

# NEW YORK UNIVERSITY LAW REVIEW

---

VOLUME 95

MAY 2020

NUMBER 2

---

## ARTICLES

### AN EMPIRICAL STUDY OF STATUTORY INTERPRETATION IN TAX LAW

JONATHAN H. CHOI\*

*A substantial academic literature considers how agencies should interpret statutes. But few studies have considered how agencies actually do interpret statutes, and none has empirically compared the methodologies of agencies and courts in practice. This Article conducts such a comparison, using a newly created dataset of all Internal Revenue Service (IRS) publications ever released, along with an existing dataset of court decisions. It applies natural language processing, machine learning, and regression analysis to map methodological trends and to test whether particular authorities have developed unique cultures of statutory interpretation.*

*It finds that, over time, the IRS has increasingly made rules on normative policy grounds (like fairness and efficiency) rather than merely producing rules based on the “best reading” of the relevant statute (under any interpretive theory, like purposivism or textualism). Moreover, when the IRS does focus on the statute, it has grown much more purposivist over time. In contrast, the Tax Court has not grown more normative and has followed the same trend toward textualism as most other courts. But although the Tax Court has become more broadly textualist, it prioritizes different interpretive tools than other courts, like Chevron deference and holistic-textual canons of interpretation. This suggests that each authority adopts its own flavor of textualism or purposivism.*

*These findings complicate the literature on tax exceptionalism and the judicial nature of the Tax Court. They also inform ongoing debates about judicial deference and the future of doctrines like Chevron and Skidmore deference. Most broadly,*

---

\* Copyright © 2020 by Jonathan H. Choi, Fellow, New York University School of Law. Thanks to Anne Alstott, Aaron Bruhl, Richard Epstein, Bill Eskridge, Ryan Fackler, Jacob Goldin, Andy Grewal, Lilai Guo, Kristin Hickman, Hayes Holderness, Ariel Kleiman, David Louk, Joshua Macey, Florencia Marotta-Wurgler, and Julian Nyarko, as well as the participants in the Harvard Law School Caselaw Access Project Summit, the Conference on Empirical Legal Studies, the Association of American Law Schools annual conference, the New York University School of Law Lawyering Scholarship Colloquium, and the Critical Tax Conference.

*they provide an empirical counterpoint to the existing theoretical literature on statutory interpretation by agencies.*

INTRODUCTION .....	365
I. KEY QUESTIONS .....	371
A. <i>Statutory Judgments or Normative Policymaking?</i> ...	371
B. <i>Textualism or Purposivism?</i> .....	375
C. <i>Cohesion Among Courts or Among Specialists?</i> .....	379
II. EMPIRICAL METHODS .....	382
A. <i>Natural Language Processing</i> .....	382
B. <i>Machine Learning</i> .....	386
C. <i>Regression Analysis</i> .....	388
D. <i>Limitations</i> .....	389
1. <i>Term Frequency as a Proxy for Methodology</i> ...	389
2. <i>Doing Different Things, Doing Things Differently</i> .....	390
III. RESULTS .....	391
A. <i>The IRS Has Become More Normative and Less Statutory</i> .....	392
B. <i>The Tax Court Has Maintained the Same Proportion of Statutory and Normative Terms</i> .....	395
C. <i>The IRS Has Become More Purposivist and Less Textualist</i> .....	397
D. <i>The Tax Court Has Become More Textualist and Less Purposivist</i> .....	398
E. <i>The Tax Court Has Developed a Unique Interpretive Methodology Relative to Other Courts</i> .....	401
F. <i>Democratic Judges Are More Purposivist and Republican Judges Are More Textualist at the Tax Court</i> .....	404
G. <i>Case Outcomes Do Not Statistically Significantly Predict Interpretive Methodology at the Tax Court</i> ..	407
IV. ROBUSTNESS CHECKS .....	410
A. <i>Reading Cases to Confirm Term Frequency Results</i> .	410
B. <i>Bootstrapped Confidence Intervals for Machine Learning</i> .....	410
C. <i>Validating OCR Quality over Time</i> .....	412
D. <i>Confirming that Results Are Not Driven by Changes in Terminology</i> .....	413
CONCLUSION .....	415
APPENDIX .....	416
A. <i>Data Sources</i> .....	416
1. <i>IRS Publications</i> .....	416

2.	<i>Court Opinions</i> .....	418
3.	<i>Excluding Non-Substantive Opinions</i> .....	418
B.	<i>Terms Analyzed</i> .....	419
1.	<i>Purposivist Terms</i> .....	420
2.	<i>Textualist Terms</i> .....	421
3.	<i>Statutory Terms</i> .....	422
4.	<i>Normative Terms</i> .....	423
5.	<i>Substantive Canons</i> .....	423
C.	<i>Non-Normal Distribution of Term Frequencies in Tax Court Opinions</i> .....	424
D.	<i>Tf-idf Transformation and Classification in Machine Learning</i> .....	427
E.	<i>Regression Analysis of Tax Court Opinions</i> .....	428
1.	<i>Ordinary Least Squares Regression Model</i> .....	430
2.	<i>Two-Part Regression Model</i> .....	432
F.	<i>Log-Transformed Charts</i> .....	437
G.	<i>Bootstrapped Confidence Intervals for Machine Learning</i> .....	440

INTRODUCTION

After decades of debate over statutory interpretation by courts, scholars have more recently turned to the interpretive practices of agencies. Many have argued that agencies have relatively greater expertise in assessing statutory purpose, concluding that they ought to be more purposivist<sup>1</sup> than courts.<sup>2</sup> More fundamentally, many have

<sup>1</sup> Textualists generally emphasize the plain meaning of statutory text and eschew legislative history. Purposivists generally look to all available evidence, including legislative history. The methodological distance between purposivists and textualists is often overstated, since all sides generally attempt to reconstruct statutory purpose and merely differ in the tools that they use to do so. For instance, although textualists are often presented as the foil to purposivists, modern textualists will also generally consider nontextual indicia of statutory purpose when statutory text is unclear. *See, e.g.,* John F. Manning, *What Divides Textualists from Purposivists?*, 106 COLUM. L. REV. 70, 84–85 (2006) (“[T]extualists generally forgo reliance on legislative history as an authoritative source of [the statute’s apparent overall] purpose, but that reaction goes to the reliability and legitimacy of a certain type of evidence of purpose . . . . [W]hen semantic ambiguity creates the necessary leeway, textualists will try to construct a plausible hypothetical purpose (if possible) . . . .”).

<sup>2</sup> *See, e.g.,* William N. Eskridge Jr., *Expanding Chevron’s Domain: A Comparative Institutional Analysis of the Relative Competence of Courts and Agencies to Interpret Statutes*, 2013 WIS. L. REV. 411, 434 (“[A]gencies interpret statutes purposively, and that is on the whole a good impulse in the modern regulatory state. A consequence of a purposivist approach to statutes is that the interpreter will read the statute dynamically, to reach beyond the original problems that were the basis of congressional deliberation.”); Michael Herz, *Purposivism and Institutional Competence in Statutory Interpretation*, 2009 MICH. ST. L. REV. 89, 92 (“In general, my conclusion is that agencies make more

suggested that judicial deference regimes, like *Chevron* deference,<sup>3</sup> empower agencies to make rules based on normative policy concerns, rather than merely seeking the “best reading” of a statute (using purposivism, textualism, or any other methodology).<sup>4</sup>

But despite a large theoretical literature on how agencies ought to interpret statutes, little scholarship has considered how they actually *do* interpret statutes.<sup>5</sup> Past work has focused on agency practice

---

respectable and less problematic purposivists than do judges.”); Jerry L. Mashaw, *Norms, Practices, and the Paradox of Deference: A Preliminary Inquiry into Agency Statutory Interpretation*, 57 ADMIN. L. REV. 501, 511 (2005) (“In some instances, only the skillful deployment of legislative history will permit agencies to fulfill their constitutional role as faithful agents in the statute’s implementation.”); Cass R. Sunstein & Adrian Vermeule, *Interpretation and Institutions*, 101 MICH. L. REV. 885, 928 (2003) (“[A]gencies are likely to be in a better position to decide whether departures from the text actually make sense.”). But see Richard J. Pierce, Jr., *How Agencies Should Give Meaning to the Statutes They Administer: A Response to Mashaw and Strauss*, 59 ADMIN. L. REV. 197, 202 (2007) (“[T]he agency should use the same ‘traditional tools of statutory construction’ that it expects a reviewing court to use. If the agency uses a different method of interpretation . . . it increases significantly the risk of judicial reversal without good reason.”). Pierce’s suggestion that agencies should follow the interpretive practices of courts only applies to interpretation carried out in *Chevron* step one. With respect to *Chevron* step two, Pierce believes (as do many others) that agencies ought to select the best policy rather than relying on any conventional interpretive norms. See *infra* note 4 and accompanying text.

<sup>3</sup> *Chevron U.S.A. Inc. v. Nat. Res. Def. Council, Inc.*, 467 U.S. 837 (1984).

<sup>4</sup> See *infra* Section I.A.; see also *Chevron*, 467 U.S. at 842, 843 & n.9, 845 (holding that an agency’s interpretation of an ambiguous statute warrants deference so long as it represents a “reasonable policy choice”); cf. *Covad Commc’ns Co. v. FCC*, 450 F.3d 528, 537 (D.C. Cir. 2006) (requiring the agency to “articulate a satisfactory explanation for its action including a rational connection between the facts found and the choice made” (internal citations and quotations omitted)); Pierce, *supra* note 2, at 200 (arguing that, under *Chevron*, agencies can choose among permissible interpretations of a statute “only by engaging in a policymaking process”). But see Aaron Saiger, *Agencies’ Obligation to Interpret the Statute*, 69 VAND. L. REV. 1231, 1231 (2016) (“An agency that commands deference bears a duty to adopt what it believes to be the best interpretation of the relevant statute.”). One could theoretically defend normative rulemaking by arguing that a reasonable legislator would have preferred the normatively best policy to prevail, and that therefore the best means for the agency to act as the “faithful agent” of the legislator is to prioritize policy concerns. This might be considered a particularly expansive form of purposivism, reminiscent of T. Alexander Aleinikoff’s “nautical” approach, which “understands a statute as an on-going process (a voyage) in which both the shipbuilder and subsequent navigators play a role.” T. Alexander Aleinikoff, *Updating Statutory Interpretation*, 87 MICH. L. REV. 20, 21 (1988). That said, scholars have generally accepted the distinction between the pursuit of the “best reading” of a statute and the “best policy.” See *infra* Section I.A.

<sup>5</sup> See Amy Semet, *An Empirical Examination of Agency Statutory Interpretation*, 103 MINN. L. REV. 2255 (2019) (considering statutory interpretation in decisions by the National Labor Relations Board); Christopher J. Walker, *Inside Agency Statutory Interpretation*, 67 STAN. L. REV. 999 (2015) (surveying attitudes toward statutory interpretation among agency administrators). However, Walker’s survey did not consider any actual decisions by agencies, and Semet’s empirical work may be specific to the NLRB, due to its unusually intense partisanship. See Semet, *supra*, at 2280 (“Board voting is highly ideological . . . . Often, the Board reverses many of the decisions of the prior

within a relatively narrow period,<sup>6</sup> making it impossible to evaluate how agency practice differed over time (especially before and after *Chevron*). Moreover, no empirical work has compared how agencies and courts differ while interpreting the same statutes.

