (2016).

⁵ Patent litigation is an unusual area of U.S. law in which juries are regularly called upon to decide invalidity disputes, despite the fact that such decisions involve significant questions of law. See Paul R. Gugliuzza, *Law, Fact, and Patent Validity*, 106 IOWA L. REV. 607 (2021).

Studies of “patent quality”—that is, the likelihood that a patent will be found valid if challenged⁶—are voluminous.⁷ Quality has long been a central theoretical and empirical question for observers of the patent system, many of whom argue that low-quality patents (*i.e.*, patents which are unlikely to be upheld) are the cause of many of the system’s ills.⁸ Legal scholars and economists have used litigation to search for reliable proxies for quality,⁹ but measuring the concept has proven somewhat elusive—and the selection effects inherent in any study of litigation complicate efforts to generalize the findings to non-litigated patents.¹⁰ In his study of patent litigation involving non-practicing entities, Michael Risch concludes that “there were unobserved patent quality factors that

-
- 6 There is no single definition of a patent’s “quality,” but the one that I adopt here is commonly used in the literature. *See, e.g.*, John R. Allison, *Patent Value*, in 2 RESEARCH HANDBOOK ON THE ECONOMICS OF INTELLECTUAL PROPERTY LAW 47, 52 (Peter Menell & David Schwartz eds., 2019); Michael Risch, *A Generation of Patent Litigation*, 52 SAN DIEGO L. REV. 67, 68 (2015); Brian J. Love, Shawn P. Miller & Shawn Ambwani, *Determinants of Patent Quality: Evidence from Inter Partes Review Proceedings*, 90 COLO. L. REV. 67 (2019). For an alternative (and more capacious) view of patent quality, *see* Christi J. Guerrini, *Defining Patent Quality*, 82 FORDHAM L. REV. 3091 (2013).
- 7 *See, e.g.*, Lee Petherbridge, *On Addressing Patent Quality*, 158 U. PA. L. REV. PENNUMBRA 13 (2009); Michael D. Frakes & Melissa F. Wasserman, *Does the US Patent and Trademark Office Grant Too Many Bad Patents?: Evidence From a Quasi-Experiment*, 67 STAN. L. REV. 613 (2015); Ronald J. Mann & Marian Underweiser, *A New Look At Patent Quality: Relating Patent Prosecution to Validity*, 9 J. EMPIRICAL LEGAL STUD. 1 (2012).
- 8 *See, e.g.*, R. Polk Wagner, *Understanding Patent-Quality Mechanisms*, 157 U. PA. L. REV. 2135, 2141 (2009) (arguing that “the current patent system has too much uncertainty, and that low patent quality bears substantial responsibility for it”). *But see* Dan L. Burk & Mark A. Lemley, *Fence Posts or Sign Posts? Rethinking Patent Claim Construction*, 157 U. PA. L. REV. 1743, 1751–1753 (2009) (arguing, *inter alia*, that patent claims may be inherently indeterminate).
- 9 A significant body of research, mostly in economics, focuses on observable characteristics such as forward citations (references to the patent in future applications) as measures of a patent’s contribution to the general stock of knowledge. *See generally* ADAM B. JAFFE & MANUEL TRAJTENBERG, *PATENTS, CITATIONS, AND INNOVATIONS: A WINDOW ON THE KNOWLEDGE ECONOMY* (2002). Another line of research connects quality to *value*, searching for proxies and signals such as the payment of renewal fees or data from auctions. *See, e.g.*, Christina Odasso, Giuseppe Scellato & Elisa Ughetto, *Selling Patents at Auction—An Empirical Analysis of Patent Value*, 24 INDUS. & CORP. CHANGE 417 (2015). However, in general, the extent to which a patent’s quality aligns with its value (either to society or to the patentee) is unclear. *See* Allison, *supra* note __ at 52 (contending that “quality in many cases is a necessary but insufficient condition for value”).
- 10 Allison, *supra* note __ at 73 (“...there are the many issues with selection effects and unmeasurable explanatory variables that plague studies of . . . the outcome of litigation”).

affected whether to mount a challenge and whether that challenge was successful, but those quality factors are [not] on the face of the patent . . .”¹¹ It seems likely that factors affecting patent quality are indeed important determinants of eventual success or failure in litigation, but whether or not we can observe them remains an open question.

In this Article I demonstrate that many of those patent quality factors are, in fact, on the face of the patent, and that with modern computational text analysis tools we can identify them directly. To test this question, I combine an extant dataset of patent litigation outcomes with the claim text of litigated patents. Using a machine-learning technique,¹² I train a statistical model to distinguish between valid and invalid patents based on the prevalence of specific words in their claims.¹³ The results are striking, showing that court decisions about patent validity are predictable to a much greater degree than has previously been recognized. With just the text of patent claims, a machine learning model can correctly predict 73% of invalidity outcomes in unseen “test” cases.¹⁴ Combining the claim text with other characteristics of the patent and the litigation further improves the model’s accuracy, although the text alone appears to encode much of the relevant information about the patent’s likelihood of success.

Developing a better understanding of predictability has significant implications for a patent system that has long been decried for its lack of certainty. Complaints about unpredictability are legion: academics and commentators alike have criticized nearly every area of patent law as uncertain, including claim construction and the scope of issued patents,¹⁵ jury decisions,¹⁶ and (os-

11 Risch, *supra* note __ at 122.

12 Specifically, I use an algorithm called a “random forest,” which is a popular tool for classification in statistics but is relatively novel in empirical legal analysis. *See infra* Part II.

13 *See infra* Parts III.A and III.B.

14 *See infra* Part III.D.

15 *See, e.g.,* Wagner, *supra* note __; Harry Surden, *Efficient Uncertainty in Patent Interpretation*, 68 WASH. & LEE L. REV. 1737, 1740 (2011) (“Scope uncertainty is problematic because it can reduce net innovation in particular fields of endeavor.”); Oskar Liivak, *The Unresolved Interpretive Ambiguity of Patent Claims*, 49 UC DAVIS L. REV. 1851 (2016); Gretchen Ann Bender, *Uncertainty and Unpredictability in Patent Litigation: The Time is Ripe for a Consistent Claim Construction Methodology*, 8 J. INTELL. PROP. L. 175 (2001).

16 Gugliuzza, *supra* note __ at 607 (“jury decisions on technologically complex questions of patentability can be unpredictable”).

tensibly) arbitrary case outcomes.¹⁷ And while some level of unpredictability is inevitable, there is reason to think that too much uncertainty is harmful to the innovation system as a whole.¹⁸ Widespread uncertainty about patent validity imposes costs on all innovators,¹⁹ and may chill expensive investments in research and development.²⁰ Uncertainty about validity can also hamper the market for litigation insurance,²¹ and as the cost of litigating a patent case continues to rise,²² uncertainty about the validity of asserted patents can make it more difficult to defend against low-probability nuisance suits. Similarly, on the plaintiff side, where contingent-fee representation has become a more significant part of the patent litigation landscape,²³ being able to make an accurate

-
- 17 John F. Luman III & Christopher L. Dodson, *No Longer a Myth, the Emergence of the Patent Troll: Stifling Innovation, Increasing Litigation, and Extorting Billions*, 18 INTELL. PROP. & TECH. L.J. (2006) (asserting that the outcome of patent cases is inherently unpredictable). There are, however, notable dissenters to the prevailing view. See, e.g., Ted M. Sichelman, *Myths of (Un)Certainty at the Federal Circuit*, 43 LOY. L.A. L. REV. 1161 (2010) (arguing that Federal Circuit jurisprudence is not excessively unpredictable or panel-dependent); Michael J. Mazzeo, Jonathan Hillel & Samantha Zyontz, *Explaining the “Unpredictable”: An Empirical Analysis of U.S. Patent Infringement Awards*, 35 INT’L REV. L. & ECON. 58 (2013) (finding that damages awards are predictable based on characteristics of the patent and litigation).
 - 18 Wagner, *supra* note __ at 2140–41 (arguing that “[p]articularly compelling is the recognition that a patent system characterized by low patent quality sows substantial uncertainty at all levels of the patent system”).
 - 19 Peter S. Menell & Michael J. Meurer, *Notice Failure and Notice Externalities*, 5 J. LEGAL ANALYSIS 1, 5 (2013) (“The proliferation of such rights, as well as the uncertainty regarding their validity and scope, can impose substantial costs on other creators, thereby undermining progress.”).
 - 20 Christopher M. Holman, *Unpredictability in Patent Law and Its Effect on Pharmaceutical Innovation*, 76 MO. L. REV. 645, 648 (2011) (noting a pervasive view in the pharmaceutical industry that patent enforcement has become “uncertain and unpredictable”).
 - 21 See Colleen V. Chien, *Predicting Patent Litigation*, 90 TEX. L. REV. 283, 294–97 (2011); JAMES BESSEN & MICHAEL J. MEURER, *PATENT FAILURE: HOW JUDGES, BUREAUCRATS, AND LAWYERS PUT INNOVATORS AT RISK* 51–54 (2008).
 - 22 See, e.g., Anne S. Layne-Farrer, *The Cost of Doubling Up: An Economic Assessment of Duplication in PTAB Proceedings and Patent Infringement Litigation*, 10 LANDSLIDE 1 (2018) (noting that even for low-stakes patent litigation—with potential damages less than \$10 million—the average cost to litigate to judgment was \$2 million, with half of that amount incurred after discovery).
 - 23 See David L. Schwartz, *The Rise of Contingent Fee Representation in Patent Litigation*, 64 ALA. L. REV. 335 (2012) (detailing the growth of a contingent fee patent litigation practice that is surprisingly diverse across firm types).

forecast of a patent's validity is an essential tool for practitioners making decisions about which cases to take.

In this paper, then, I make a series of contributions to the legal literature. First, I offer a new perspective on patent litigation, demonstrating for the first time that validity decisions are more predictable than previously thought. Second, I shed light on a novel way to measure the quality of a patent—a long-held goal for researchers. Although a full picture of quality is limited by selection concerns, the findings here show that a patent's claim text contains significant information about the likelihood that it will survive a validity challenge. Finally, I contribute to an emerging literature in empirical legal studies that harnesses the potential of automated text analysis to reveal previously undiscovered patterns in legal texts.

The balance of the paper proceeds as follows. In the next part, I discuss the current state of the empirical literature on patent litigation and patent quality. In Part II, I describe the rise of computational methods in legal analysis, and situate this project within the emerging field. In Part III, I describe the data, my empirical strategy and the machine learning tools that I use, and demonstrate that text can be used to more confidently predict validity outcomes in litigation. Part IV concludes with a discussion of the implications of this work for both patent scholarship and litigation, along with some important caveats and limitations.