This Article contributes to this conversation by studying a fertile area for agency-court comparisons: federal tax law. Because the IRS is one of the largest government agencies,<sup>7</sup> and because its Internal Revenue Bulletin has been published so consistently (weekly)<sup>8</sup> for so long (since 1919),<sup>9</sup> it provides ample material for a longitudinal study of interpretive methodology over time. Similarly, the Tax Court handles the vast majority of federal tax cases (roughly ninety-seven percent)<sup>10</sup> and has operated since 1942,<sup>11</sup> again producing a large amount of source material.

It was previously difficult or impossible to analyze such large bodies of documentation, not least because they were not readily accessible by researchers. This Article addresses this problem by creating a new dataset of all Internal Revenue Bulletins ever published, which it analyzes along with a dataset recently launched by Harvard Law School's Caselaw Access Project.<sup>12</sup> Between these two sources,

---

administration when a new partisan majority takes gains [sic] control of the Board.”); Ronald Turner, *Ideological Voting on the National Labor Relations Board*, 8 U. PA. J. LAB. & EMP. L. 707, 712 (2006) (arguing the same point). More broadly, without considering comparable judicial practice, it is difficult to say how much her results were driven by the NLRB's status as an agency versus how much they were driven by issues unique to labor relations law.

<sup>6</sup> Walker's article relies on a single survey conducted in 2013. Walker, *supra* note 5, at 1015. Semet's article considers NLRB decisions between 1993 and 2017. Semet, *supra* note 5, at 2282. *Chevron* was decided in 1984. See *Chevron*, 467 U.S. 837.

<sup>7</sup> The IRS had 74,454 employees as of fiscal year 2019, forming the vast majority of the Treasury Department's staff. See INTERNAL REVENUE SERV., DEP'T OF THE TREASURY, CONGRESSIONAL BUDGET JUSTIFICATION AND ANNUAL PERFORMANCE REPORT AND PLAN: FISCAL YEAR 2020, at 1 (2019); see also Renu Zaretsky, *America, We Have a Problem: The IRS Brain Drain*, TAX POL'Y CTR.: TAXVOX (Feb. 6, 2019), <https://www.taxpolicycenter.org/taxvox/america-we-have-problem-irs-brain-drain> (“[O]ver 80 percent of Treasury Department employees work at the IRS.”).

<sup>8</sup> E.g., 2020-1 I.R.B., intro. (“The Internal Revenue Bulletin is the authoritative instrument of the Commissioner of Internal Revenue for announcing official rulings and procedures of the Internal Revenue Service . . . . It is published weekly.”).

<sup>9</sup> See 1 C.B. i (1919).

<sup>10</sup> Elizabeth Chao & Andrew R. Roberson, *Overview of Tax Litigation Forums*, TAX CONTROVERSY 360 (Apr. 21, 2017), <https://www.taxcontroversy360.com/2017/04/overview-of-tax-litigation-forums>.

<sup>11</sup> The U.S. Board of Tax Appeals, the predecessor to the Tax Court, was founded by the Revenue Act of 1924. Revenue Act of 1924, Pub. L. No. 68-176, § 900, 43 Stat. 253, 336. The Board of Tax Appeals was restructured and renamed the U.S. Tax Court by the Revenue Act of 1942. Revenue Act of 1942, Pub. L. No. 77-753, § 504(a), 56 Stat. 798, 957.

<sup>12</sup> See *infra* Appendix Section A.

this Article analyzes 182,535 pages of Internal Revenue Bulletins and 470,099 court opinions.<sup>13</sup>

Broadly, this Article asks four main questions. First, how have interpretive methods evolved at the Tax Court and the IRS *within* each institution? Second, what is the difference *between* institutions—do agencies interpret statutes differently from courts? Third, what is the difference *between subject areas*—does the Tax Court interpret statutes differently from other federal courts (both Article I and Article III courts)? Fourth, what are the implications of interpreters' choices *between methods*—do they vary by party, or are particular methods associated with particular outcomes (either pro- or anti-taxpayer)?

To answer these questions, this Article uses “natural language processing” (algorithmic analysis of large bodies of text)<sup>14</sup> to assess how the IRS, the Tax Court, and other courts have used different tools in their decisions over time: statutory versus normative<sup>15</sup> and textualist versus purposivist. It measures the frequency with which authorities cite these tools—for example, textualists citing dictionaries or purposivists citing legislative history—to map methodological trends. It then uses machine learning for more granular analysis,<sup>16</sup> by training algorithms to distinguish between court opinions based on interpretive methodology alone. This allows the algorithm to identify which specific terms, if any, are most strongly associated with the Tax Court and with the district courts, providing a more nuanced account of the kind of purposivism or textualism each court applies. Finally, the Article uses regression analysis to test whether methodology can be

---

<sup>13</sup> See *infra* Appendix Section A.

<sup>14</sup> See CHRISTOPHER D. MANNING & HINRICH SCHÜTZE, FOUNDATIONS OF STATISTICAL NATURAL LANGUAGE PROCESSING, at xix (1999).

<sup>15</sup> This Article describes decisionmaking as “statutory” when it reflects a decisionmaker’s attempt to act as a “faithful agent” of the legislature, archaeologically discerning a statute’s true meaning while abstaining from value judgments. A statutory approach, under this definition, may follow any interpretive method, including textualism, purposivism, or pragmatism. In contrast, a “normative” approach reflects a decisionmaker’s attempt to create rules *de novo* based on its own policy preferences. There is a broader sense in which any decision by a court or agency could be described as “statutory” if it concerns a statute; this Article does not use the term in that sense. The statutory and normative perspectives will often overlap and will often be considered simultaneously, especially since the normative desirability of a particular interpretation might be considered a factor in favor of its statutory validity. See also *infra* note 4 (observing that a nonstandard view of interpretation might hold that the statutory and normative viewpoints are identical).

<sup>16</sup> Machine learning uses computer algorithms in order to accomplish a particular task without human instructions. This Article primarily uses machine learning based on statistical inference. See *infra* Section II.B.

predicted based on certain case characteristics, such as the party of the trial judge or the case outcome.

The main results are as follows. First, over time, the IRS has increasingly issued guidance based on normative preferences rather than statutory evidence.<sup>17</sup> In contrast, the Tax Court has used roughly the same proportion of normative and statutory language since it was founded in 1942.<sup>18</sup>

Second, the IRS became much more purposivist and less textualist from the 1920s to approximately 1950, but has retained the same relatively purposivist posture since then.<sup>19</sup> On the other hand, the Tax Court has followed the general judicial movement of the past four decades away from purposivism and toward textualism.<sup>20</sup> The combination of these first two results suggests greater methodological cohesion among courts than among tax specialists.

Third, the machine learning results reveal that Tax Court opinions can be distinguished from those of district courts and the Court of Federal Claims based on the specific interpretive tools each employs. As compared to district courts, the Tax Court favors congressional reports (especially reports from the Congressional Budget Office and the Joint Committee on Taxation) over hearings, holistic-textual canons (those emphasizing a cohesive reading of the tax code) over language canons, and *Chevron* deference over constitutional canons.<sup>21</sup> This complicates the conventional story that all courts have become more textualist; while this is true in broad terms, the precise flavor of each court's interpretive methods differs in the details.

Fourth, regression analysis indicates that Tax Court judges appointed by Democratic presidents are more likely to use purposivist terms and less likely to use textualist terms than Republican appointees.<sup>22</sup> However, substantive outcomes (whether the court rules for or against the taxpayer) do not have a statistically significant relationship with interpretive methodology.<sup>23</sup>

Apart from theoretical interest, the findings in this Article have important practical implications. By underscoring agencies' shift toward normative decisionmaking, this Article is consistent with the widespread belief that *Chevron* permits agencies to make their own policy judgments rather than merely restating Congress's. Some

---

<sup>17</sup> See *infra* Section III.A.

<sup>18</sup> See *infra* Section III.B.

<sup>19</sup> See *infra* Section III.C.

<sup>20</sup> See *infra* Section III.D.

<sup>21</sup> See *infra* Section III.E.

<sup>22</sup> See *infra* Section III.F.

<sup>23</sup> See *infra* Section III.G.

scholars view this as a feature of judicial deference, and some view it as a bug. Either way, this finding informs the positions taken by *Chevron*'s critics and its supporters.

The findings also suggest that tax exceptionalism—the widespread belief that tax statutes are or ought to be interpreted differently<sup>24</sup>—may be overstated in some respects and understated in others. Overstated, in that the Tax Court methodologically hews closer to other courts than to the IRS, despite the Tax Court and the IRS's shared subject matter. So the conventional story that tax experts are exceptional because they are more purposivist may be incorrect. Understated, at the same time, in that the Tax Court does differ from other courts in its particular selection of textualist tools, suggesting that a more nuanced form of exceptionalism may apply.

Finally, the findings support controlling but controversial case law indicating that the Tax Court plays an “exclusively judicial role.”<sup>25</sup> The conclusion in this Article that the Tax Court interprets statutes more like other courts than like the IRS undermines the claims of some scholars that Tax Court opinions should be subject to judicial deference, much like agency pronouncements.<sup>26</sup> Instead, at least on the key dimension of interpretive methodology, the Tax Court

---

<sup>24</sup> See *infra* notes 76–79 and accompanying text.

<sup>25</sup> *Freytag v. Comm'r*, 501 U.S. 868, 892 (1991). But see *Kuretski v. Comm'r*, 755 F.3d 929, 932 (D.C. Cir. 2014) (appearing to reach the opposite conclusion); Brant J. Hellwig, *The Constitutional Nature of the United States Tax Court*, 35 VA. TAX REV. 269, 326 (2016) (“The exercise of attempting to definitively locate the United States Tax Court in a particular branch of government proves difficult at best, and at times feels like a hopeless exercise.”).

<sup>26</sup> Some scholars have argued that Tax Court opinions ought to be entitled to *Chevron* deference. See David F. Shores, *Deferential Review of Tax Court Decisions: Dobson Revisited*, 49 TAX LAW. 629, 671 (1996) (noting that it is “difficult to imagine that the body of federal tax law would suffer were courts of appeals to affirm a Tax Court decision based on a reasonable interpretation of the statute”); David F. Shores, *Rethinking Deferential Review of Tax Court Decisions*, 53 TAX LAW. 35, 42 (1999) (“I continue to believe that deferential review would improve the tax litigation system.”); Andre L. Smith, *Deferential Review of the United States Tax Court: The Chevron Doctrine*, 37 VA. TAX REV. 75, 75 (2017) (arguing that “[t]he Tax Court is eligible for *Chevron* deference . . . because it is still within the Executive Branch”). Others have disagreed. See Steve R. Johnson, *The Phoenix and the Perils of the Second Best: Why Heightened Appellate Deference to Tax Court Decisions Is Undesirable*, 77 OR. L. REV. 235, 237 (1998) (arguing that adopting additional deference in the tax system “without any structural change in the tax litigation process would make a flawed system even worse”); Leandra Lederman, *(Un)Appealing Deference to the Tax Court*, 63 DUKE L.J. 1835, 1835 (2014) (“Contrary to some scholarship, this Article argues that, as a doctrinal matter, no vestige of the *Dobson* rule remains and that courts of appeals must apply the same standard of judicial review that they apply to district courts in nonjury cases.”). As a practical matter, decisions of the Tax Court do not currently receive *Chevron* deference.



behaves like other courts, suggesting that de novo review may be appropriate.<sup>27</sup>

Part I discusses the key questions that this Article seeks to answer. Part II describes data and empirical methods. Part III presents results and explanations for those results. Part IV conducts robustness checks to provide assurance that these results are correct. The Conclusion considers possible implications of the results. The Appendix provides additional detail on methods and data.

## I

### KEY QUESTIONS

#### A. *Statutory Judgments or Normative Policymaking?*

*Chevron* famously held that an agency's interpretation of an ambiguous statute warrants deference so long as it reflects a "reasonable policy choice."<sup>28</sup> Many have concluded from this that agencies should make rules based on normative considerations, rather than merely aiming at the "best reading" of a statute. E. Donald Elliott recounts from his tenure at the Environmental Protection Agency's (EPA) Office of General Counsel that, before *Chevron*, the EPA had treated each statute as a "prescriptive text having a single meaning, discoverable by specialized legal training and tools."<sup>29</sup> After *Chevron*, it treated statutes as creating "a range of permissible interpretive discretion," within which "[t]he agency's policy-makers, not its lawyers, should decide which of several different but legally defensible interpretations to adopt."<sup>30</sup>

Peter Strauss put forward an influential version of this view with his idea of "*Chevron* space." He argues that *Chevron* creates a zone of agency discretion for readings of the statute that are "permissible" but not "necessary" under ordinary rules of statutory interpretation.<sup>31</sup> When confronted with several such plausible alternative readings, an

---

<sup>27</sup> This issue has a chicken-and-egg quality, in that the Tax Court likely uses textualist methodology at least in part to follow reviewing courts, since the Tax Court would risk reversal if it remained purposivist like the IRS. In contrast, if Tax Court decisions were to receive deference, the Tax Court would have more freedom to use purposivist methodology with less risk of reversal. So the Tax Court may presently behave like a court because it is treated like a court, without judicial deference. See *infra* notes 74–75 and accompanying text.

<sup>28</sup> *Chevron U.S.A. Inc. v. Nat. Res. Def. Council, Inc.*, 467 U.S. 837, 842, 845 (1984).

<sup>29</sup> E. Donald Elliott, *Chevron Matters: How the Chevron Doctrine Redefined the Roles of Congress, Courts and Agencies in Environmental Law*, 16 VILL. ENVTL. L.J. 1, 11 (2005).

<sup>30</sup> *Id.* at 12; see also Mashaw, *supra* note 2, at 532–33 & nn.71, 73 (discussing the EPA's use of *Chevron* deference).

<sup>31</sup> Peter L. Strauss, "Deference" Is Too Confusing—Let's Call Them "*Chevron Space*" and "*Skidmore Weight*," 112 COLUM. L. REV. 1143, 1163–64 (2012).

agency may select among them, whether on normative policy grounds or statutory grounds, without judicial interference.<sup>32</sup> A number of other scholars have created models following this approach, emphasizing the tradeoff between courts' statutory focus and agencies' normative focus.<sup>33</sup>

Some have pushed back. In particular, Aaron Saiger has argued that agencies "must reject interpretations that [they] conclude[] are interpretively suboptimal, notwithstanding that an ethical, law-abiding reviewing court would acquiesce in those interpretations."<sup>34</sup> In his view, judicial deference to agencies requires those agencies to take on the mantle of the court, which has a duty to "reach the best account it can of what a statute means."<sup>35</sup>

This Article takes no position on whether a normative shift would be appropriate or not. It only remarks that a shift toward normative decisionmaking has been posited much more often than it has been demonstrated. The widespread belief in this normative shift has been supported primarily by anecdote,<sup>36</sup> which is troubling given that it is the main basis for the critique of *Chevron* leveled by current Justices of the Supreme Court.<sup>37</sup>

<sup>32</sup> *Id.*

<sup>33</sup> See Yehonatan Givati, *Strategic Statutory Interpretation by Administrative Agencies*, 12 AM. L. & ECON. REV. 95, 96 (2010) ("In the model, the agency, which maximizes some objective function, adopts a rule that interprets a statute . . ."); Matthew C. Stephenson, *The Strategic Substitution Effect: Textual Plausibility, Procedural Formality, and Judicial Review of Agency Statutory Interpretations*, 120 HARV. L. REV. 528, 535–36, 544 (2006) (assuming that agencies are "interpretive instrumentalists, attaching no intrinsic importance to textual fidelity or analogous concerns" but instead attempting to "secure whatever interpretation would best advance [their] substantive policy agenda"); John R. Wright, *Ambiguous Statutes and Judicial Deference to Federal Agencies*, 22 J. THEORETICAL POL. 217, 226 (2010) (modelling agency action as a function of policy goals).

<sup>34</sup> Saiger, *supra* note 4, at 1233.

<sup>35</sup> *Id.* at 1234.

<sup>36</sup> See David S. Tatel, *The Administrative Process and the Rule of Environmental Law*, 34 HARV. ENVTL. L. REV. 1, 2 (2010) ("[I]t looks for all the world like agencies choose their policy first and then later seek to defend its legality."); Brett M. Kavanaugh, *Fixing Statutory Interpretation*, 129 HARV. L. REV. 2118, 2150 (2016) (book review) ("From my more than five years of experience at the White House, I can confidently say that *Chevron* encourages the Executive Branch (whichever party controls it) to be extremely aggressive in seeking to squeeze its policy goals into ill-fitting statutory authorizations and restraints."); *supra* notes 29–30 and accompanying text.

<sup>37</sup> See, e.g., *Michigan v. EPA*, 135 S. Ct. 2699, 2713 (2015) (Thomas, J., concurring) (complaining that *Chevron* empowers agencies "not to find the best meaning of the text, but to formulate legally binding rules to fill in gaps based on policy judgments made by the agency rather than Congress"); Kavanaugh, *supra* note 36, at 2151 ("*Chevron* invites an extremely aggressive executive branch philosophy of pushing the legal envelope . . . . After all, an executive branch decisionmaker might theorize, 'If we can just convince a court that the statutory provision is ambiguous, then our interpretation of the statute should pass muster as reasonable.'").

Of course, the application of *Chevron* deference to traditional regulatory rulemaking is only part of the story.<sup>38</sup> Subregulatory guidance (including, for the IRS, revenue rulings and revenue procedures<sup>39</sup>) is instead subject to *Skidmore*<sup>40</sup> deference, under which “courts are obliged to take an agency’s view about statutory meaning into account when interpreting statutes the agency administers.”<sup>41</sup> Scholars have disagreed about the implications of *Skidmore* deference for statutory interpretation. Peter Strauss has suggested that it should be rebranded “*Skidmore* weight,” since it is not deference so much as a factor that courts are obliged to consider in their decisions.<sup>42</sup> Connor Raso and William Eskridge describe it as just “mildly deferential,” or

---

<sup>38</sup> This is a relatively recent development with respect to the IRS. Prior to the Supreme Court’s 2011 ruling in *Mayo Foundation for Medical Education and Research v. United States*, it was unclear whether all IRS regulations were subject to *Chevron* deference or whether some might be subject to (weaker) *Skidmore* deference instead. See Mayo Found. for Med. Educ. & Research v. United States, 562 U.S. 44, 56 (2011) (“We see no reason why our review of tax regulations should not be guided by agency expertise pursuant to *Chevron* to the same extent as our review of other regulations.”); see also MICHAEL I. SALTZMAN & LESLIE BOOK, IRS PRACTICE AND PROCEDURE ¶ 3.02[4] (rev. 2d ed. 2002 & Supp. 2019) (describing the rise and fall of tax exceptionalism in judicial deference to IRS regulations); Michael Hall, Note, *From Muffler to Mayo: The Supreme Court’s Decision to Apply Chevron to Treasury Regulations and Its Impact on Taxpayers*, 65 TAX LAW. 695 (2012) (same). If the IRS expected weak *Skidmore* deference rather than stronger *Chevron* deference for some of its regulations prior to *Mayo*, then we might expect the shift toward normative decisionmaking discussed in Section III.A to be even more pronounced at other agencies, where *Chevron* always applied across the board.

<sup>39</sup> While it is widely believed that *Skidmore* deference applies to IRS subregulatory guidance, Kristin Hickman has argued that “because Treasury has construed penalty provisions in the I.R.C. as extending to taxpayer noncompliance with . . . IRB guidance documents, those agency actions carry the force of law” and that therefore “courts should evaluate the legal interpretations advanced in these formats using the *Chevron* standard.” Kristin E. Hickman, *Unpacking the Force of Law*, 66 VAND. L. REV. 465, 471 (2013).

<sup>40</sup> *Skidmore v. Swift & Co.*, 323 U.S. 134, 140 (1944) (holding that subregulatory guidance, “while not controlling upon the courts by reason of their authority, do[es] constitute a body of experience and informed judgment to which courts and litigants may properly resort for guidance”). Despite the terminology, the concept that agency statutory interpretation might be “entitled to very great respect” precedes *Skidmore*. *Edwards’ Lessee v. Darby*, 25 U.S. (12 Wheat.) 206, 210 (1827) (“In the construction of a doubtful and ambiguous law, the contemporaneous construction of those who were called upon to act under the law, and were appointed to carry its provisions into effect, is entitled to very great respect.”); see also, e.g., *Fawcus Mach. Co. v. United States*, 282 U.S. 375, 378 (1931) (holding that contemporaneous construction of an administering agency is “entitled to respectful consideration”); *Swendig v. Wash. Water Power Co.*, 265 U.S. 322, 331 (1924) (same).

<sup>41</sup> Strauss, *supra* note 31, at 1153; see also SALTZMAN & BOOK, *supra* note 38, ¶ 3.03[1][b] (“Prior to the Supreme Court’s decision in *Mead*, some courts applied *Chevron* deference to revenue rulings while others gave no deference whatsoever. After *Mead*, the general consensus is that *Skidmore* is the more appropriate standard . . . . The Supreme Court itself, however, has not expressly ruled on the question . . .”).

<sup>42</sup> Strauss, *supra* note 31, at 1146.

merely “a judicial willingness to go along.”<sup>43</sup> Kristin Hickman and Matthew Krueger argue, based on an empirical study of circuit court cases, that “*Skidmore*’s standard is, as a whole, surprisingly deferential, with courts applying *Skidmore*’s standard to accept agencies’ views at a higher rate than was previously assumed by some scholars.”<sup>44</sup> Finally, Saiger considers this question in the alternative: If *Skidmore* requires courts to give deference, then agencies have a duty (which they may or may not fulfill in practice) to produce subregulatory guidance that is grounded in statutes rather than normative goals.<sup>45</sup> And, in Saiger’s view, even if *Skidmore* does not demand deference, agencies would still be “wise” to emphasize interpretation in order to avoid reversal by courts.<sup>46</sup>

An additional wrinkle specific to tax law is that some tax regulations—those relying on the IRS’s general authority to promulgate regulations,<sup>47</sup> rather than some specific grant of regulatory power in the tax code<sup>48</sup>—were historically thought by some courts to be subject to a lesser degree of deference, known as *National Muffler* deference.<sup>49</sup> *National Muffler* deference was named for a 1979 Supreme Court decision, which held that these tax regulations would receive deference if they “implemented the congressional mandate in some reasonable manner.”<sup>50</sup> Some courts held that *National Muffler* was superseded by *Chevron*,<sup>51</sup> and some held that *National Muffler* and *Chevron* were indistinguishable.<sup>52</sup> But other courts held that *National*

<sup>43</sup> Connor N. Raso & William N. Eskridge, Jr., *Chevron as a Canon, Not a Precedent: An Empirical Study of What Motivates Justices in Agency Deference Cases*, 110 COLUM. L. REV. 1727, 1737, 1744 (2010).

<sup>44</sup> Kristin E. Hickman & Matthew D. Krueger, *In Search of the Modern Skidmore Standard*, 107 COLUM. L. REV. 1235, 1238 (2007).

<sup>45</sup> See Saiger, *supra* note 4, at 1281.