I. PRIOR STUDIES OF PATENT LITIGATION

Two decades ago, John Allison and Mark Lemley published one of the first rigorous and comprehensive empirical examinations of patent litigation outcomes. They showed that around 46% of challenged patents were found invalid,²⁴ and that non-obviousness and novelty provided the statutory bases for invalidity in the majority of decisions. Perhaps surprisingly, they also found no significant differences in invalidity rates between fields of invention. Since then, empirical studies have blossomed, exploring almost every aspect of litigation. Some of Allison and Lemley's initial findings have proven remarkably robust; others have been overturned with the times.

24 John R. Allison & Mark Lemley, *Empirical Evidence on the Validity of Litigated Patents*, 26 AIPLA Q.J. 185 (1998).

Scholars have since probed a wide variety of potential determinants of outcomes in patent litigation.²⁵ In an early study, two economists test the relationship of a variety of patent attributes (number of claims, forward and backward citations, and patentee portfolio size) with both decisions to litigate and litigation outcomes.²⁶ They find that “win rate outcomes are almost completely independent of observed characteristics of patents and their owners.”²⁷ Legal scholars have paid closer attention to the characteristics of the process itself. For example, Cotropia, Lemley, and Sampat find that characteristics of the patent examination process are predictive: patents that were subject to re-examination are more likely to fail at litigation, all else equal.²⁸ Similarly, Lemley, Li, and Urban consider the role of judicial experience, and note that district court judges who have more experience with patent cases are more likely to find non-infringement—though, interestingly, not invalidity.²⁹ Another notable change from the early studies is a growing discrepancy between industries and technology areas. In a 2012 paper, for example, Allison et al. find that “internet patents” (a subset of software patents) fare remarkably poorly, with overall patentee win rates of just 3% in their decade-long dataset.³⁰

Allison and Lemley, this time with David Schwartz, returned to these questions in a trio of articles that explore patent litigation that commenced at

25 This is a large field, and the highly selective survey here serves to give an overview of the most important findings. For a broad recent summary of the literature, see Ronald Mann & Christopher Cotropia, *Empirical Studies in Patentability*, in 2 RESEARCH HANDBOOK ON THE ECONOMICS OF INTELLECTUAL PROPERTY LAW 281 (Peter Menell & David Schwartz eds., 2019).

26 Jean O. Lanjouw & Mark Schankerman, *Enforcement of Patent Rights in the United States*, in PATENTS IN THE KNOWLEDGE-BASED ECONOMY 145 (Wesley M. Cohen & Stephen A. Merrill eds., 2003).

27 *Id.* at 172.

28 Christopher A. Cotropia, Mark Lemley & Bhaven Sampat, *Do Applicant Patent Citations Matter?*, 42 RES. POL’Y 844 (2013).

29 Mark Lemley, Su Li & Jennifer Urban, *Does Familiarity Breed Contempt Among Judges Deciding Patent Cases?*, 66 STAN. L. REV. 1121 (2014).

30 John R. Allison, Emerson H. Tiller, Samantha Zyontz & Tristan Bligh, *Patent Litigation and the Internet*, 2012 STAN. TECH. L. REV. 3. See also Lemley et al., *supra* note __ at 1144–50.

the end of the 2000s.³¹ Using a careful hand-coding of all merits decisions in cases initially filed in 2008 and 2009, they find evidence that jurisdiction makes a difference to case outcomes, as do some characteristics of the lawyers and judges involved.³² They confirm prior wisdom that certain kinds of patents are less likely to be found invalid—particularly those relating to pharmaceuticals—and also observe that patentees lose nearly 75% of cases that are litigated to judgment.³³ In the third study, the authors analyze the role of non-practicing entities (NPEs) in litigation, finding that operating companies generally experience better outcomes in litigation than pure patent-assertion entities.³⁴ However, the authors caution that much of the effect is down to software patents, where NPEs are overrepresented as litigants, but no patentee fares well. Beyond this study, the role of NPEs has been the subject of significant recent scholarly attention,³⁵ with a mixed picture being drawn. Michael Risch, for example, finds that claims asserted by NPEs are settled at much higher rates than those by practicing entities, and that those that make it to a merits hearing see their patents held invalid at higher rates.³⁶

Finally, in work most closely related to this Article, several papers have used litigation outcomes to explore the question of patent quality. Jonathan

31 John R. Allison, Mark Lemley & David L. Schwartz, *Understanding the Realities of Modern Patent Litigation*, 92 TEX. L. REV. 1769 (2014); John R. Allison, Mark Lemley & David L. Schwartz, *Our Divided Patent System*, 82 U. CHI. L. REV. 1073 (2015); John R. Allison, Mark Lemley & David L. Schwartz, *How Often Do Non-Practicing Entities Win Patent Suits?*, 32 BERKELEY TECH. L.J. 235 (2017).

32 *Id.*, *Understanding the Realities*, at 1799.

33 Allison et al., *Divided Patent System*, *supra* note __ at 1112.

34 Allison et al., *Non-Practicing Entities*, *supra* note __. Non-practicing entities are often referred to as “patent-assertion entities” (PAEs) or “patent trolls.” The umbrella terms represent a range of entities, including universities, individual inventors, and failed startups. Of particular note are patent-holding companies who purchase patents and subsequently assert them in litigation, but who are not themselves involved in research and development or other innovative activities. There is some controversy over the social costs and benefits of such entities. *See, e.g.*, Christopher A. Cotropia, Jay P. Kesan & David L. Schwartz, *Unpacking Patent Assertion Entities (PAEs)*, 99 MINN. L. REV. 649 (2014). Regardless of one’s position on this debate, from the perspective of litigation, there seems good reason to expect that assertion entities would behave in ways that are systematically different to other types of patent owner.

35 *See, e.g.*, Shawn Miller et al., *Who’s Suing Us? Decoding Patent Plaintiffs Since 2000 with the Stanford NPE Litigation Dataset*, 21 STAN. TECH. L. REV. 234 (2018).

36 Risch, *supra* note __.

Ashtor develops a measure of the “informational content” of a patent, which combines various facially observable characteristics of the patent: number of claims, length of written description and abstract, number of inventors, and so on.³⁷ First, he creates a set of weights by relating these variables to the number of forward citations (that is, citations to the patent by *future* applicants—a traditional measure of patent value for economists). He then tests that weighted index against validity decisions and other factors, and finds that patents with more “informational content” are slightly less likely to be invalidated.³⁸

Meanwhile, Ronald Mann and Marian Underweiser analyze validity decisions at the Federal Circuit.³⁹ Operationalizing a patent’s quality as its likelihood of being held valid in federal court, the authors show that features of a patent’s prosecution history are correlated with judicial outcomes. For example, they find a patent is more likely to be found invalid on appeal if its application involved many rounds (versus those which were approved more quickly) or if it has high rates of prior-art references that were added by the examiner (versus those disclosed by the applicant herself).⁴⁰ A particularly interesting feature of this paper is its text-based measure of the “alignment” between the patent’s written description and its claims.⁴¹ Although this metric has a somewhat limited application, because the alignment between a patent’s written description and its claims are only relevant to a small subset of decisions about validity,⁴² it is an intriguing early approach to incorporating patent text directly into a model of litigation.

So, with all these studies, what do we know? While different scholars have emphasized quite different characteristics, we can draw some general conclusions. Certain factors do seem reliably predictive of litigation outcomes. There are clear industry effects—for example, pharmaceutical patents are struck down at much lower rates than others, all else equal.⁴³ Jurisdiction effects, too, have

37 Jonathan H. Ashtor, *Does Patented Information Promote the Progress of Technology?*, 113 NW. U. L. REV. 943 (2019).

38 *Id.* at 980–85.

39 Mann & Underweiser, *supra* note ____.

40 *Id.* at 20–21.

41 See *infra* Part II for a general discussion of text-based “similarity” measures.

42 35 U.S.C. § 112 (2011). See John R. Allison & Lisa Larrimore Ouellette, *Patent Definiteness and Disclosure*, 65 DUKE L.J. 609, 617–19 (2016) (describing the “enablement” and “written description” requirements).

43 Allison et al., *Divided Patent System*, *supra* note ____ at 1114.

been shown to be important, though the specifics are somewhat mixed.⁴⁴ On the other hand, studies have generally found that forward citations provide little information about invalidity (or infringement) decisions.

Most importantly for the project at hand, these studies generally find that observable attributes of the patent itself provide little information about the success or failure of any given patent at trial. As Risch summarizes: “. . . predicting which patents were invalidated had more to do with case-specific factors, such as the number of defendants, than with objectively measurable patent quality indicators.”⁴⁵ The other feature that unites nearly all of these models, however, is that they are mostly not very successful at explaining the data. Despite significant data collection and a host of theoretically well-justified covariates that have been tested across multiple different studies, the explanatory power of the regression models is consistently poor. As Allison, Lemley, and Schwartz conclude: “The pseudo R²s in our regressions . . . are very low, revealing that most of the variation in patent litigation outcomes is not predictable, at least based upon the extensive variables we captured.”⁴⁶

All this prompts several questions. Are patent litigation outcomes truly unpredictable? Are we missing potentially valuable sources of information? In the next Part, I describe the computational methods that offer new ways to answer these questions.

44 *Id.* For example, the Eastern District of Texas has long been thought to be patentee friendly. See Ofer Eldar and Neel Sukhatme, *Will Delaware Be Different? An Empirical Study of TC Heartland and the Shift to Defendant Choice of Venue*, 104 CORNELL L. REV. 101, 110–18 (2018). But see Lemley et al., *supra* note __ (finding that, for patent owners, D. Del. has generally been a *more* favorable venue than E.D. Tex., after conditioning on judicial experience).

As Eldar and Sukhatme describe, a recent Supreme Court decision limits plaintiff forum selection in patent suits and appears likely to change this effect in the future. Its full effect remains to be seen. See *TC Heartland, LLC v. Kraft Food Brands Grp. LLC*, 137 S. Ct. 1514 (2017) (requiring that patent infringement suits be brought in the defendant company’s district of incorporation).

45 See, e.g., Risch, *supra* note __ at 131.

46 Allison et al., *Modern Litigation*, *supra* note __ at 1799. Interestingly, one feature of patent litigation—damages awards—does seem to be reasonably predictable, based on observable features of the patent and litigants. See Mazzeo et al., *supra* note __ at 67.