<sup>46</sup> *Id.* at 1283 (“If courts defer under *Skidmore* to agency interpretations they think are interpretively suboptimal, then agencies . . . must promulgate the interpretation they think is interpretively the best. If courts will not accept interpretations with which they do not agree, agencies are . . . usually wise to privilege the courts’ anticipated interpretation over their own . . .”).

<sup>47</sup> I.R.C. § 7805(a) (2018).

<sup>48</sup> The classic example is section 1502 of the Code, which authorizes the Secretary of the Treasury to “prescribe such regulations as he may deem necessary” regarding the taxation of consolidated corporate groups.

<sup>49</sup> Nat’l Muffler Dealers Ass’n v. United States, 440 U.S. 472, 476–77 (1979) (giving deference to an IRS regulation on the definition of a “business league” exempt from federal income tax, because the agency interpretation “harmonizes with the plain language of the statute, its origin, and its purpose”).

<sup>50</sup> *Id.* at 476 (quoting *United States v. Cartwright*, 411 U.S. 546, 550 (1973)).

<sup>51</sup> *Cf.*, e.g., *Hosp. Corp. of Am. & Subsidiaries v. Comm’r*, 348 F.3d 136, 140–41 (6th Cir. 2003) (applying *Chevron* deference instead of *National Muffler* deference).

<sup>52</sup> *Swallows Holding, Ltd. v. Comm’r*, 126 T.C. 96, 131 (2006) (opining that the result under a *National Muffler* analysis would not differ from that under a *Chevron* analysis),

*Muffler* deference continued to apply,<sup>53</sup> essentially as an intermediate level of deference between *Chevron* and *Skidmore*.<sup>54</sup> This view was common until 2011, when the Supreme Court's decision in *Mayo* rejected *National Muffler* deference, conclusively holding that *Chevron* was the appropriate standard.<sup>55</sup>

Here, then, is the overall picture. Trial courts always write decisions with the underlying statutes in mind, not least because they know that reviewing courts will do so. Agencies are generally thought to have greater flexibility to issue regulations and other guidance based on normative criteria than courts, although there is debate over the legitimacy of this approach and historical unclarity about the precise degree of deference accorded to certain tax regulations and sub-regulatory guidance. And, if this theoretical account is descriptively correct, we might expect to see shifts toward normative decision-making after 1944 (*Skidmore*), 1979 (*National Muffler*), 1984 (*Chevron*), and 2011 (*Mayo*), each of which arguably increased the amount of deference accorded to tax regulations.

### B. Textualism or Purposivism?

Once a particular authority has decided to engage in statutory interpretation, the next question will be what *kind* of interpretation it should conduct. Here, the key questions have been whether particular interpreters are more textualist or purposivist and how their practices have changed over time.<sup>56</sup>

---

*vacated*, 515 F.3d 162, 167–68 (3d Cir. 2008) (holding that the Tax Court's application of *National Muffler* was erroneous because it differed from an application of *Chevron*).

<sup>53</sup> See *Snowa v. Comm'r*, 123 F.3d 190, 197–98 (4th Cir. 1997) (holding that *National Muffler* applies to regulations that provide interpretations of congressional language, as opposed to regulations that fill gaps in legislation); *Schuler Indus., Inc. v. United States*, 109 F.3d 753, 754–55 (Fed. Cir. 1997).

<sup>54</sup> Kristin E. Hickman, *The Need for Mead: Rejecting Tax Exceptionalism in Judicial Deference*, 90 MINN. L. REV. 1537, 1557 (2006) (“Although the practical difference is not always apparent, in [jurisdictions according some Treasury regulations only *National Muffler* deference, rather than *Chevron* deference], specific authority regulations are given ‘controlling weight’ pursuant to *Chevron* while general authority regulations promulgated under I.R.C. § 7805(a) are given only ‘considerable weight’ under *National Muffler*.”).

<sup>55</sup> *Mayo Found. for Med. Educ. & Research v. United States*, 562 U.S. 44, 57 (2011) (holding that “*Chevron* and *Mead*, rather than *National Muffler* and *Rowan*, provide the appropriate framework”).

<sup>56</sup> See, e.g., James J. Brudney & Corey Ditslear, *The Warp and Woof of Statutory Interpretation: Comparing Supreme Court Approaches in Tax Law and Workplace Law*, 58 DUKE L.J. 1231 (2009) (evaluating textualism and purposivism at the Supreme Court); Aaron-Andrew P. Bruhl, *Statutory Interpretation and the Rest of the Iceberg: Divergences Between the Lower Federal Courts and the Supreme Court*, 68 DUKE L.J. 1 (2018) (evaluating textualism and purposivism in the Supreme Court, circuit courts, and district courts); Anita S. Krishnakumar, *Statutory Interpretation in the Roberts Court's First Era:*

To place the relationship between textualism and purposivism in context, consider the best-known trend in statutory interpretation: the rise and fall of purposivism at the Supreme Court. The standard story is that modern purposivism took root around 1940, tracking President Franklin Roosevelt's appointment of purposivist Justices and the development of new judicial methodologies to complement the expanded administrative state.<sup>57</sup> Purposivism continued its ascent into the 1970s, which have been described as the "heyday of purposive analysis."<sup>58</sup> But after peaking in the 1970s, purposivism at the Supreme Court sharply declined, thanks to the appointment of textualist Justices by Republican Presidents (especially Justice Scalia in 1986).<sup>59</sup>

Figure 1 illustrates the conventional story, using the same methodology that this Article applies to the IRS and Tax Court below.<sup>60</sup> Each point in the Figure represents the average term frequency of purposivist terms or textualist terms among all Supreme Court cases for the relevant year,<sup>61</sup> normalized to avoid inappropriately emphasizing the absolute magnitude of term frequencies.<sup>62</sup> Because term frequency is inevitably based on the subjective choice of particular terms,

---

*An Empirical and Doctrinal Analysis*, 62 HASTINGS L.J. 221 (2010) (evaluating textualism and purposivism in the Roberts Court).

<sup>57</sup> Nicholas R. Parrillo, *Leviathan and Interpretive Revolution: The Administrative State, the Judiciary, and the Rise of Legislative History, 1890–1950*, 123 YALE L.J. 266, 266 (2013) ("[T]his Article reveals that judicial use of legislative history became routine quite suddenly, in about 1940. The key player in pushing legislative history on the judiciary was the newly expanded New Deal administrative state."); see also, e.g., JOHN W. JOHNSON, *THE DIMENSIONS OF NON-LEGAL EVIDENCE IN THE AMERICAN JUDICIAL PROCESS: THE SUPREME COURT'S USE OF EXTRA-LEGAL MATERIALS IN THE TWENTIETH CENTURY* 187 (1990); Jorge L. Carro & Andrew R. Brann, *Use of Legislative Histories by the United States Supreme Court: A Statistical Analysis*, 9 J. LEGIS. 282, 285 (1982); Nancy Staudt et al., *Judging Statutes: Interpretive Regimes*, 38 LOY. L.A. L. REV. 1909, 1945 (2005).

<sup>58</sup> Anita S. Krishnakumar, *Backdoor Purposivism*, 69 DUKE L.J. 1275, 1277 (2020).

<sup>59</sup> See, e.g., John Calhoun, Note, *Measuring the Fortress: Explaining Trends in Supreme Court and Circuit Court Dictionary Use*, 124 YALE L.J. 484, 498 (2014) ("[T]he sharpest increase in the use of dictionaries began in the mid-1980s, around the time Justice Scalia arrived at the Court."); Paul Clement, Opinion, *Arguing Before Justice Scalia*, N.Y. TIMES (Feb. 17, 2016), <https://www.nytimes.com/2016/02/17/opinion/arguing-before-justice-scalia.html> (describing 1987 as "when Justice Scalia started writing opinions for the court emphasizing the importance of statutory text and the unreliability of legislative history, and that made all the difference").

<sup>60</sup> See *infra* Section II.A (discussing empirical methods in greater detail).

<sup>61</sup> All the Figures in this Article were produced calculating the average of the term frequencies for all judicial opinions (or regulatory documents) for that year, weighted based on the word count of each document. For example, in calculating the textualist score for each year, a Tax Court opinion that is twice as long will count twice as much toward that score.

<sup>62</sup> See *infra* Section II.A (discussing the problems with comparisons of absolute term frequency magnitudes between interpreters).

as explained in greater detail below,<sup>63</sup> the absolute magnitudes of term frequencies are less important than relative magnitudes over time.

For ease of reading, the points are used to generate a trend line, with a ninety-five percent confidence interval represented by the shaded area.<sup>64</sup> These charts are presented as exploratory data analysis rather than reflecting causal inferences, since the year an opinion was written is likely not the primary driver of interpretive methodology so much as it is correlated with deeper shifts in judicial philosophy.

---

<sup>63</sup> See *infra* Section II.A.

<sup>64</sup> The trend lines are generated using locally estimated scatterplot smoothing (LOESS), a non-parametric form of local regression that fits a smooth curve to data points. See WILLIAM S. CLEVELAND, *THE ELEMENTS OF GRAPHING DATA* 168–73 (rev. ed. 1994) (describing LOESS); Bruhl, *supra* note 56, at 57 n.189 (applying LOESS to a similar analysis of term usage, but using a smoothing factor of 0.33 rather than 0.5, resulting in a more tightly fitted curve). I use a smoothing factor of 0.5. *Smoothed Conditional Means*, GGPLOT2, [https://ggplot2.tidyverse.org/reference/geom\\_smooth.html](https://ggplot2.tidyverse.org/reference/geom_smooth.html) (last visited Nov. 8, 2019).

The confidence intervals in Figures 1 through 7, 11, and 20 through 26, are all calculated using bootstrapping. The bootstrapping process used is analogous to the ones described in Section IV.B and Section G of the Appendix. Given a sample of data points (in this case, with years and the term frequency of a particular methodology for that year), bootstrapping recreates a sample of the same size by randomly sampling (with replacement) from the original sample. This is repeated a number of times, here one thousand times, and LOESS curves are recalculated with respect to each bootstrapped sample. For each point on the graph's x-axis (here, each point in time), the values of each bootstrapped LOESS curve are stored and then used to calculate a confidence interval.

The confidence intervals follow the basic bootstrap (also known as the “reverse percentile,” “pivotal,” or “empirical” bootstrap) equation, such that at each point on the x-axis, where  $\theta$  is the LOESS value in the original sample,  $\theta_{0.025}^*$  is the 2.5th-percentile bootstrapped value, and  $\theta_{0.975}^*$  is the 97.5th-percentile bootstrapped value, the confidence interval equals:

$$(2\theta - \theta_{0.975}^*, 2\theta - \theta_{0.025}^*)$$

A.C. DAVISON & D.V. HINKLEY, *BOOTSTRAP METHODS AND THEIR APPLICATION* 194 (1997). Note that the confidence intervals are the confidence intervals of the *curve*, not confidence intervals of *observations*. That is, within each interval with respect to a given point on the x-axis, there is a ninety-five percent probability that the true regression line lies within that interval. But this does *not* imply that there is a ninety-five percent probability that any observation will lie within that interval. The latter probability would be captured by a prediction interval, which would take into account both uncertainty regarding the regression line as well as pointwise variance in the distribution of observations.