II. COMPUTATIONAL ANALYSIS IN LAW

Legal analytics is a growing subfield within law generally, and intellectual property specifically. The use of machine learning tools for text analysis has dramatically increased across the social sciences, but the tools hold particular promise in empirical legal studies—a field where we have large troves of text that, although highly salient, have traditionally been difficult to study empirically.⁴⁷ A set of machine learning approaches, referred to as natural language processing (NLP), enables rapid and flexible processing of large quantities of text, which allows researchers to detect new—and sometimes counterintuitive—patterns in various kinds of legal documents.⁴⁸

Within intellectual property scholarship, scholars have begun to employ creative uses of NLP. One common goal is to conceptualize and measure underlying legal doctrines. For example, a team of computer scientists and business school faculty have developed a measure of “novelty” that is based on the first entry of a word into the patent corpus.⁴⁹ Similarly, Jeffrey Kuhn and Neil Thompson propose a measure of patent scope that counts the number of words in the first claim of a patent.⁵⁰ The development of “similarity” measures is another active area of research. These metrics are based on the geometric distance between bodies of text, which are represented as (vast) vectors of words.⁵¹ For example, Laura Pedraza-Fariña and Ryan Whalen compute scores of the mathematical similarity between the texts of thousands of patent pairs, and use those scores to derive a network-based measure of the non-obviousness doc-

47 For a broad introduction to this emerging field, see *LAW AS DATA: COMPUTATION, TEXT, AND THE FUTURE OF LEGAL ANALYSIS* (Michael A. Livermore & Daniel N. Rockmore eds., 2019).

48 See generally Jens Frankenreiter & Michael Livermore, *Computational Methods in Legal Analysis*, 16 ANN. REV. L. & SOC. SCI. 39 (2020).

49 Benjamin Balsmeier et al., *Machine Learning and Natural-Language Processing on the Patent Corpus: Data, Tools, and New Measures*, 27 J. ECON. & MGMT. 535 (2018).

50 Intuitively, more words narrow the scope of a claim because—in general—a competing product must practice every element of the claim to infringe. Jeffrey M. Kuhn & Neil C. Thompson, *How to Measure and Draw Causal Inferences with Patent Scope*, 26 INT’L J. ECON. BUS. 5 (2019).

51 For a thorough explanation of this approach as applied to patents, see Kenneth A. Younge and Jeffrey M. Kuhn, *Patent-to-Patent Similarity: A Vector Space Model* (working paper, 2016), <https://dx.doi.org/10.2139/ssrn.2709238>.

trine.⁵² Jonathan Ashtor uses a similar paired approach in an effort to create a synthetic, algorithmic alternative to traditional citations.⁵³ Finally, Janet Freilich uses a computational analysis of 25,000 patents (primarily non-litigated) to show that as many of 25% of the nouns in patent claims are “ancillary,” in the sense that they appear less than twice in the patent’s specification.⁵⁴

My goal in this paper is somewhat different, and has more in common with traditional empirical studies of litigation. Rather than recovering algorithmic measures of particular patent doctrines, or trying to extract semantic meaning from the claim language, my interest is in exploring the variability of patent litigation outcomes. I treat the question as one of prediction, and ask: given information about a patent’s claim text, how well can we predict its likely validity? This is a different application of computational analysis to prior work in IP,⁵⁵ but the underlying tools are similar.

Typically, empirical studies of litigation use regression analysis of some kind in order to describe the correlates of an outcome or to draw causal infer-

52 Laura G. Pedraza-Fariña & Ryan Whalen, *A Network Theory of Patentability*, 87 U. CHI. L. REV. 63 (2020). Whalen in particular is at the leading edge of network-based computational approaches to innovation law. See, e.g., Ryan Whalen, *Boundary Spanning Innovation and the Patent System: Interdisciplinary Challenges for a Specialized Examination System*, 47 RES. POL’Y 1334 (2018); Ryan Whalen, *Legal Networks: The Promises and Challenges of Legal Network Analysis*, 2016 MICH. ST. L. REV. 539.

53 See Jonathan Ashtor, *Investigating Cohort Similarity as an Ex Ante Alternative to Patent Forward Citations*, 16 J. EMPIRICAL LEGAL STUD. 848 (2019).

54 This is important because, in the parlance of patent law, the elements of a claim must be enabled by and described in the specification. 35 U.S.C. § 112. Janet Freilich, *Patent Clutter*, 103 IOWA L. REV. 925 (2018). This finding is particularly interesting for the present study, given that it suggests that a significant portion of patent claim language is in fact irrelevant to the claimed invention. That said, to the extent that an excess of “ancillary” language might drive invalidation decisions in district court (which are not the focus of Freilich’s study), it should do so primarily on the basis of Sec. 112, which comprise a minority of the decisions in the 2009–10 window.

55 There is, however, some intriguing early work in computer science which incorporates text features in a model of selection into litigation. See Papis Wongchaisuwat, Diego Klabjan & John McGinnis, *Predicting Litigation and Time to Litigate*, PROC. 16TH ANN. CONF. INT’L ASS’N ARTIFICIAL INTELLIGENCE & L. 257 (2017).

ences.⁵⁶ By contrast, the use of machine learning approaches to prediction and classification is quite novel in empirical legal scholarship.⁵⁷ One technique that has been deployed fruitfully is a classification tree (a type of decision tree).⁵⁸ Unlike a traditional logistic regression, a decision tree doesn't require any *ex ante* information about the structure of the relationship between the outcome and the predictors of interest. Instead, it splits the outcome data into a set of endpoints (or leaves, to borrow from the tree metaphor) which are connected

56 Traditionally, researchers specify a fixed (usually linear) model of the relationship between the outcome and a set of inputs—for example, between the plaintiff win-rate and the jurisdiction hearing the case. (Usually this supposed relationship is provided by a theoretical model or some other domain-specific knowledge.) They then collect a sample of data, feed it to the model to produce estimates, and assess how well the model performs (that is, how confident they are) using various goodness-of-fit tests. By contrast, the “algorithmic” culture of much modern statistics eschews theory, treating the structure of the relationship between inputs and outcomes as *a priori* unknown. Instead, researchers allow a statistical algorithm to learn the best model, and then validate performance by challenging the calibrated algorithm to predict outcomes in new, unseen data. This is a deep methodological divide. For a lively and classic discussion of the “two cultures” within (and beyond) statistics, see Leo Breiman, *Statistical Modeling: The Two Cultures*, 16 STATISTICAL SCI. 199 (2001).

57 The critical literature on machine learning in doctrinal and theoretical scholarship, on the other hand, is young but voluminous. Scholars at the intersection of law and technology, discrimination, and the criminal legal system have been particularly trenchant critics of the rise in algorithmic decision making and its troubling social implications. See, e.g., Solon Barocas & Andrew Selbst, *Big Data's Disparate Impact*, 104 CALIF. L. REV. 671 (2016) (noting the disparate impact of algorithmic approaches to employment discrimination which are based on historic patterns of prejudice); Rashida Richardson, Jason Schultz & Kate Crawford, *Dirty Data, Bad Predictions: How Civil Rights Violations Impact Police Data, Predictive Policing Systems, and Justice*, 94 NYU L. REV. ONLINE 192 (arguing that widely deployed predictive policing algorithms learn from, and help to perpetuate, past bias); Danielle Keats Citron & Frank Pasquale, *The Scored Society: Due Process for Automated Predictions*, 89 WASH. L. REV. 1 (2014) (arguing that individuals should be afforded meaningful opportunities to challenge the harmful impacts of algorithmic risk scores); FRANK PASQUALE, *THE BLACK BOX SOCIETY: THE SECRET ALGORITHMS THAT CONTROL MONEY AND INFORMATION* (2015). But see Talia B. Gillis & Jann L. Spiess, *Big Data and Discrimination*, 86 U. CHI. L. REV. 459 (2019) (arguing that machine-derived decision rules can provide a framework for testing algorithmic discrimination); W. Nicholson Price & Arti Rai, *Clearing Opacity Through Machine Learning*, 106 IOWA L. REV. 775 (2021) (arguing that machine learning has the potential to shed light on non-intuitive, complex systems in the biomedical sciences).

58 See Jonathan P. Kastellec, *The Statistical Analysis of Judicial Decisions and Legal Rules with Classification Trees*, 7 J. EMPIRICAL LEGAL STUD. 202 (2010).

by branches. These branches represent the flow of a set of if-then logical decision rules.⁵⁹

A landmark study uses decision trees to classify the outcome of Supreme Court cases, based on characteristics of each case and its procedural history.⁶⁰ The model was remarkably successful, correctly predicting 75% of case outcomes in the 2002 term—and significantly outperforming the panel of expert lawyers and academics who attempted the task in parallel. In the IP context, Tammy Cowart, Roger Lirely, and Sherry Avery demonstrate the potential of classification trees to improve the clarity of models of patent litigation.⁶¹ However despite their flexibility, decision trees have a tendency to “overfit” to the data, consequently failing to predict future cases. In more recent work, a legal analytics research team use a methodological extension of decision trees, known as random forests, to forecast the resolution of employment law litigation using observable characteristics of the cases.⁶² The latter project is also interesting for its incorporation of case features that are extracted directly from the text of docket sheets. This is an increasingly common approach.

Beyond the IP examples mentioned above, applications of automated text analytics are finding increased use in empirical legal studies across domains. For example, Julian Nyarko, David Pozen, and Eric Talley use the text of the Congressional Record to demonstrate the changing partisanship of “constitutional” political language over two centuries.⁶³ Jonathan Choi uses the text of

⁵⁹ *Id.* at 209–13.

⁶⁰ Theodore W. Ruger, Pauline T. Kim, Andrew D. Martin & Kevin M. Quinn, *The Supreme Court Forecasting Project: Legal and Political Science Approaches to Predicting Supreme Court Decisionmaking*, 104 COLUM. L. REV. 1150 (2004).

⁶¹ Tammy Cowart, Roger Lirely & Sherry Avery, *Two Methodologies for Predicting Patent Litigation Outcomes: Logistic Regression Versus Classification Trees*, 51 AM. BUS. L.J. 843 (2014) (highlighting the interpretative and flexibility advantages of classification tree analysis).

⁶² Charlotte S. Alexander, Khalifeh al Jadda, Mohammad Javad Feizollahi & Anne M. Tucker, *Using Text Analytics to Predict Litigation Outcomes*, in LAW AS DATA: COMPUTATION, TEXT, AND THE FUTURE OF LEGAL ANALYSIS 271–308 (Michael A. Livermore & Daniel N. Rockmore eds., 2019).

⁶³ David Pozen, Eric Talley & Julian Nyarko, *A Computational Analysis of Constitutional Polarization*, 105 CORNELL L. REV. 1 (2020). Scholars have also applied NLP to judicial writing. For example, Elliott Ash and various collaborators use the text of written opinions to explore and classify judicial ideology. See, e.g., Elliott Ash & Daniel L. Chen, *What Kind of Judge is Brett Kavanaugh?*, 2018 CARDOZO L. REV. DE NOVO 70; Carina I. Hausladen, Marcel H. Schubert & Elliott Ash, *Text Classification of Political Ideology Labels in Judicial Opinions*, 62 INT. REV. L. ECON. 1 (2020).