FIGURE 1. PURPOSIVIST AND TEXTUALIST TERMS IN SUPREME COURT OPINIONS

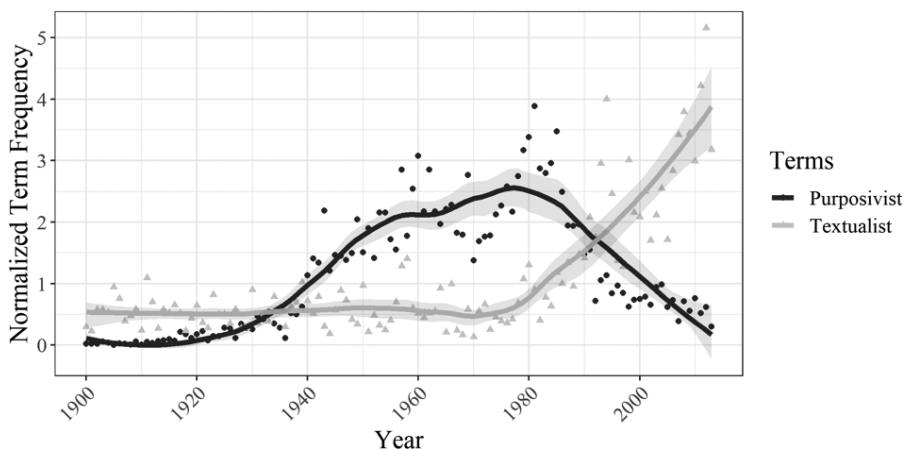
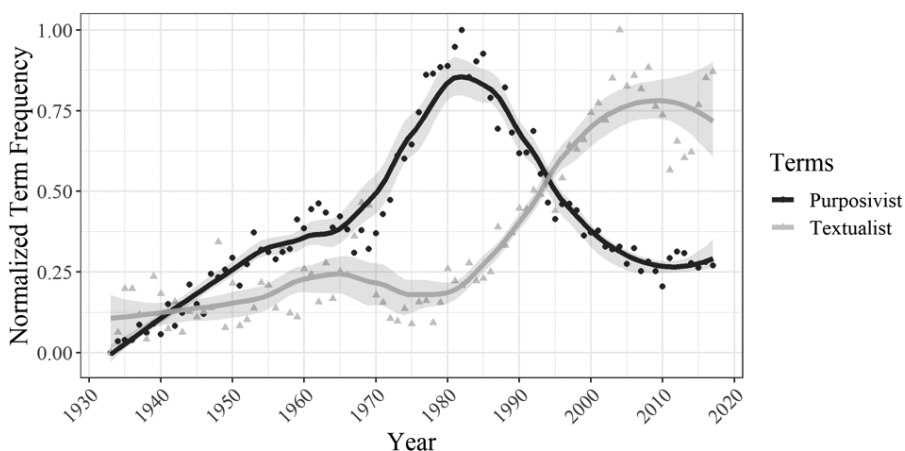


Figure 1 squares neatly with existing literature, showing the same rise in purposivism during the 1930s and 1940s, the peak in the 1970s, and the subsequent decline to the present, accompanied by a sharp uptick in textualism. The fact that Figure 1 is consistent with past scholarship is an early reassurance of the validity of the methods in this Article. Prior empirical research has also concluded that appellate and district courts have followed the same rough directional trend as the Supreme Court, albeit less dramatically and with a slight lag.<sup>65</sup> My methodology again generally confirms this result in Figure 2:

<sup>65</sup> See Bruhl, *supra* note 56, at 1 (“[A]ll federal courts have shifted toward more frequent use of textualist tools in recent decades. However, that shift has been less pronounced as one moves down the judicial hierarchy.”); see also Lawrence Baum & James J. Brudney, *Two Roads Diverged: Statutory Interpretation by the Circuit Courts and Supreme Court in the Same Cases*, 88 *FORDHAM L. REV.* 823, 824 (2019) (generally confirming Bruhl’s findings). But see Abbe R. Gluck & Richard A. Posner, *Statutory Interpretation on the Bench: A Survey of Forty-Two Judges on the Federal Courts of Appeals*, 131 *HARV. L. REV.* 1298, 1309–15 (2018) (arguing that judicial methodology in federal appellate courts is more complicated than the traditional textualist/purposivist divide, but acknowledging the general shift in recent decades toward textualist methods, even by those judges unwilling to self-identify as textualists).



FIGURE 2. PURPOSIVIST AND TEXTUALIST TERMS IN DISTRICT COURT OPINIONS



The decline in purposivism and the ascent of textualism both begin slightly later at the district courts. But the overall modern trend, away from purposivism and toward textualism, is clearly visible at both levels of court.

The crucial empirical question for this Article is whether agencies have followed the courts in their move toward textualism. Most scholars have argued that agencies should remain purposivist as a normative matter, although some have disagreed.<sup>66</sup> But whether they have actually done so is an open question and one that past studies have not attempted to answer.<sup>67</sup> This Article will address this question empirically, exploring more than a century of IRS guidance.

### C. Cohesion Among Courts or Among Specialists?

The pattern of the purposivist/textualist shift at agencies and courts presents competing hypotheses with respect to the Tax Court. If the IRS and generalist courts differ methodologically (and this Article concludes that they do), which will more strongly influence the Tax Court: cohesion with the IRS or cohesion with generalist courts?

The Tax Court handles almost all federal tax cases,<sup>68</sup> operating much like a centralized federal trial court. It takes cases after administrative adjudication by the IRS's internal Office of Appeals,<sup>69</sup> and, if

<sup>66</sup> See *supra* note 2 and accompanying text.

<sup>67</sup> Semet, in particular, did not analyze trends over time, since her study was a snapshot of a fourteen-year period, too short to illustrate long-term methodological trends. Semet, *supra* note 5, at 2282.

<sup>68</sup> See *supra* note 10 and accompanying text.

<sup>69</sup> See generally 26 C.F.R. § 601.106 (2019) (describing the procedures for the Office of Appeals).

cases are appealed from the Tax Court, they are reviewed de novo by the circuit court that had jurisdiction over the taxpayer.<sup>70</sup> Although the Tax Court is an Article I court, the Supreme Court ruled in *Freytag v. Commissioner*<sup>71</sup> that it “exercises judicial power to the exclusion of any other function . . . in much the same way as the federal district courts exercise theirs,”<sup>72</sup> concluding that the Tax Court’s “exclusively judicial role distinguishes it from other non-Article III tribunals that perform multiple functions.”<sup>73</sup>

Given the Tax Court’s judicial role, there is reason to suspect that it would follow the general federal judicial trend toward textualism. More pragmatically, because Tax Court cases are reviewed de novo by circuit courts<sup>74</sup> and because the Tax Court “follows the law of the circuit in which a taxpayer’s appeal would lie,”<sup>75</sup> the Tax Court has every incentive to conform its interpretive practice to that of the courts of appeals. If the Tax Court had remained purposivist, it might have found itself reversed with increasing frequency by textualist-leaning circuit courts.

On the other hand, scholars have long observed that tax law operates differently than other fields of law. In particular, “tax exceptionalists” have argued that federal tax statutes must be read in a more purposivist manner than other federal statutes, due to idiosyncrasies of the tax code or the tax legislative process.<sup>76</sup> Corey Ditslear

---

<sup>70</sup> I.R.C. § 7482(a)(1) (2018); Smith, *supra* note 26, at 78 (“[D]e novo review represents the status quo . . .”).

<sup>71</sup> 501 U.S. 868 (1991).

<sup>72</sup> *Id.* at 891.

<sup>73</sup> *Id.* at 892 (ruling in the context of a dispute over the method for appointing special trial judges). This conclusion is, however, somewhat controversial. The D.C. Circuit’s *Kuretski* ruling appears to cut the other way. See *Kuretski v. Comm’r*, 755 F.3d 929, 932 (D.C. Cir. 2014) (identifying the Tax Court as executive in nature); see also Hellwig, *supra* note 25, at 326 (“The exercise of attempting to definitively locate the United States Tax Court in a particular branch of government proves difficult at best, and at times feels like a hopeless exercise.”).

<sup>74</sup> See *supra* note 70 and accompanying text.

<sup>75</sup> Amandeep S. Grewal, *The Un-Precedented Tax Court*, 101 IOWA L. REV. 2065, 2078 (2016). This is known as the “Golsen rule.” See *Golsen v. Comm’r*, 54 T.C. 742, 757 (1970) (“[W]here the Court of Appeals to which appeal lies has already passed upon the issue before us, efficient and harmonious judicial administration calls for us to follow the decision of that court.”), *aff’d on other grounds*, 445 F.2d 985 (10th Cir. 1971).

<sup>76</sup> See, e.g., Bradford L. Ferguson, Frederic W. Hickman & Donald C. Lubick, *Reexamining the Nature and Role of Tax Legislative History in Light of the Changing Realities of the Process*, 67 TAXES 804, 806–07 (1989) (citing the Code’s complexity, age, extensive legislative history, specialized nature, and specialized drafting process); Mary L. Heen, *Plain Meaning, the Tax Code, and Doctrinal Incoherence*, 48 HASTINGS L.J. 771, 786 & n.73, 818–19 (1997) (arguing against textualism in tax law); Michael Livingston, *Congress, the Courts, and the Code: Legislative History and the Interpretation of Tax Statutes*, 69 TEX. L. REV. 819, 822 (1991) (“The Article argues that the unique characteristics of tax law render generalized theories of interpretation inadequate for tax

and James Brudney have found that, as a descriptive matter, the Supreme Court has been more purposivist in its tax opinions than in other opinions, although they largely attribute this to the influence of Justice Blackmun.<sup>77</sup> And Steve Johnson has directly speculated that the Tax Court's subject matter expertise might free it to apply purposivist techniques, much like the IRS.<sup>78</sup> Since the Tax Court and the IRS are both populated by tax experts, known for their cultural insularity, one might expect them to converge in their interpretive techniques.<sup>79</sup>

Moreover, despite the Supreme Court's view that the Tax Court is "exclusively judicial,"<sup>80</sup> the Tax Court's status as an Article I court carries some distinctions from district courts. Tax Court judges are specialists,<sup>81</sup> they are appointed for limited fifteen-year terms (although they are frequently reappointed),<sup>82</sup> and they can be removed (for cause) by the President, whereas Article III judges must be impeached.<sup>83</sup> Practically speaking, circuit courts might be more reluctant to overturn the judgments (including the purposivist judgments) of specialists than generalists. Consequently, the Tax Court might also differ from other courts for procedural, rather than substantive, reasons.

If the tax exceptionalists are right, or if Article I courts tend to be distinct, then the Tax Court should resist the trend toward textualism and remain purposivist, like the IRS. But if cohesion among courts is the stronger force, then we should expect the Tax Court to trend toward textualism, like other federal courts.

---

cases."); Clinton G. Wallace, *Congressional Control of Tax Rulemaking*, 71 TAX L. REV. 179, 183 (2017) (arguing for a "JCT Canon" under which tax statutes would be interpreted with a special eye toward legislative history generated by the Joint Committee on Taxation (JCT)). Some scholars have resisted the notion of tax exceptionalism. See Paul L. Caron, *Tax Myopia, or Mamas Don't Let Your Babies Grow Up to Be Tax Lawyers*, 13 VA. TAX REV. 517, 518 (1994) (accusing the tax bar and tax scholars of "tax myopia"); Michael Livingston, *Practical Reason, "Purposivism," and the Interpretation of Tax Statutes*, 51 TAX L. REV. 677, 710 (1996) (criticizing "the myth of tax essentialism").

<sup>77</sup> Brudney & Ditslear, *supra* note 56, at 1270–75. Ditslear and Brudney note that "[a]fter Blackmun departed . . . the Court's willingness to invoke legislative history in its tax majorities significantly declined." *Id.* at 1274.