IRS publications and tax-court opinions to show that tax courts have become more textualist over time, even as the IRS has adopted a more purposivist approach to interpretation in rulemaking.⁶⁴ And in the corporate context, Adam Badawi demonstrates that the text of first-filed complaints can be used to predict the outcome of securities litigation.⁶⁵

The approach I use here is in line with these recent applications of NLP to law. I treat the claim text of a patent as a proxy for a measure of quality, and explore its relationship with invalidity decisions. In brief, my strategy is to break down the patent claims into their constituent words, and use the relative prevalence of those words in a given patent (with some adjustments) to predict the likelihood of success or failure in a validity adjudication. In Part III, I describe these steps in more detail, and introduce the dataset of litigation that I use.

III. PREDICTING INVALIDITY

A. Data sources

The litigation data used in this analysis come from a set of recent studies of patent litigation conducted by John Allison, Mark Lemley, and David Schwartz. I refer the reader to those authors—in particular the first paper in their trilogy, *Understanding the Realities of Modern Patent Litigation*—for a full description of their careful hand-coding procedure.⁶⁶

⁶⁴ Jonathan H. Choi, *An Empirical Study of Statutory Interpretation in Tax Law*, 95 NYU L. REV. 363 (2020).

⁶⁵ Adam Badawi, *How Informative is the Text of Securities Complaints?* (unpublished manuscript) (October 1, 2019) (on file with author). See also Adam Badawi & Elisabeth de Fontenay, *Contractual Complexity in Debt Agreements: The Case of EBITDA* (Duke Law School Public Law & Legal Theory Series No. 2019-67, 2019), <https://dx.doi.org/10.2139/ssrn.3455497> (using NLP to analyze EBITDA definitions in credit agreements).

⁶⁶ Allison et al., *Understanding the Realities*, *supra* note _____. For a fuller discussion of some of the complexities of coding patent litigation data, see Jason Rantanen, *Empirical Analyses of Judicial Opinions: Methodology, Metrics, and the Federal Circuit*, 49 CONN. L. REV. 227 (2016).

Table 1: Summary of outcomes (validity decisions in case-patent pairs)

	<i>Ruled not invalid</i>	<i>Ruled invalid</i>	<i>Total</i>
<i>Observations</i>	205	159	367

The dataset includes every outcome from a merits decision in patent cases filed in U.S. district court in 2008 and 2009.⁶⁷ The unit of analysis in the data is the patent-case combination. Patents are often litigated more than once, and if so they would show up appear multiple times in the data.⁶⁸ Similarly, if a case includes judgments on multiple discrete patents, the disposition of each patent is coded individually.

The key outcome variable is the result of an adjudication of invalidity (“valid” or “invalid”).⁶⁹ This includes all doctrinal grounds (eligible subject matter, novelty, obviousness, and so on), and decisions that occur at any stage in the litigation.⁷⁰ If a district court judgment was subsequently overturned on appeal, then the invalidity outcome reflects that. Importantly, the outcome includes judgments of partial (in)validity. In other words, if only a subset of the claims was at issue, but all of those were found invalid, the outcome is coded as a finding of “invalidity” on the patent.

⁶⁷ The cases under study commenced over a decade ago. Given the long pendency of litigation, some cases did not reach a definitive conclusion until as late as 2015, but it is nevertheless reasonable to suspect that at some of the specific patterns detected in this period would not hold true today. Patent law has experienced significant doctrinal and statutory upheaval in the past ten years. I discuss the generalizability of the core result in Part IV, *infra*.

⁶⁸ A patent (or individual patent claim) that is found invalid by a court is permanently lost by the patentee. A patent that has been subject to an affirmative finding of “no invalidity” or, more commonly, simply not ruled invalid remains at risk in future litigation. See Mark A. Lemley, *The Fractioning of Patent Law*, in *INTELLECTUAL PROPERTY AND THE COMMON LAW* 504 (Shyamkrishna Balganesh ed., 2013) (noting the asymmetry in which a patentee has to win every time a patent is challenged).

⁶⁹ Strictly speaking, courts find that a patent is not invalid.

⁷⁰ I am essentially averaging over different grounds for invalidity, but one might suspect that this approach would work better for some than others. For example, § 112 challenges on the basis of inadequate written description or lack of enablement are primarily “internal” to the patent, whereas judgments about novelty and obviousness implicate, to various extents, the external prior art, which is not explicitly included in the model. Unfortunately, I do not have sufficient data to train models for the individual statutory provisions, but this is an interesting area for future research.

In one important respect, I treat the data differently to the previous authors. Fifteen cases (out of 216 unique cases in total) involved a split decision on a particular patent, in which some claims were found invalid while others were upheld. Although this is not a modeling problem in principle, it does make it difficult to compare the text of these patents to the other examples, all of which are coded at the case-patent level (even if only a subset of claims was at issue). To ensure all the outcomes are comparable, I drop these observations from the analysis.

B. Preparing the text for analysis

I supplement the litigation data with the full claim text of each patent, which I obtain from the USPTO.⁷¹ Because the outcomes are coded at the level of the patent rather than the individual claims, I concatenate the text of all claims into a single “document” for each patent.⁷² Then, to make the text amenable to quantitative analysis, I put it through a series of typical “preprocessing” steps.⁷³

First, I break down the patent into its constituent words, convert the words to lowercase, and remove all punctuation, numbers, and other symbols. Note that in the process, I discard all information about word order and syntactical structure. This is a deliberately simple but important choice, which warrants further explanation. This representation of a text (commonly known as a “bag of words”) may appear quite reductive, but prior research in computational social science has repeatedly shown that these parsimonious representations

71 USPTO Patents View, *Data Download* (Dec. 31, 2019), <https://www.patentsview.org/download>. The post-processed claim text will be available in an online appendix.

72 In other words, I use the combined text of every claim in a given patent to predict a decision on that patent, even when a court is asked to consider only a subset of the claims. This obviously risks introducing a certain amount of noise—in the sense of irrelevant data—into the model. On the other hand, this is a conservative choice: if the words in non-adjudicated claims are not relevant to the validity decision, they should not contribute useful information to the prediction model.

73 The method I set out is a common approach to handling text of this kind. See, e.g., Eric Talley & Drew O’Kane, *The Measure of a MAC: A Machine-Learning Protocol for Analyzing Force Majeure Clauses in M&A Agreements*, 168 J. INSTITUTIONAL & THEORETICAL ECON. 181 (2012); Gabriel Rauterberg & Eric Talley, *Contracting Out Of The Fiduciary Duty Of Loyalty: An Empirical Analysis Of Corporate Opportunity Waivers*, 117 COLUM. L. REV. 1075 (2017).

are remarkably rich “summaries” of documents.⁷⁴ Nevertheless, the approach may still seem counterintuitive in *this* area. Patent language, after all, is arcane, technical, and interpreted in a self-consciously contextual way.⁷⁵ New tools in NLP offer increasingly sophisticated ways to model and represent documents, many of which maintain explicit information about semantic relationships and word co-occurrence.⁷⁶ However, these more elaborate models come at a cost: they are highly opaque. In contrast, the traditional bag-of-words approach has the signal advantage of being very transparent. In particular, it allows us to observe which words are important to the prediction and which are not. And as I show in Part III, *infra*, I find that the simple approach performs extremely well in this context.

Next, a common (though not universal) technique is to remove “stop-words” from the corpus. These are short phrases such as “a,” “the,” “or,” and “to,” that are very common in language but do not usually convey any information about the meaning of the text. Again, however, the patent context is somewhat unusual. Many traditional stopwords are themselves terms of art, and have been subject to significant litigation about their proper interpretation.⁷⁷ The contested nature of these words may well add complexity to more elaborate semantic models of patent text, but for our purposes, the choice about whether to include them is essentially an empirical question—which is more predictive? In Table 3, below, I present results both with and without stopwords included.⁷⁸

74 Justin Grimmer & Brandon M. Stewart, *Text as Data: The Promise and Pitfalls of Automatic Content Analysis for Political Texts*, 21 POL. ANALYSIS 1, 6–7 (2013); *supra* notes 62–64 and accompanying text.

75 Patent doctrine explicitly instructs courts to interpret the meaning of terms and claims in the context of the whole patent. *See, e.g.*, Phillips, 415 F.3d at 1313 (“[p]roperly viewed, the “ordinary meaning” of a claim term is its meaning to the ordinary artisan after reading the entire patent”).

76 In particular, modern NLP applications often use word or document “embeddings,” which represent text as low-dimensional numerical vectors that incorporate information about the co-occurrence rates of words within the corpus. *See, e.g.*, Elliott Ash, Daniel L. Chen & Arianna Ornaghi, *Gender Attitudes in the Judiciary: Evidence from U.S. Circuit Courts*, (working paper, 2020), http://elliottash.com/wp-content/uploads/2019/11/200205_Ash-Chen-Ornaghi.pdf.

77 *See, e.g.*, Burk & Lemley, *supra* note __ at 1751–1753.

78 I use the Snowball stopword list. *See* Stopwords 2.0, <https://stopwords.quanteda.io>. The full list is 175 words long, and includes pronouns, prepositions, and conjunctive words, as well as common contractions (“shan’t,” “won’t,” and so on).

Finally, to avoid overfitting the model to idiosyncratic terms, I filter the dataset to the 300 words that appear most commonly across the entire corpus (that is, from the set of all words that appear at least once in any patent).⁷⁹ Using this filtered set of words, I count the number of times each term appears in each patent. The output of this process is set of vectors, one per patent, which contain a count of the number of times that each of the 300 terms appears in that patent. To ensure that the results are not driven primarily by the length of a patent, I normalize each word by dividing its frequency by the sum of all word counts in that patent.⁸⁰

Taken together, all these steps result in a grid of normalized word frequencies, with a row for each patent, and a column for every word in the corpus. This representation of the text is known as a document-term matrix (“DTM”). Table 2 shows the first few rows and columns of the transformed DTM for this patent corpus.⁸¹ To create the final analysis dataset, I merge the DTM with the litigation outcomes and metadata based on the involved patent in each adjudication.⁸²

C. Methodological approach

To accommodate the textual data, I adopt a somewhat different empirical strategy to prior work in this area. The typical approach is to use a logistic regression, which is a generalization of linear regression designed for binary (rather than continuous) outcomes. A traditional linear model is problematic

79 The choice of 300 terms is driven by the relatively small size of the training and test data that I have available. Without sufficient training data, increasing the number of words used in building the model is likely to reduce its out-of-sample performance. In Appendix Part A, *infra*, I show robustness checks for different vector sizes. In summary, while 300 words performs well, the decrease in performance for other choices is relatively slight, and the core results are not driven by this choice.

80 Of course, we can also control for the length of the patent in the analysis. This step simply ensures that the effect of the words themselves and the overall length are considered separately.

81 The grid is largely populated with zeroes; in statistical language, it is extremely sparse. Data in this form can be computationally challenging to work with, and for this reason, researchers often take steps to further reduce the dimensionality of the data, such as computing the singular value decomposition of the matrix. See, e.g., Talley & O’Kane, *supra* note __. For purposes of this project, I am particularly interested in preserving information about the distinct words, so I leave the DTM as is.

82 Note that in the final dataset, each row of the DTM will appear as many times as its patent was litigated.

Table 2: Transformed document-term matrix

	<i>absorbent</i>	<i>acceptable</i>	<i>access</i>	<i>accordance</i>	<i>according</i>	<i>acid</i>	...
<i>patent1</i>	0	0	0	0	0	0	...
<i>patent2</i>	0	0	0	0	0	0	...
<i>patent3</i>	0	0	0	0	0	0	...
<i>patent4</i>	0	0.031	0	0	0.062	0	...
<i>patent5</i>	0	0	0	0	0	0	...
<i>patent6</i>	0	0	0	0	0.048	0	...
⋮	⋮	⋮	⋮	⋮	⋮	⋮	

here for two reasons. First, textual data is inherently very “high-dimensional”—in most cases I have nearly as many variables as observations, which renders linear and logistic regression unreliable at best, and impossible at worst. Second, all these variables (text and otherwise) interact with each other in highly complex, non-linear ways, which are not generally easy to anticipate or specify *ex ante*. Machine learning offers promising solutions to both these problems.

To model the data, I use a statistical tool known as a “random forest.”⁸³ Random forests are a generalization of the basic decision tree introduced above.⁸⁴ Decision trees are easily interpretable and fit the data in a very flexible way, but they also tend to *overfit* and are quite sensitive to small perturbations in the data, such that the output can be quite variable from tree to tree. To circumvent this, a random forest grows thousands of trees in parallel, using randomly chosen subsets of data and covariates each time. The algorithm then combines their predictions by aggregating the “votes” of each individual tree, resulting in a better prediction from the whole ensemble together.

⁸³ See Leo Breiman, *Random Forests*, 45 MACHINE LEARNING 5 (2001). For the canonical textbook treatment, see TREVOR HASTIE, ROBERT TIBSHIRANI & JEROME FRIEDMAN, *THE ELEMENTS OF STATISTICAL LEARNING* 587–601 (2d ed., 2009). Note that random forests are only one of a broad class of “machine learning” algorithms. They perform well in this application, but we should not expect them to dominate in every setting. I argue that the results in this study suggest that we should make greater use of sophisticated and flexible classification tools in empirical legal research, but *not* that random forests are uniformly superior to any other approach.

⁸⁴ See footnote ___ and accompanying text.

For each model, I follow a common procedure to train the algorithm and produce predictions.⁸⁵ First, I split the observations at random into a “training set” (90% of the data) and a “test set” (10%). Then, I estimate a statistical model that best fits the training data.⁸⁶ In each case, the outcome of interest is whether or not a patent is held valid. Then, using this trained model, I predict the outcomes in the previously unused test set.⁸⁷ I repeat this procedure nine more times, each time using a distinct set of test data, such that I end up with an “out of sample” prediction of invalidity or validity for every observation in the data.⁸⁸ We can then compare these predictions to the hand-coded “true” outcomes in order to assess how well our algorithm performs at predicting court decisions that it has not previously encountered.

To evaluate the models, I use two common measures of performance. The first is “accuracy”—or equivalently, correct classification rate (CCR)—which is simply the number of “successes” divided by the total number of predicted outcomes. In other words, if a model correctly classified 75 out of 100 outcomes and got the other 25 wrong, its CCR would be 75%. The CCR is useful and intuitive, but comes with an important caveat. In practice, the model produces a probabilistic, rather than definitive, prediction of whether a patent will be found invalid. By default, the CCR assumes that 50% is the correct cutoff—that is, any probability above 50% indicates (in our case) invalidity, and anything below 50% indicates validity. In some cases—and especially when the data are evenly balanced between the two classes—this can be a reasonable assumption. However, there are sometimes reasons to prefer a different cutoff, and so it is

85 See, e.g., Pozen et al., *supra* note __ at 31.

86 Note that in principle this could be *any* statistical model, from a simple linear regression to a complex neural network.

87 Traditional measures of statistical performance assess how well a model performs “within sample.” Perhaps the most well-known example is R-squared, which is used in classical linear regression to measure the proportion of variation in the data that is explained by a given model. (There are analogous metrics for logistic and other generalized models.) However, optimizing for in-sample performance risks overfitting the data—that is, producing a model that works well in one context or in a particular sample, but is not readily generalizable. Applying an algorithm to unseen data provides a more challenging and realistic assessment of its predictive power. See generally GARETH JAMES, DANIELA WITTEN, TREVOR HASTIE & ROBERT TIBSHIRANI, *AN INTRODUCTION TO STATISTICAL LEARNING* 29–37 (2013).

88 This procedure is known as k-fold cross-validation, where “k” is equal to the number of splits of the data. Because I have relatively few observations, I use ten folds to ensure a reasonable number of cases in each training set.

common to report a more flexible score known as the “area under the curve” (AUC).⁸⁹

D. Results

I begin with a high-level test: how predictive are the words of the claims on their own? It is important to have a baseline against which to compare the performance of our classification algorithm. An obvious possibility is simply to predict outcomes at random—that is, to assign a 50% chance of invalidity to every case. (Intuitively, a predictive model that correctly classifies fewer than 50% of outcomes is worse than guessing.) Alternatively, a more conservative choice is to assign the *most common* outcome to every case. In the dataset at hand, around 56% of patents were actually held not invalid, so a naive classifier that simply predicted that every single case would end in a finding of validity would be correct 56% of the time. Clearly, neither of these strategies produce insightful predictions—like a stopped clock, they are both right some of the time—but they provide a floor for the quality of any model.

Second, it is important to compare the performance of the new approach to a more traditional, theory-driven model. As I discuss above, empirical scholars of patent litigation have tested a huge range of potential correlates of validity, and there is no universally agreed-upon set of controls. I construct an arbitrary “canonical” logistic regression model, which includes controls for a range of attributes that are generally thought be related to validity outcomes or which

⁸⁹ Rather than assuming a 50% threshold, AUC varies the cutoff from 0% to 100%, and computes the true positives (correct prediction of invalidity) and false positives (incorrect prediction of invalidity) at each point. We plot each of these pairs of values on a graph, and calculate the area under the resulting curve. An AUC of 0.5 indicates a completely uninformative classifier (that is, true and false positives are equal at every possible cutoff), while an AUC of 1 reflects a model that perfectly discriminates. An example of these curves is shown in the Appendix.

have proven informative in earlier studies. I include metadata related to the patent, litigation, technology class, jurisdiction, and type of assertion entity.⁹⁰

The main results, presented in Table 3, are striking. The random forest model correctly predicts 71.1% of validity decisions, using only the text of the patent claims. Meanwhile, the theory-driven model achieves an accuracy of just 62% on the same data. In other words, the text-only model achieves an eight percentage-point improvement on the more complex theoretical model.⁹¹ Interestingly, the inclusion of stopwords (as described above) slightly increased the

Table 3: Overall performance

	Guess the most common outcome	“Canonical” logistic regression	Text-only random forest	
			<i>Without stopwords</i>	<i>With stopwords</i>
Correctly classified	56.3%	62.0%	68.2%	71.1%
AUC	-	0.68	0.78	0.78

90 The full logistic specification uses the following covariates:

(1) Patent metadata: foreign origin; adjusted number of citations; total prior art references; and age at filing of the litigation.

(2) Litigation metadata: number of patents asserted in the suit and the number of defendants.

(3) Dummies for five primary technology areas (one excluded), and an indicator for patents that are implicated in abbreviated new drug applications. I adopt the technology areas used in the original data: mechanical, electronics, optics, biotechnology, chemical, and software.

(4) Mutually exclusive dummies for individual, failed startup, university, and patent-assertion entity (operating company is the excluded category).

(5) Dummies for three of the most important jurisdictions—the Eastern District of Texas, the Northern District of California, and the District of Delaware (all others excluded for statistical power).

91 Note also that the text-only model performs well despite being “handicapped.” Several patents appear multiple times in the analysis, but the model can only produce one prediction for the text of a given patent. For example, if a patent is upheld twice before ultimately being found invalid in a third case, the model will be wrong at least once (if it predicts validity) or twice (if it predicts invalidity). In the results reported in Table 5, *infra*, I relax this constraint by including additional case-level variables.

performance of the model, suggesting that they do help to discriminate between valid and invalid patents.

Table 4 shows the breakdown of correct and incorrect predictions for the text model. It correctly identified 151 of the true decisions in favor of the patentee (74%), but incorrectly predicted the remaining 54. The model performs slightly worse for findings of invalidity, correctly predicting 97 out of 159 decisions (61%).⁹²

Table 4: True outcomes versus predictions (text-only model, stopwords removed)