<sup>78</sup> Steve R. Johnson, *The Canon that Tax Penalties Should Be Strictly Construed*, 3 NEV. L.J. 495, 518 (2003) ("It may well be that the Tax Court, as a result of its greater expertise, feels greater confidence in applying the copious interpretive materials that, I have argued, should be the proper bases for construing tax penalty statutes.").

<sup>79</sup> See Caron, *supra* note 76, at 519–31.

<sup>80</sup> Freytag v. Comm'r, 501 U.S. 868, 891 (1991).

<sup>81</sup> Lederman, *supra* note 26, at 1880 ("[T]he Tax Court is specialized—its judges only decide tax cases—and accordingly has greater expertise in tax matters than do other courts.").

<sup>82</sup> See *infra* note 164.

<sup>83</sup> See Smith, *supra* note 26, at 95–96.

## II EMPIRICAL METHODS

To answer these questions, I created a new dataset of all IRS publications ever released, dating back to 1919. The dataset includes all regulatory rulemaking and published subregulatory guidance, but excludes unpublished, non-precedential guidance provided directly to specific taxpayers. The publications were converted to plain text using optical character recognition (OCR) and then cleaned both manually<sup>84</sup> and using computer code—for example, by spell-checking, regularizing whitespace, and removing sections of the publications irrelevant to this Article’s analysis. In addition, I downloaded court data from Harvard Law School’s Caselaw Access Project, a high-quality dataset that includes almost every court case ever decided in the United States until 2015. Section A of the Appendix provides additional detail on the data used in this Article.

### A. Natural Language Processing

The primary measure of interpretive methodology in this Article is the frequency with which agencies and courts cite particular tools, such as legislative history, dictionaries, or canons of construction. This is the dominant approach in existing literature, and maps closely onto conventional conceptions of textualism and purposivism.<sup>85</sup> For

---

<sup>84</sup> In particular, I read through the plain text of each Cumulative Internal Revenue Bulletin to ensure that my code had correctly removed legislative history that did not represent original IRS writing. See *infra* Appendix Section A.1.

<sup>85</sup> See, e.g., Bruhl, *supra* note 56, at 29 (“To a significant degree, the observable difference between competing interpretive approaches lies in which tools they prioritize and emphasize. A judge that uses linguistic canons and dictionaries extensively but uses legislative history sparingly is more textualist than a judge who displays the opposite tendencies.”); Lawrence M. Solan, *Private Language, Public Laws: The Central Role of Legislative Intent in Statutory Interpretation*, 93 GEO. L.J. 427, 453–55 (2005) (citing judicial references to legislative intent as primary evidence of judicial intentionalism). See generally James J. Brudney & Lawrence Baum, *Dictionaries 2.0: Exploring the Gap Between the Supreme Court and Courts of Appeals*, 125 YALE L.J.F. 104 (2015) (studying the frequency of dictionary citations by the Supreme Court and courts of appeals, using word searches); Calhoun, *supra* note 59 (same). For other applications of term frequency analysis not limited to statutory interpretation methodology, see, for example, Keith Carlson, Michael A. Livermore & Daniel Rockmore, *A Quantitative Analysis of Writing Style on the U.S. Supreme Court*, 93 WASH. U. L. REV. 1461, 1478–80 (2016) (using term frequency to evaluate judicial “friendliness”). See generally Daniel Martin Katz et al., *Legal N-Grams? A Simple Approach to Track the Evolution of Legal Language*, in LEGAL KNOWLEDGE AND INFORMATION SYSTEMS: JURIX 2011: THE TWENTY-FOURTH ANNUAL CONFERENCE 167 (Katie M. Atkinson ed., 2011) (using n-gram analysis to track the evolution of legal language); David E. Pozen, Eric L. Talley & Julian Nyarko, *A Computational Analysis of Constitutional Polarization*, 105 CORNELL L. REV. (forthcoming 2020) (using textual analysis to analyze constitutional polarization). More generally, term frequency underlies the “bag-of-words” model that is one of the most common classification schemes used in

example, if a particular document had sixteen phrases relating to legislative history out of 8000 words, the term frequency score for the document with respect to legislative history would be:

$$\frac{16}{8000}$$

A single document might have a positive term frequency score for both textualism and purposivism, or both statutory and normative decisionmaking. Judicial decisions sometimes weigh both textualist and purposivist considerations in the alternative, so this is not uncommon.<sup>86</sup>

Different scholars have different specific definitions of textualism and purposivism, and this Article does not simplistically argue that purposivism is merely the act of using legislative history to interpret statutes. Nevertheless, textualists' skepticism toward legislative history and the general view of purposivism as a philosophy in opposition to textualism makes the use of legislative history a useful proxy for purposivist methodology.<sup>87</sup> In contrast, textualist judges are typically distinguished by their emphasis on the "plain meaning" of statutory text,<sup>88</sup> the use of dictionaries to determine plain meaning,<sup>89</sup> and canons of interpretation.<sup>90</sup>

---

natural language processing and machine learning. See MICHAEL MCTEAR, ZORAIDA CALLEJAS & DAVID GRIOL, *THE CONVERSATIONAL INTERFACE: TALKING TO SMART DEVICES* 167 (2016). This Article uses the bag-of-words model to implement machine learning, which is the standard approach. See, e.g., Pozen, Talley & Nyarko, *supra* (manuscript at 21) (analyzing the frequency with which terms are used without taking into account the context in which they are used). Term frequency also underlies many measures of "similarity" between different documents. See, e.g., Carlson, Livermore & Rockmore, *supra*, at 1483–86 (measuring divergence in judicial writing styles); Elliott Ash & Omri Marian, *The Making of International Tax Law: Empirical Evidence from Natural Language Processing* 16 (Univ. of Cal. Irvine Sch. of Law Legal Studies Research Paper Series, Paper No. 2019-02, 2019), [https://papers.ssrn.com/sol3/papers.cfm?abstract\\_id=3314310](https://papers.ssrn.com/sol3/papers.cfm?abstract_id=3314310).

<sup>86</sup> See, e.g., *Whistleblower 21276-13W v. Comm'r*, 147 T.C. 121, 124 n.8, 128 (2016); *Gardner v. Comm'r*, 145 T.C. 161, 164, 176, 179 (2015).

<sup>87</sup> Bruhl, *supra* note 56, at 29.

<sup>88</sup> See William N. Eskridge, Jr., *The New Textualism*, 37 UCLA L. REV. 621, 623–25 (1990) ("[N]ew textualism posits that once the Court has ascertained a statute's plain meaning, consideration of legislative history becomes irrelevant.").

<sup>89</sup> Bruhl, *supra* note 56, at 29 ("A judge that uses linguistic canons and dictionaries extensively but uses legislative history sparingly is more textualist than a judge who displays the opposite tendencies.").

<sup>90</sup> See, e.g., Gluck & Posner, *supra* note 65, at 1304–05 ("Textualists advanced the canons, in particular, as a more objective and coordinating set of tools for resolving statutory disputes than alternatives like legislative history . . ."); John F. Manning, *Legal Realism & the Canons' Revival*, 5 GREEN BAG (2d ser.) 283, 290 (2002) ("Because textualists believe in a strong version of legislative supremacy, their skepticism about actual intent or purpose has predictably inspired renewed emphasis on the canons of interpretation, particularly the linguistic or syntactic canons of interpretation."). Borrowing from Aaron Bruhl, I divide canons of construction between "substantive canons," "language canons," and "holistic-textual canons." See Bruhl, *supra* note 56, at 26,

The specific terms selected, and the rationales behind them, are described in Section B of the Appendix. The full source code, including all of the phrases used as proxies in this Article, is publicly available online.<sup>91</sup> I conduct several robustness checks in Part IV to ensure that the measures used in this Article are valid; in particular, I spot-check term frequency results in Section IV.A by randomly sampling opinions containing terms I designate as textualist, purposivist, statutory, or normative, in order to ensure that they match conventional conceptions of these methodologies.

All of the analysis in this Article was conducted by downloading bulk data and using Python code to analyze text. Past research has generally relied either on manual tabulation of the occurrences of certain terms, or on searches in Westlaw or Lexis.<sup>92</sup> Programming automates these tasks and makes the analysis more flexible. This enables the application of machine learning techniques, as well as more granular detection and avoidance of false positives and negatives—for example, this Article counts appearances of the phrase “tax administration” but excludes the phrase “effective tax administration” (a term of art referring to a particular type of IRS settlement),<sup>93</sup> which would not be possible using a typical Boolean search without entirely excluding any documents that discuss “effective tax administration.”<sup>94</sup> It also permits more detailed analysis by political affiliation and case outcome<sup>95</sup> and the robustness checks in Part IV.

Most importantly, coding on raw data allows analysis of term frequency—the number of times a phrase appears in a document divided by the word count of the document—rather than a binary analysis of whether or not a phrase appears in the document at all, which is all that is feasible using a word search in Westlaw or Lexis.<sup>96</sup> Word

---

64. The language canons and holistic-textual canons are most closely associated with textualists. *See id.* at 29; *infra* Appendix Section B.2.

<sup>91</sup> *See Code*, JONATHAN H. CHOI, <https://www.jonathanhchoi.com/code-empirical-study> (last updated Mar. 31, 2020).

<sup>92</sup> Bruhl, *supra* note 56, at 30 (“[T]he analyses in this Article rely on electronic searches, primarily in Westlaw, to identify and count cases.”); Solan, *supra* note 85, at 453–54 nn.118–19 (using Lexis searches to assess methodology).

<sup>93</sup> *See* IRM 4.18.3 (Feb. 28, 2017) (defining “Effective Tax Administration Offers”).

<sup>94</sup> For example, in a Westlaw search, one could search for “tax administration,” and one could search for “tax administration % ‘effective tax administration,’” (“%” is the symbol for “not” in Westlaw searches), but the latter search would not pick up a document that *both* included a legitimate occurrence of “tax administration” *and* an occurrence of “effective tax administration.” *See* THOMSON REUTERS, *SEARCHING WITH TERMS AND CONNECTORS* 4 (2009).

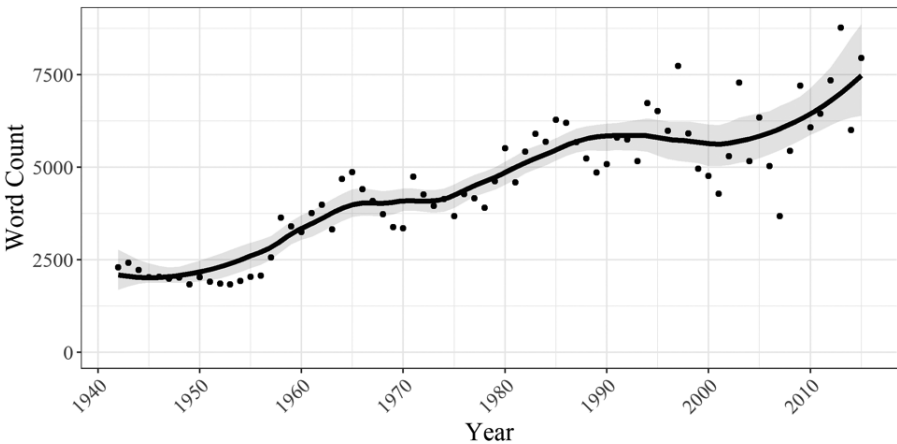
<sup>95</sup> *See infra* Sections III.F, III.G.

<sup>96</sup> Lexis and Westlaw do allow searches for documents that contain a particular term at least a certain number of times. But this would be an impracticably unwieldy method to determine term frequency count, since it would have to be run many times to determine

searches only return the raw number of documents that contain any mention of a particular search term and cannot account for characteristics of the documents retrieved. This means that they cannot consider the number of times a search term appears in the document or the length of the document.

Because the average length of judicial and administrative decisions has varied over time, certain terms might appear more or less purely as a function of greater or lesser detail, rather than due to trends in judicial methodology. For example, as Figure 3 shows, the average length of Tax Court opinions has significantly increased over time. If this phenomenon simply resulted from a trend toward more thorough descriptions of the rationales behind rulings, then a study that counted the number of opinions containing certain tools would overestimate reliance on those tools in later periods.<sup>97</sup>

FIGURE 3. AVERAGE WORD COUNT OF TAX COURT OPINIONS



In addition, a mere count of documents containing a particular phrase cannot measure the “intensity” of that phrase’s usage. It might

---

how many documents use a term at least once, at least twice, at least three times, and so on. For example, to determine normative, statutory, textualist, and purposivist scores just for the IRS would take 339,000 separate manual searches, conservatively assuming thirty occurrences per term per year and that the “at least” search function could be used with proximity searches (which it cannot). (339,000 equals 113 terms, times thirty searches per term per year, times one hundred years.)

<sup>97</sup> Note that while measuring term frequency tends to mitigate this problem, it is not a complete solution. Term frequency merely applies a linear adjustment for word count, but the relationship between interpretive depth and word count is not likely to be perfectly linear. For example, if court opinions less than three thousand words never engaged in any interpretation, but all of the words between the 3000th and the 4000th involved interpretation, then word count minus three thousand (minimum one) would be the more appropriate denominator in calculating the degree of textualism or purposivism in an opinion.

be cited once in passing, or many times as the central rationale to a ruling, but the numerical result would be the same. Term frequency addresses both the problems addressed above: It places a lower value on a phrase that only appears once in a long document, compared to a phrase that occurs many times in that same document.

### B. Machine Learning

This Article uses term frequency analysis to illustrate broad trends, like the Tax Court's movement toward textualism. For more granular analysis on specific interpretive tools, it turns to machine learning. Machine learning, broadly stated, uses algorithms based on a mathematical model to make predictions or decisions without explicit human direction.<sup>98</sup> In doing so, machine learning can uncover trends and test hypotheses that would be onerous or potentially unreliable for humans to analyze manually.

This Article uses a binary classification model in order to test whether the court that wrote a given Tax Court opinion can be identified based on methodology alone. First, each opinion in the dataset is converted from plain text into a "vector" of numbers based on the occurrences of each interpretive tool in that opinion.<sup>99</sup> The classifier must be trained to predict which court wrote a particular opinion based on its vector.<sup>100</sup> To accomplish this, the opinions are randomly divided into a "training set," consisting of 80% of the opinions in the sample, and a "test set," consisting of the other 20%. The classifier repeatedly attempts to classify the opinions in the training set, hundreds of thousands of times, with small tweaks to the classifier between each iteration. The tweaks are retained if they improve performance and discarded otherwise. With each tweak, the classifier iteratively improves (learns) until its performance reaches a maximum.<sup>101</sup>

---

<sup>98</sup> For a general explanation of machine learning methods, see TREVOR HASTIE, ROBERT TIBSHIRANI & JEROME FRIEDMAN, *THE ELEMENTS OF STATISTICAL LEARNING: DATA MINING, INFERENCE, AND PREDICTION* (1st ed. 2001).

<sup>99</sup> All of the machine learning in this Article is conducted using a bag-of-words approach (i.e., analyzing only the terms used, without regard to grammar or word order), using a Python utility provided by the Scikit-Learn project. I specifically use a count vectorizer, term frequency-inverse document frequency transformer with logarithmic term frequencies, and logistic regression (with five-fold cross-validation, five hundred maximum iterations, and refitting). See SCIKIT-LEARN, <https://scikit-learn.org/stable> (last visited Dec. 27, 2019). Section D of the Appendix discusses machine learning methodology in more detail.

<sup>100</sup> This is only a simplified description—in practice, the vector is transformed before it is used to classify data. See *infra* Appendix Section D.

<sup>101</sup> Specific algorithms will vary in how they implement the general concept of iterative improvement, often using mathematical models. See, e.g., Fabrizio Sebastiani, *Machine*



After the training is completed, the performance of the classifier is evaluated using the test set. This entire process is then repeated five times in order to ensure that the results are robust and not dependent on the specific training and test sets chosen.<sup>102</sup> By comparing the classifier's predictions for the test set with the actual classifications, we can produce various metrics of the classifier's predictive abilities. Additional technical detail on the machine learning methodology is provided in Section D of the Appendix.

One widely used measure of predictive performance is the Matthews correlation coefficient (MCC),<sup>103</sup> which produces a score between -1 and +1, where +1 represents perfect correlation (perfect prediction), -1 represents perfect inverse correlation (again, perfect prediction), and 0 represents no correlation (the worst possible score, no better than random). The interpretation of coefficients is highly subjective; however, as an extremely rough rule of thumb, a coefficient represents a weak correlation or no correlation if its absolute value is less than 0.3; a moderate correlation if between 0.3 and 0.7; and a strong correlation if above 0.7.<sup>104</sup>

For completeness, I also list each classifier's "accuracy" (also known as the "correct classification rate"<sup>105</sup>) and "F<sub>1</sub> score."<sup>106</sup> Accuracy is the most intuitive measure of predictive power, representing the percentage of all predictions that were correct.<sup>107</sup> However, it is ill-suited to imbalanced datasets—in an extreme example with ninety-nine observations in category 1, but just one observation in category 2, a classifier that always guessed category 1 would still have an accuracy

---

*Learning in Automated Text Categorization*, 34 ACM COMPUTING SURVS. 1, 10 (2002) (describing the "inductive construction of the classifiers").

<sup>102</sup> See *supra* note 99; cf. George Seif, *Why and How to Do Cross Validation for Machine Learning*, TOWARDS DATA SCI. (May 24, 2019), <https://towardsdatascience.com/why-and-how-to-do-cross-validation-for-machine-learning-d5bd7e60c189> (describing cross validation in machine learning).

<sup>103</sup> See, e.g., Davide Chicco, *Ten Quick Tips for Machine Learning in Computational Biology*, BIO DATA MINING (Dec. 8, 2017), <https://biodatamining.biomedcentral.com/articles/10.1186/s13040-017-0155-3> ("[W]e strongly encourage to evaluate [sic] each test performance through the Matthews correlation coefficient (MCC), instead of the accuracy and the F<sub>1</sub> score, for any binary classification problem." (emphasis omitted)).

<sup>104</sup> See E. GARCIA, A TUTORIAL ON CORRELATION COEFFICIENTS 8–9 (2011). Garcia notes that correlation coefficients may be weaker than initially supposed if degrees of freedom are low due to a small sample size. *Id.* at 10. This is generally not an issue for the tests in this Article, which use relatively large sample sizes.

<sup>105</sup> See, e.g., Pozen, Talley & Nyarko, *supra* note 85.

<sup>106</sup> See KEVIN P. MURPHY, MACHINE LEARNING: A PROBABILISTIC PERSPECTIVE 182–83 (2012) ("Precision measures what fraction of our detections are actually positive, and recall measures what fraction of the positives we actually detected. . . . These are often combined into a single statistic called the F score, or F<sub>1</sub> score, which is the harmonic mean of precision and recall." (emphasis omitted)).

<sup>107</sup> See *id.* at 182.

of 99%. The MCC and, to a lesser extent, the  $F_1$  score, accounts for this problem.<sup>108</sup>

The classification method I use<sup>109</sup> assigns weights to each of the terms in the vocabulary, which allows more granular analysis of how strongly each term is associated with each category—for example, to what degree the “rule of lenity” is associated with the Tax Court or with district courts. In Section III.E, I use these data to produce Figure 8, which illustrates the interpretive tools most characteristic of each court.

### C. Regression Analysis

While natural language processing and machine learning are useful in mapping general interpretive trends and identifying which courts use which particular tools, they are less appropriate in identifying the causal relationship between interpretive methodology and case characteristics. For example, as Section III.F illustrates, casual examination of cases might suggest that Democratic Tax Court judges are less likely to use purposivist terms in their opinions. But this apparent correlation could be caused by other factors, like the year an opinion was written or the year that the judge who wrote it was appointed. When controlling for these factors, the ultimate result is the reverse. Regression analysis allows separate consideration of each of these contributors to methodology.

Section E of the Appendix contains additional technical detail on regression methodology. Because the term frequencies in Tax Court cases do not follow a normal distribution, I rely on two-part regression (logit and a log-transformed generalized linear model) rather than ordinary least squares regression. Sections C through G of the Appendix present additional robustness checks in light of the distributional issues in the dataset.

---

<sup>108</sup> See GARCIA, *supra* note 104, at 8–9. In technical terms, MCC is the only one of the three measures that factors in every quadrant of the “confusion matrix”: that is, true classification into category 1, false classification into category 1, true classification into category 2, and false classification into category 2. See Pierre Baldi et al., *Assessing the Accuracy of Prediction Algorithms for Classification: An Overview*, 16 BIOINFORMATICS REV. 412, 415 (2000) (noting that MCC “uses all four numbers” and therefore “may often provide a much more balanced evaluation of the prediction”). I also correct for imbalanced datasets by undersampling from the over-represented dataset until the sample is evenly balanced between the two categories. See *infra* note 148.

<sup>109</sup> Specifically, I use logistic regression with cross-validation. See *supra* note 99.

### D. Limitations

#### 1. Term Frequency as a Proxy for Methodology

Despite its advantages, term frequency analysis has some limitations. For one, it does not capture whether courts cite a certain interpretive tool approvingly or disapprovingly. A critic might speculate that the Tax Court began to cite textualist tools not in order to follow general judicial trends, but merely to observe and criticize those trends. While reviewing sources to select the terms analyzed in this Article, as well as during the *ex post* checks in Section IV.A, I did not find this to be the case—in fact, I found no disapproving citations of either textualist tools or legislative history in any IRS or Tax Court document.<sup>110</sup> Moreover, a disapproving mention of a specific tool would still suggest that the author considers the tool important to others, even though the author disputes its validity.<sup>111</sup>

A second limitation is that term frequency will not always reflect how important a particular interpretive tool was to the judge's ultimate decision. Legislative history, for example, might be a decisive factor in a court's ruling, even though it is only mentioned once. Or it could be mentioned several times, even though the court ultimately decides the case on other grounds.

To address this limitation, this Article focuses not on absolute results, but on *relative* results. It would be problematic to use term frequency in isolation to assess how textualist a particular court opinion was, or even how textualist the Tax Court as a whole was in any particular year. Instead, this Article always asks how many textualist terms the entire Tax Court used this year compared to last year, or compared to some other court in the same year.

Imagine that dictionaries were infrequently cited by courts but that, when they were cited, they were only mentioned once and with decisive effect. This would imply that term frequency is not a reliable means to compare dictionary use with, say, legislative history—and this Article does not do so. Instead, this Article considers whether an authority cites dictionaries more *over time*. Consequently, a skeptic would need to argue that the way they are cited has changed over time. I have found no evidence of such changes while individually

---

<sup>110</sup> On the other hand, interpreters sometimes cite evidence for one view even if they ultimately decide the other way, but this is to be expected in the ordinary course of statutory interpretation, where different sources may disagree.