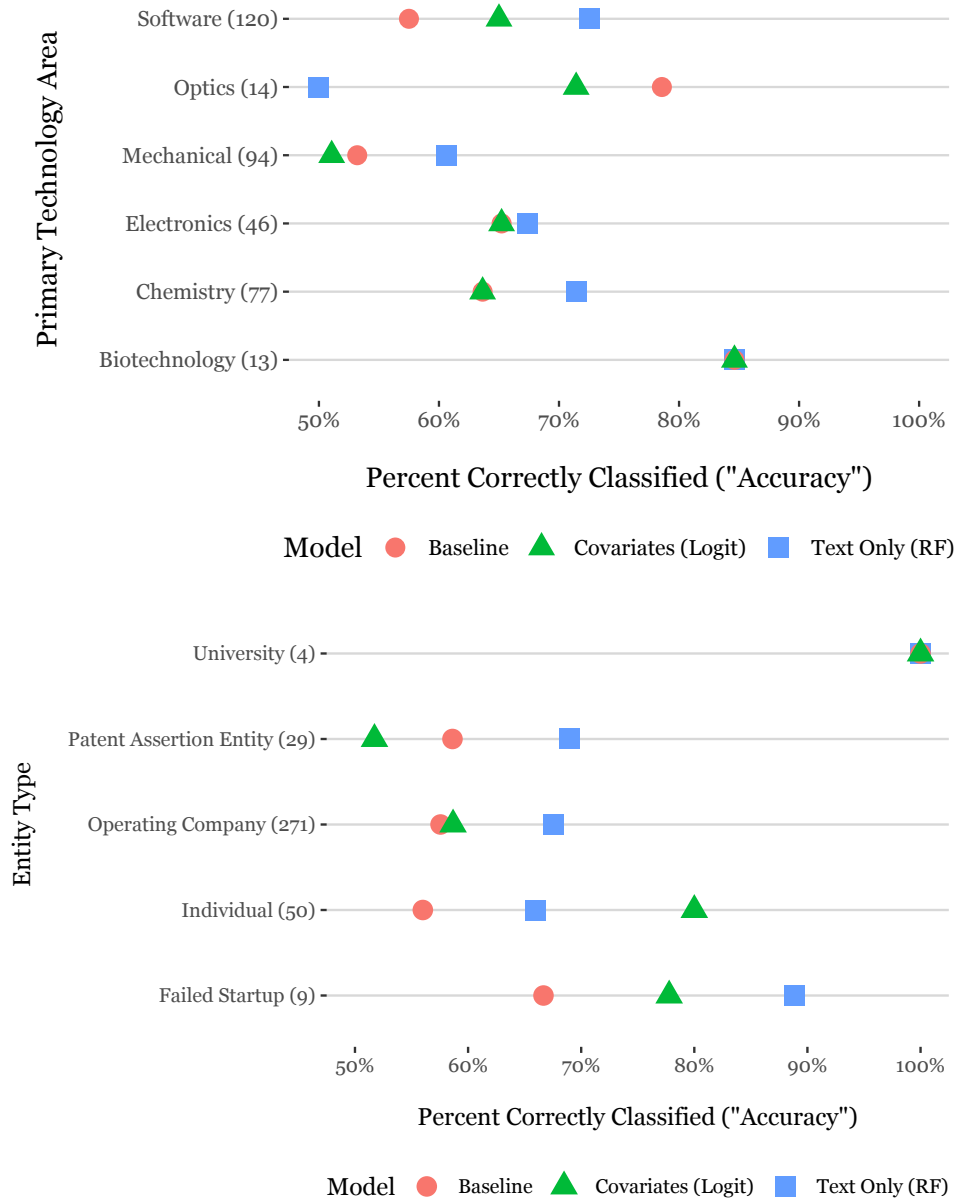
		Actual decision	
		<i>Valid</i>	<i>Invalid</i>
Model prediction	<i>Valid</i>	168	68
	<i>Invalid</i>	37	91
		205	159

Another useful way to explore the predictive accuracy of the model is to look at its success (and failure) across different groups of the data. Recall that—unlike the theory-driven model—the algorithm based on claim text is given no explicit information about the patents, their owners, or any aspects of the litigation. To the extent that there are informative differences contained in the text of electrical and software patents, for example, the model deduced those patterns itself.

Figure 1 shows a breakdown of the model’s success rate across (a) different primary technology areas and (b) the type of entity asserting the patent. The graphs compare the accuracy of the text and theory-driven models, as well as a baseline that simply guesses the most common outcome for each primary technology area or asserting entity, respectively.

In general, the text model predicts well across the whole range of industries. It is notably more accurate than the traditional approach for software,

⁹² In settings where one of the outcomes has more observations, random forests are known to favor that class in its predictions. See HASTIE ET AL., *supra* note __, at 317. To address this, I slightly upweight the “invalid” decisions at the training stage, such that the model sees more of them in each iteration. The overall accuracy varies slightly as this parameter is changed and the model favors decision of validity to a greater or lesser extent, but this does not substantively affect the conclusions (the AUC remains steady and high).



NOTE: Number of observations in parentheses. The “baseline” predicts the most common outcome within each technology area or entity. The most common outcome for biotechnology and mechanical patents was a finding of invalidity; for chemistry, electrical, optics, and software, it was the opposite. Amongst asserting entities, the most common outcome for individuals was a finding of invalidity; for all others it was validity.

chemistry, and mechanical patents. Meanwhile in two sectors—biotechnology and optics patents—the text approach performed worse or equal to the baseline. Both those technology areas, however, had relatively few challenged patents (13 and 14, respectively). If there is something unique about the claim text within each area, the algorithm likely did not have enough data to learn it for these two.

Between different types of patent owner, the text-only model had more mixed success. It is a clear improvement on past models (and on guessing) for operating companies, who make up by far the majority of patentees in this set of litigation. The accuracy amongst patents asserted by PAEs (70% correct) is a particularly striking improvement on the theoretical model. On the other hand, the text model predicted less accurately than the theory-driven model for patents that were asserted by individuals (50 observations). This is an interesting exception to the generally strong performance of the text-based approach, and what drives it is unclear. It may be that individuals assert more idiosyncratic patents, such that any patterns are hard to discern.⁹³

Finally, I ask: how much information is contained in text *alone*? To test this, I reestimate the random forest on the same data, but this time include both the words and metadata together. Table 5 reports the accuracy and AUC for the “combined” model (the first column repeats the text-only results from Table 3 for comparison purposes). The inclusion of the litigation and patent metadata brings a trivial performance improvement, taking the overall accuracy to 71.4%. (Table 6 translates this into concrete outcomes: the combined model correctly predicts one additional invalidity finding.) It seems that a substantial amount of information about the patent is already encoded in the text of the claims.

Table 5: Comparison of Random Forest model performance

	<i>Text-only random forest</i>	<i>Combined random forest</i>
Correctly classified	71.1%	71.4%
AUC	0.78	0.81

⁹³ The evidence here is in line with more recent findings that patents asserted by individuals fare poorly in patentable-subject-matter litigation. See Mark A. Lemley & Samantha Zyontz, *Does Alice Target Patent Trolls?*, 18 J. EMPIRICAL LEGAL STUD. 47 (2021).

Table 6: True outcomes versus predictions (combined model)

		Truth	
		<i>Valid</i>	<i>Invalid</i>
Model prediction	<i>Valid</i>	168	67
	<i>Invalid</i>	37	92
		205	159

E. The challenge of interpretability

While the random forest is superior to the traditional approach in terms of predictive accuracy, it is significantly less interpretable. This is the infamous “black box” character of many machine learning models, and is often an unfortunate tradeoff in this area.⁹⁴ Unlike traditional regression analysis, the random forest model does not produce specific estimates for the effect of any single variable. (Indeed, each variable might relate to the outcome in complicated and non-linear ways; this flexibility is one of the strengths of machine learning.) However, there are ways to explore what factors might be particularly important to the overall accuracy of the model.

The random forest algorithm produces “variable importance” scores, which indicate how far the accuracy of the model falls when a given variable is removed.⁹⁵ (For example, a variable importance of 0.01 would be equivalent to a 1 percentage-point decline in accuracy.) Two notes of caution are in order. First, in situations where variables are highly correlated, the variable importance scores can be artificially low. For example, if two words provide very similar

⁹⁴ As Arti Rai has described, the “black box” nature of an algorithm can be a result of various different factors. Arti Rai, *Machine Learning at the Patent Office: Lessons for Patents and Administrative Law*, 104 IOWA L. REV. 2617 (2019). In the present study, the “black box” derives from complexity. The inputs, outputs, and algorithmic steps are all quite clear, but it is difficult for humans to conceptualize the workings of 2,000 trees simultaneously. *But see* Cynthia Rudin, *Stop Explaining Black Box Machine Learning Models for High Stakes Decisions and Use Interpretable Models Instead*, 1 NATURE MACHINE INTELLIGENCE 206 (2019) (arguing that in areas where explanation itself is a particularly salient social value—for example, in predictions of recidivism—we should always prefer interpretable models).

⁹⁵ In practice, the algorithm randomly rearranges (permutes) the values of the variable in question, so that any real relationship between that variable and the outcome is destroyed. It then re-estimates the model, and notes how much predictive accuracy falls as a consequence of this change.

(but important) information about the outcome, such as “pharmaceutical” and “drug,” then overall accuracy may not decrease when only one of them is removed. Second, the precise variable importance scores can vary slightly as apparently innocuous model specifications are changed.⁹⁶ For both these reasons, we should interpret these scores lightly: a rough guide to what the model is thinking, rather than a precise estimate of any particular word’s effect.

With those cautions in mind, figure 2 shows the fifteen most predictive words for the text-only model, estimated on the entire dataset. A few points stand out. “Amount” and “pharmaceutically” are both amongst the most predictive words. On closer inspection, they correlate very strongly with patents involved in Abbreviated New Drug Applications (ANDA) and the chemical technology area, indicating that they are serving as a proxy for drug patents. Ditto “circuit” and electrical patents, as well as “service” and software. It is striking that the model picks up these relationships so clearly, and their predictive importance reinforces the salience of technology areas that has been highlighted in prior work.⁹⁷

However, it’s not the case that the model is only identifying “easy cases” based on industry. Recall that the algorithm is able to discriminate quite successfully between valid and invalid patents *within* broad technological areas. For example, in this dataset, the word “service” is closely related to a class of software patents that involve cellphone and telecommunications protocols. Unlike software patents generally, the vast majority of this subset of patents were upheld, and the algorithm draws a strong inference from this. Similarly, while the algorithm is very successful at identifying pharmaceuticals as a class, it *also* distinguishes between high and low quality drug patents.

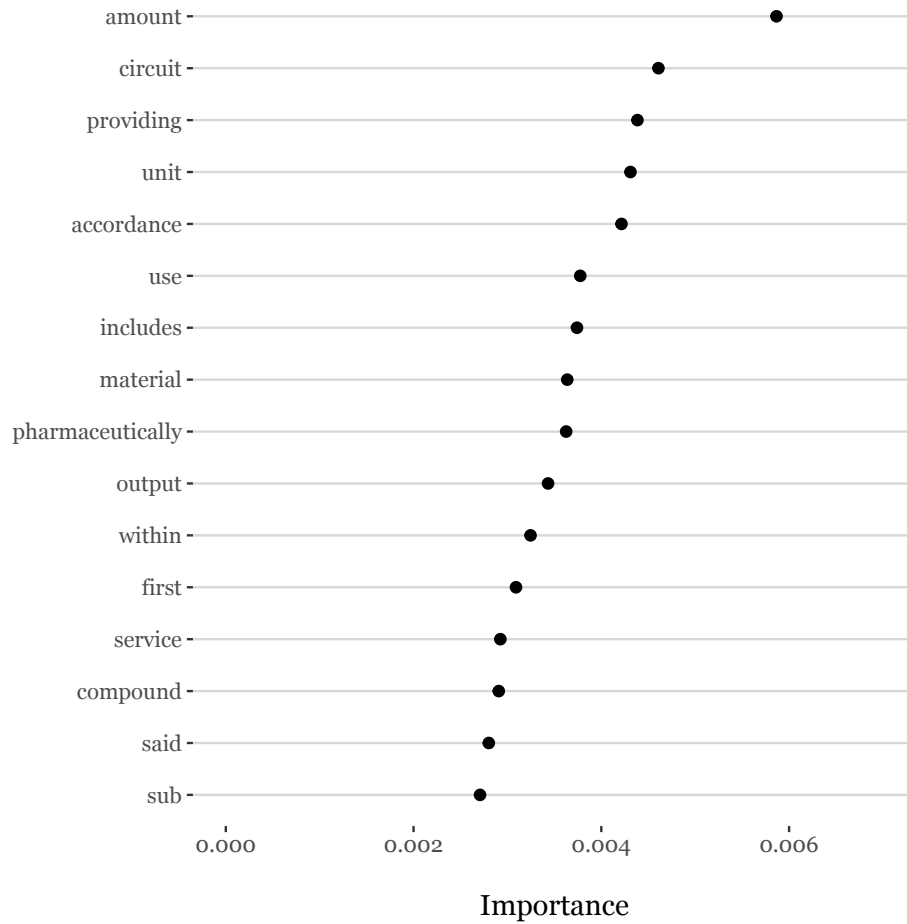
Interestingly, there are relatively few terms of art in the list of most predictive words. Of course, all of these effects are quite small—no single word dominates. This seems appropriate and indicates, unsurprisingly, that the words are related to the outcome (and to each other) in complex ways.

One further note of caution is warranted: nothing about these findings implies that claim text has a causal effect on outcomes. The words proxy for some underlying measure of quality, but they don’t *cause* invalidity decisions in

96 One example of such a parameter is the rate at which we sample from “invalid” and “valid” decisions when building the model. The potential for a multiplicity of models and algorithms to perform at similar levels of accuracy is a broader problem in machine learning, which Breiman refers to as the “Rashomon effect.” Breiman, *supra* note ____.

97 Allison et al., *Divided Patent System*, *supra* note ____.

Figure 2: Top fifteen most predictive words for text-only model



and of themselves. Consider a word that is used almost entirely in the context of pharmaceutical patents. We know from past research that those patents are held invalid at lower rates than other types of utility patent. A good algorithm would therefore be likely to find, *inter alia*, that words that were strongly associated with drug patents were contributors to predictive accuracy (as indeed this model does). But it would make no sense to try to deploy this language outside of the pharmaceutical context. The model finds patterns in the world, not causal relationships.

IV. IMPLICATIONS

My central findings are twofold. First, that patent litigation is considerably more predictable than previous research has realized. And second, that the words on the face of the patent contain significant information about the patent's propensity to be found invalid in court. Even with a relatively small dataset—just 364 adjudications of invalidity—the machine learning algorithm is able to correctly predict a substantial fraction of validity outcomes using no more than the text of the claims. These findings have important but distinct implications for patent research and the practice of patent litigation.

A. New approaches to research on patent quality

On the academic side, there are multiple avenues for future research. While the text of claims appears to encode much of the information that is also contained in the traditional correlates of litigation outcomes, it appears that the claim text contains additional information about the likelihood of a validity decision, above and beyond the usual metadata. Although the goals of prediction and explanation are in some tension here, there are potential ways to incorporate textual measures of quality into more traditional, interpretable models of litigation, particularly given sufficient data.⁹⁸

What can we learn about patent quality more generally? The possibility of directly extrapolating this result to a broader set of non-litigated patents is hampered by serious selection concerns. It is widely understood the patents that are subject to a decision in litigation are not representative of all patents, and it would be a mistake to assume a textual model of the quality of a disputed patent could be applied directly to the broader population.⁹⁹ However, the methods presented here *can* be applied to help better understand the patents at different stages in their lifecycle. An obvious possibility is to use analysis of the text of litigated and non-litigated patents to improve models of selection into litigation. This is an area of significant interest to academics and industry alike. As Colleen Chien has argued, our limited understanding of propensity for litigation renders patent litigation mostly uninsurable, and drives companies to

⁹⁸ For example, with a larger set of litigated patents to draw on, it might be possible to compute a low-dimensional characterization of the description and claim text that could be incorporated into linear regression, as a text-based control for patent quality. *See, e.g.*, Badawi, *supra* note ____.

⁹⁹ *See, e.g.*, Michael D. Frakes & Melissa F. Wasserman, *Do Patent Law Suits Target Invalid Patents?*, in *SELECTION AND DECISION IN JUDICIAL PROCESS AROUND THE WORLD: EMPIRICAL INQUIRIES* 6 (Yun-chien Chang ed., 2019); Allison, *supra* note ____ at 56–7.

amass large patent portfolios as defensive arsenals to ward off the threat of it.¹⁰⁰ But despite various academic efforts,¹⁰¹ there remains a great deal of ambiguity about exactly which characteristics of a patent reliably predict entry into litigation.¹⁰² The results in this Article suggest that the text of patents can provide important insight on this question.

There are, however, some caveats. Beyond the general concern about the limits of explainability discussed in Part III, there are some specific limitations to the findings here. A major caveat of the study is its relatively short time window. The data I use here are extremely high quality, and permit us to compare the performance of the computational text analysis approach to that of existing studies, but they only account for cases filed in 2008 and 2009.¹⁰³ In the past ten years, patent law has been in significant flux. There have been important changes to the doctrines of patentable subject matter¹⁰⁴ and written description,¹⁰⁵ and the statute itself has been substantially revised.¹⁰⁶ Given that

¹⁰⁰ See Chien, *supra* note __ (finding that after-acquired characteristics are somewhat predictive of selection into litigation).

¹⁰¹ *Id.* See also John Allison, Mark A. Lemley, Kimberly A. Moore & Derek Trunkey, *Valuable Patents*, 92 GEO. L.J. 435 (2004) (litigated patents are younger, contain more claims, and are more highly cited than a matched group of non-litigated patents); Jean O. Lanjouw & Mark Schankerman, *Characteristics of Patent Litigation: A Window on Competition*, 32 RAND J. ECON. 129 (2001) (litigated patents are, *inter alia*, more highly cited and have more claims).

¹⁰² See Lee Petherbridge, *On Predicting Patent Litigation*, 90 TEX. L. REV. SEE ALSO 75, 79 (2011) (noting the need for a more “specific and sensitive test” given that litigation is a rare event relative to the population of issued patents).

¹⁰³ The average duration of patent litigation filed in this period was just over three years, such that most cases were resolved by the end of 2012, although certain cases which were appealed were pending for considerably longer.

¹⁰⁴ See, e.g., *Mayo Collaborative Servs. v. Prometheus Labs., Inc.*, 566 U.S. 66 (2012) (mere applications of laws of nature are not eligible subject matter in the absence of some additional inventive step); *Ass’n for Molecular Pathology v. Myriad Genetics, Inc.*, 569 U.S. 576 (2013) (isolated DNA sequences that occur naturally are not eligible subject matter); *Alice Corp. v. CLS Bank Int’l*, 573 U.S. 208 (2014) (abstract ideas are not rendered patentable merely through a generic computational implementation). The patentable subject matter cases in particular represent an important change in patent litigation. See Lemley & Zyontz, *supra* note __.

¹⁰⁵ *Nautilus, Inc. v. Biosig Instruments, Inc.*, 572 U.S. 898 (2014) (patent claims must inform those skilled in the art of the scope of the invention, with reasonable certainty).

¹⁰⁶ Leahy-Smith America Invents Act, Pub. L. No. 112-29, 125 Stat. 284 (2011). See also Mark A. Lemley, *The Surprising Resilience of the Patent System*, 95 TEX. L. REV. 1, 1–6 (2016) (describing some of the major changes to and in the patent system in the last three decades).

aspects of the law have changed in important ways, it would be a mistake to extrapolate this *particular* model to a case filed a decade later and expect it to predict as well.¹⁰⁷ However, the core contribution of this Article is conceptual. At this level, there appears to be little reason to expect that the text of patent claims would become systematically less informative as a result of the intervening legal changes, even as the precise structure of the patterns has undoubtedly changed.¹⁰⁸

At the same time, however, the nature of patent litigation itself may have changed. For example, there has been a significant increase in the volume of litigation brought by NPEs.¹⁰⁹ In the 2009-10 study window, relatively few filed cases fell into this category. But the years that followed played host to a stark change in filing patterns: Cotropia et al. report that in 2012, nearly 50% of patent litigation was brought by NPEs of some kind, an increase of nearly 30 percentage-points since 2010.¹¹⁰ For the purposes of prediction, this change—which was driven primarily by patent holding companies—may be consequential to the extent that we think the patents asserted by these new PAEs are different *in kind* to other patents. For example, it could change the makeup of cases if PAEs were more likely to advance low-probability cases to merits decisions. Still, there is no reason to assume *a priori* that the text-based model would perform poorly in this context. (And indeed, it was remarkably successful at discriminating between the valid and invalid patents asserted by PAEs in 2008 and 2009.)

¹⁰⁷ In the machine learning literature, this phenomenon is known as “concept drift.” See, e.g., João Gama et al., *A Survey on Concept Drift Adaption*, 46 ACM COMPUTER SURV. (2014). For an example of significant change that occurred just before the study window see Ryan T. Holte & Ted Sichelman, *Cycles of Obviousness*, 105 IOWA L. REV. 107 (2019) (finding a notable increase in invalidations for obviousness after the Supreme Court’s landmark decision in *KSR v. Teleflex*).

¹⁰⁸ In one important regard, this claim might be controversial. The recent line of Supreme Court cases addressing patentable subject matter (described in note __, *supra*) has been widely criticized for creating an insolubly ambiguous standard for patent eligibility. But recent survey evidence indicates that patent prosecutors (though not litigators!) are still able to predict outcomes with some degree of confidence. See Jason D. Reinecke, *Is the Supreme Court’s Patentable Subject Matter Test Overly Ambiguous? An Empirical Test*, 2019 UTAH L. REV. 581. This is an obvious area for future research.

¹⁰⁹ Cotropia et al., *supra* note __.

¹¹⁰ *Id.* at 674.

B. What does this mean for patent litigation?

Although there are valuable insights in this work for scholars, perhaps the most fertile ground for the kind of predictive analytics technology demonstrated here is in the profession. A central component of litigation practice lies in answering what are, essentially, questions of prediction. How likely are we to succeed on this eligibility question? How strong is this patent? What are the chances of a successful invalidity counterclaim?¹¹¹ The results in this paper suggest that validity decisions can be forecast with greater confidence than previously thought.

To be sure, uncertainty remains: even for this model, many predictions fall into a grey area.¹¹² But the higher-probability predictions allow patentees to make more informed, confident choices about their litigation strategy, and to reduce overall litigation costs.¹¹³ Equally important, from the perspective of alleged infringers, better *ex ante* prediction tools can also serve as a prophylactic against strategies which assert many patents in a single suit—some of which may be of questionable validity—in order to overwhelm defendants.¹¹⁴ Finally, more accurate prediction is highly consequential on the plaintiff side, given the growth of contingent fee representation in patent litigation.¹¹⁵

These changes sit in the context of broader developments in law and legal analytics.¹¹⁶ A burgeoning literature on the use of predictive analytics in legal

111 It is important to note that the predictive approach demonstrated here is related to, but distinct from, computational tools that guide practitioners as to *how* to calibrate legal arguments. See, e.g., *infra* notes 122–25 and accompanying text.

112 The graphs in Appendix Part C give a visual sense of how well the various models succeed in separating true valid and true invalid decisions. In cases where the text model gives a predicted probability of validity above 70% (or below 30%), there are very few false positives.

113 See *supra* note ____.

114 See Schwartz, *supra* note ____ at 375–6.

115 *Id.*

116 Engstrom and Gelbach catalog some of the many companies now operating in the “legal tech” space. Most provide sophisticated data aggregation and presentation (e.g., Lex Machina; Gavelytics), but a number focus explicitly on questions of prediction, including Blue J Legal (tax) and Colossus (insurance). See David Freeman Engstrom & Jonah B. Gelbach, *Legal Tech, Civil Procedure, and the Future of Adversarialism*, 169 U. PA. L. REV. ____ (2021).

practice evinces both enthusiasm from proponents of “legal tech,”¹¹⁷ and reservations from those concerned about its impact on the profession and the fair administration of justice.¹¹⁸ Although concerns about fairness and access to justice are extremely important (particularly in criminal and consumer law), the meta-conversation may have run ahead of our current technological capacity in some respects.¹¹⁹ David Engstrom and Jonah Gelbach argue that predictive analytics will likely remain too superficial to account for the range of nuance that is implicit in much legal reasoning, which incorporates both “synoptic” and “subtle” judgments.¹²⁰ Put differently, legal decisions simultaneously involve two types of analysis. First, the identification of regular, observable regularities—e.g., what are the broad patterns of decisions about particular types of patent, or within specific jurisdictions? Second, the parsing of “shades of grey”—e.g., how do the factual and linguistic idiosyncrasies distinguish *this* patent from previous cases? Current NLP technology excels primarily at the former, and is less well suited to the latter.

The evidence in this paper broadly supports this supposition. In the context of a small sample and with no specific instruction (or professional expertise), NLP proved remarkably good at identifying consistent patterns in the data. But while the models in this paper are a significant improvement on the state of the art, there remains a clear grey area: even with the text of the claims, litigant information, and patent metadata all combined, the model still misclassified one-quarter of the decisions. More training data will no doubt improve the accuracy of machine predictions of the sort that I demonstrate here, but the

¹¹⁷ See, e.g., Daniel Martin Katz, *Quantitative Legal Prediction—Or—How I Learned to Stop Worrying and Start Preparing for the Data-Driven Future of the Legal Services Industry*, 62 EMORY L.J. 910 (2013); Benjamin Alarie, *The Path of Law: Towards Legal Singularity*, 66 U. TORONTO L.J. 443 (2016); John O. McGinnis & Russell G. Pearce, *The Great Disruption: How Machine Intelligence Will Transform the Role of Lawyers in the Delivery of Legal Services*, 82 FORDHAM L. REV. 3041 (2014).

¹¹⁸ See, e.g., Frank Pasquale & Glyn Cashwell, *Prediction, Persuasion, and the Jurisprudence of Behaviourism*, 68 U. TORONTO L.J. 63 (2018); Emily S. Taylor Poppe, *The Future is Bright Complicated: AI, Apps & Access to Justice*, 72 OKLA. L. REV. 185 (2019); Citron & Pasquale, *supra* note ____.

¹¹⁹ See, e.g., Dana Remus & Frank Levy, *Can Robots Be Lawyers? Computers, Lawyers, and the Practice of Law*, 30 GEO. J. LEGAL ETHICS 501 (2017) (surveying the capacity of existing technology and its likely effect on different areas of legal practice).

¹²⁰ Engstrom & Gelbach, *supra* note ____ (distinguishing between the “synoptic” and “subtle” aspects of legal judgment, and arguing that machine learning is well suited to the former, but not to the latter).

most effective forecasting in this area will likely remain a human-guided process for the foreseeable future.¹²¹

Finally, a more troubling implication of natural language approaches to patent litigation lies in the future of patent prosecution. To the extent that future applicants are able to identify specific linguistic techniques which tend to find favor with courts, but which have little to do with the disclosure of information about the underlying innovation, we might be concerned about further divorcing patent doctrine from the law's underlying social purpose. At this stage, such a concern is probably premature. Patent doctrine is a dynamic area of law, responding to constant changes in technology and social norms, such that particular patterns in patent language seem unlikely to be stable for too long.¹²²

However, to the extent that computational approaches are able to identify durable patterns in courts' decision-making, it is probably more appropriate to view this as evolution rather than revolution. Indeed, similar tools are already more commonly deployed in patent prosecution.¹²³ Computational tools can be used to automate prior-art searching,¹²⁴ calibrate applications to the particular examiner and PTO Art Unit,¹²⁵ and even draft the application itself.¹²⁶ Moreover,

121 See also Eric Talley, *Is the Future of Law a Driverless Car? Assessing How the Data-Analytics Revolution Will Transform Legal Practice*, 174 J. INSTITUTIONAL & THEORETICAL ECON. 183, 185 (2018) (arguing that law is “irreducibly complex” and will continue to need “significant human input”).

122 Fagan and Levmore argue that “machine learning is less impressive when the past is unlike the future.” Frank Fagan & Saul Levmore, *The Impact of Artificial Intelligence on Rules, Standards, and Judicial Discretion*, 93 S. CAL. L. REV. 1, 1 (2019). This is true, of course, of any forecasting exercise. And since patent law is a dynamic area, the most successful applied uses of the approach shown here would undoubtedly need to be updated regularly with contemporaneous court decisions. Fortunately, this is an area in which comprehensive data collection is quite feasible. Patent litigation is relatively low volume (particularly in terms of cases which receive a decision on the merits), and case data are comprehensively collected by popular third-party analytics services such as Lex Machina and Docket Navigator.

123 Tabrez Y. Ebrahim, *Automation & Predictive Analytics in Patent Prosecution: USPTO Implications & Policy*, 35 GA. ST. L. REV. 1185 (2019).

124 See, e.g., GOOGLE PATENTS, <https://patents.google.com> (last visited Feb. 9, 2021); PQAI, <https://projectpq.ai> (last visited Feb. 9, 2021).

125 See, e.g., *AI for Patent Drafting*, ROWAN PATENTS, <https://rowanpatents.com/drafting> (last visited Feb. 9, 2021).

126 This final example is perhaps the most ambitious use of natural language tools, but has already been reduced to practice in the area of specification drafting. See, e.g., SPECIFIO, <https://specif.io> (last visited Feb 9, 2021); Ebrahim, *supra* note __ at 1197–99.

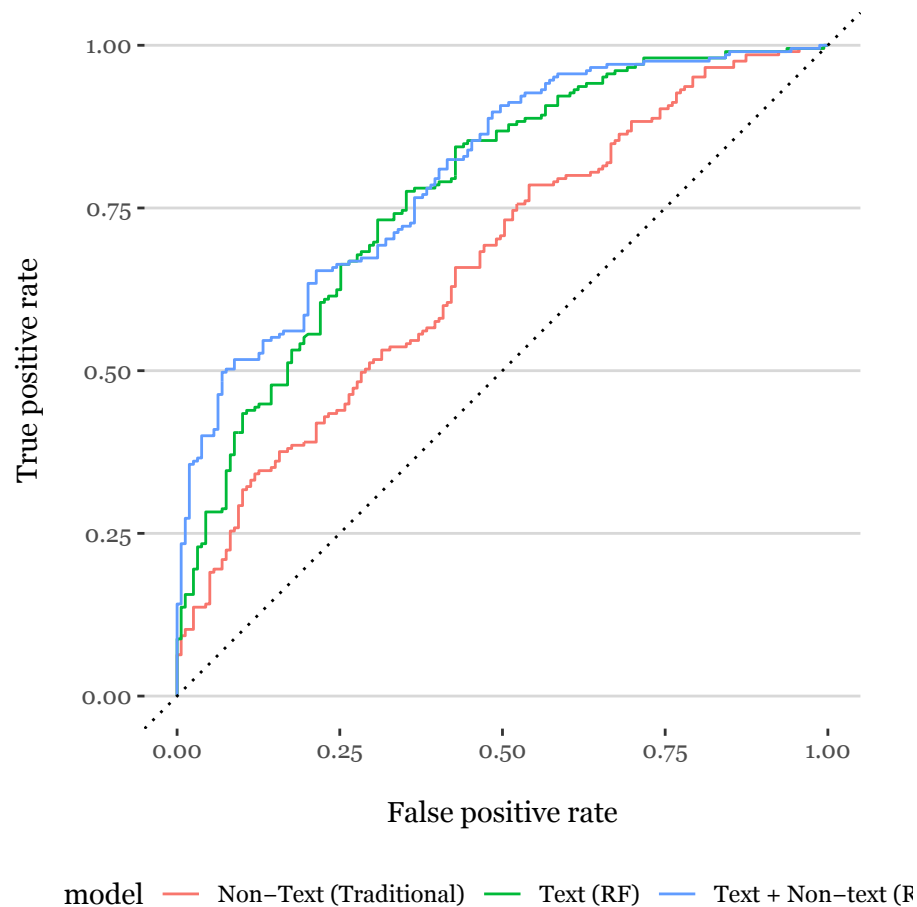
like many legal professionals, patent prosecutors are to some degree pattern-recognition experts. Claim drafting has already evolved into a highly specialized art and language, as patent attorneys respond to (and try to anticipate) decisions by patent examiners, the PTAB, and federal courts on a range of esoteric but consequential questions of linguistic interpretation. Machine learning tools may accelerate this process, but they are ultimately different in capacity, rather than different in kind.

APPENDIX

A. Receiver operating characteristic (ROC) curves

Logistic regression and random forests, like many other classification models, generate a probability for every outcome, and then use this probability to assign a predicted class (in this case, “valid” or “invalid”). By default, most classifiers use 50% as the assignment cutoff, but depending on the balance of the underlying data—or cost of false positives/negative—it is often desirable to set a different threshold. An ROC curve shows the effect of varying this cutoff.

The ROC varies the threshold from 0 to 100, and at each level calculates the rate of “true positives” (in our case, correct prediction of invalidity) and “false positives” (incorrect predictions of invalidity). The dotted diagonal line is the baseline: every point on this line is equivalent to guessing the outcome with 50% probability. Curves that are closer to the top left corner indicate an algorithm that is more successful at discriminating between outcomes. (The AUC measure, described in Part III.C., refers to the area under this curve.)



B. Separation graphs

The graph below shows the density of predictions for each model, split into valid and invalid patents. The blue curve shows the “true” invalid patents, while the red curve plots the “true” valid patents. The bottom axis shows the model’s predicted probability of validity. The traditional model is notably worse at separating the valid and invalid decisions—a substantial fraction of the true invalid patents have a validity prediction greater than 50%.