<sup>111</sup> See Anita S. Krishnakumar & Victoria F. Nourse, *The Canon Wars*, 97 TEX. L. REV. 163, 182 (2018) (“It is not necessarily the case, for example, that the most frequently invoked interpretive rule is also the most universally accepted. Nevertheless, frequency of judicial invocation does capture an important aspect of what it means to be well-established and entrenched in the legal community.”).

reading cases and agency guidance to validate the terms selected. Moreover, it is reassuring that the term frequency metrics in this Article frequently move in opposite directions. So any hypothesis as to why textualist terms have become more commonly cited at the Tax Court would need to explain why, during the same period of time, purposivist terms have declined at the Tax Court and textualist terms have declined at the IRS.

More broadly, by examining long-term methodological trends, averaged over many different documents and many consecutive years, this Article avoids the idiosyncrasies of single documents and single authors. This again reduces the influence of outlier administrators or judges. The challenge must not merely be that one judge varies her usage between periods, but that all judges vary their usage on average between periods for some reason other than methodological shifts.

## 2. *Doing Different Things, Doing Things Differently*

Differences between the IRS and the Tax Court might arise not from differences in methodology, but in subject matter. For example, perhaps more complex issues inherently demand more purposivist methodology, and perhaps the IRS generally handles more complex matters than the Tax Court. Consequently, methodological divergence between the IRS and the Tax Court may not reflect a difference in their dispositions toward the same interpretive questions, but merely that the IRS and Tax Court serve fundamentally different roles. To borrow Aaron Bruhl's terminology, the IRS and the Tax Court may be both "doing different things" and "doing things differently."<sup>112</sup>

It is undoubtedly true that the IRS and Tax Court do different things in a broad sense. Published Tax Court decisions focus on novel and substantive legal questions, whereas many IRS publications focus on procedural issues. The initial visual presentation of methodological trends in Sections III.A through III.D therefore do not conduct comparisons of the absolute frequencies of particular terms between authorities. Instead, these Sections focus on relative methodological changes *within* various authorities *over time*. In doing so, they avoid the difficulties attending comparisons between different authors.

This approach ameliorates but does not eliminate the problem. Especially over long periods, any interpreter might both change the statutes it interprets (as the statutes themselves are amended) and its interpretive preferences holding statutes constant. For example, Section III.A suggests that the IRS may have become less focused on

---

<sup>112</sup> Bruhl, *supra* note 56, at 6 ("[C]ourts at different levels of the system are both *doing different things* and *doing things differently*.").

statutes as the tax code itself expanded, leaving less statutory ambiguity for the IRS to resolve (doing different things). But Section III.A also suggests that the IRS may have relied on judicial deference to take a more normative approach in reading the tax code (doing things differently).

Similarly, the machine learning analysis in Section III.E compares Tax Court methodology in tax cases with the methodology of District Courts and the Court of Federal Claims across those courts' complete dockets. Again, the distinction between doing different things and doing things differently is blurred; as Section III.E notes, it is likely that both differences play a role in the ability of the algorithm to distinguish opinions written by the various courts. While terms specific to tax law are not included in the analysis, it is hardly surprising that, for example, an area of law dominated by the practice of an agency (the IRS) tends to cite *Chevron* more often.<sup>113</sup> And this finding does not imply that the Tax Court would be more likely to cite *Chevron* than a District Court if they were both interpreting the same statute.

Ultimately, these examples illustrate the difficulty of drawing causal inferences from descriptive term frequency statistics alone. It would be risky to attempt to assess the extent to which different authorities are “doing different things” or “doing things differently” based solely on term frequencies. Instead, I try to tease out causal explanations using historical and primary sources.

The task of this Article is a more modest one, set against the virtual absence of any existing empirical evidence. This Article merely asks whether different courts methodologically differ, for whatever reason. In doing so, it sets a baseline by suggesting that there are indeed substantial differences in interpretive style between different courts. Whether different courts would use different interpretive methodologies when confronted with the *same* statutes remains an important question for future research.

### III RESULTS

Part I describes two major methodological dichotomies: normative policymaking versus statutory fidelity, and textualism versus purposivism. Agencies and courts make “normative” decisions when they justify their decisions on policy grounds, like “fairness” or “efficient administration.”<sup>114</sup> They make “statutory” decisions when they

---

<sup>113</sup> See *infra* fig.8.

<sup>114</sup> See *infra* Appendix Section B.4.

justify their decisions by reference to statutes, as when they “interpret the Code.”<sup>115</sup>

Once an agency or court decides to engage in statutory interpretation, it may further decide to use textualist tools—like dictionaries—or purposivist tools—like legislative history.<sup>116</sup> Finally, an authority that leans either textualist or purposivist might still use different specific statutory tools—one purposivist might emphasize committee hearings, for example, and another might emphasize committee reports. This Part examines variation among authorities and over time, along all these dimensions.

### A. *The IRS Has Become More Normative and Less Statutory*

As between normative and statutory decisionmaking, the IRS has substantially moved over the past century to justify its rulings on normative policy grounds, rather than statutory grounds.

FIGURE 4. STATUTORY AND NORMATIVE TERMS IN IRS PUBLICATIONS

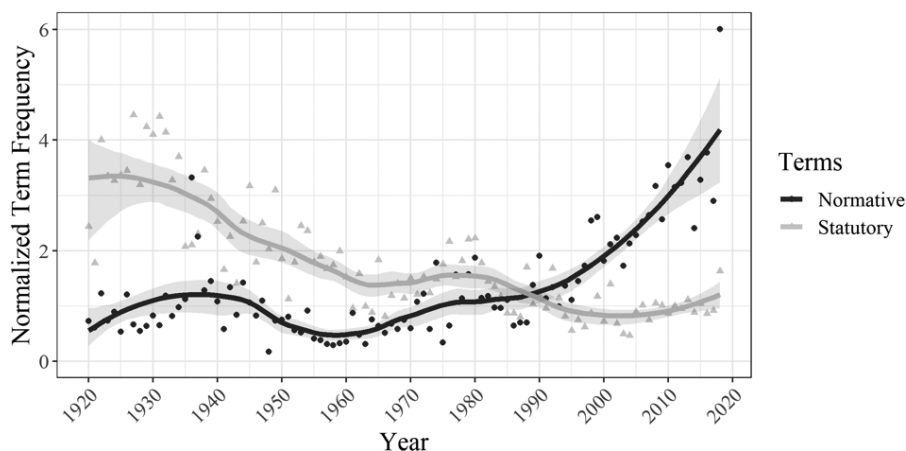


Figure 4 is consistent with the view that the IRS shifted toward a normative perspective as it received more judicial deference over time. The IRS’s use of normative terms markedly accelerated after the 1980s, following *National Muffler* and *Chevron*, peaking after the Supreme Court’s 2011 ruling in *Mayo* (which confirmed that *Chevron* deference should apply to all tax regulations).<sup>117</sup> One could argue that

<sup>115</sup> See *infra* Appendix Section B.3.

<sup>116</sup> See *infra* Appendix Sections B.1, B.2.

<sup>117</sup> See *supra* note 55 and accompanying text.

the normative shift would not have been possible if agencies had continued to constrain themselves strictly to statutory matters.<sup>118</sup>

At the same time, judicial deference is not the sole plausible causal explanation for the rise in normative terms—many significant events in administrative tax law occurred during the 1980s, like the institution of cost-benefit analysis in 1981,<sup>119</sup> the rise of textualism throughout the 1980s (especially the appointment of Justice Scalia in 1986), and the continuing popularization of law and economics through the 1980s.

The shift toward normative decisionmaking also may have been driven by institutional changes specific to tax administration. Prior to 1980, tax regulations were drafted by a group of attorneys in the IRS's Legislation and Regulation Division.<sup>120</sup> These attorneys were experts in tax law, statutory interpretation, and administrative law, but were generally not specialists in specific areas of tax law, like partnership tax or international tax.<sup>121</sup> During the 1980s, this arrangement was flipped, so that tax regulations were generally drafted by subject matter specialists more versed in the practical application of tax regulations than in general principles of statutory interpretation.<sup>122</sup> Moreover, tax legislation during the late 1970s and 1980s increasingly provided the IRS with specific grants of regulatory authority,<sup>123</sup> arguably increasing Treasury's latitude to promulgate rules as it saw fit. Both these changes may have contributed to a shift away from an emphasis on statutes and toward an emphasis on policy concerns.

A related hypothesis is that the IRS has gained expertise over time. The IRS today employs a variety of technical experts, including statisticians, economists, and computer researchers.<sup>124</sup> These specialists might provide the IRS the means to make more sophisticated normative judgments, including more accurately estimating the real-world impact of particular tax policies.

In contrast to the sudden uptick in normative terms, how can we explain the long decline in statutory terms that predated even *Skidmore*? One potential explanation is that as the tax code matured,

---

<sup>118</sup> See *supra* notes 29–30 and accompanying text.

<sup>119</sup> Exec. Order No. 12,291, 46 Fed. Reg. 13,193 (Feb. 17, 1981).

<sup>120</sup> Kristin E. Hickman, *Coloring Outside the Lines: Examining Treasury's (Lack of) Compliance with Administrative Procedure Act Rulemaking Requirements*, 82 NOTRE DAME L. REV. 1727, 1796–98 (2007).

<sup>121</sup> *Id.*

<sup>122</sup> *Id.* at 1798.

<sup>123</sup> *Id.* at 1797 (describing “Congress’s addition of more and more specific authority grants into the I.R.C.”).

<sup>124</sup> See *Research & Analysis*, INTERNAL REVENUE SERV., <https://www.jobs.irs.gov/resources/job-descriptions/research-analysis> (last visited Dec. 30, 2019).

there was less and less statutory ambiguity to be resolved by regulations and rulings. When federal income tax laws were first passed, a greater part of the IRS's work consisted of basic statutory issues, deciding on the correct reading of this or that section of the Code.<sup>125</sup> As the interstices of the Code were filled, the IRS shifted to more granular policymaking details, beginning to offer clarifications of its own regulations rather than original interpretations of statutes.

Regardless of the precise explanation, historical documents reflect the overall narrative that the IRS has grown more normative and less statutory over time. The IRS itself was concerned with a declining focus on faithful interpretation as early as the 1960s. In 1964, it issued a Revenue Procedure stating, in part:

At the heart of administration is interpretation of the Code. It is the responsibility of each person in the Service, charged with the duty of interpreting the law, to try to find the true meaning of the statutory provision and not to adopt a strained construction in the belief that he is "protecting the revenue." The revenue is properly protected only when we ascertain and apply the true meaning of the statute.<sup>126</sup>

This statement was reproduced at the front of every Cumulative Internal Revenue Bulletin from 1970 to 1999, "to emphasize [its] importance to all employees of the Internal Revenue Service."<sup>127</sup> The IRS's stress on faithful interpretation may be responsible for the bump in statutory terms during this period.

But the shift away from statutory decisionmaking has resumed in the past few decades. A survey of recent trends in IRS policy reflects this. The IRS was reformed in the mid-1990s to have an increased emphasis on service to taxpayers and taxpayer rights.<sup>128</sup> In the 2000s, the IRS shifted its focus to the proliferation of abusive multibillion-

---

<sup>125</sup> Cf. Jonathan H. Choi, *The Substantive Canons of Tax Law*, 72 STAN. L. REV. 195, 243 (2020) (describing how, "[d]uring the infancy of the federal income tax . . . statutes were relatively sparse and agency practice was relatively uncertain"). See generally Lawrence A. Zelenak, *Leaving It Up to Treasury: Congressional Abdication on Major Policy Issues in the Early Years of the Income Tax*, 81 LAW & CONTEMP. PROBS. 137 (2018) (describing how the early income tax code was silent or ambiguous on a number of essential issues, leaving them to be resolved at the discretion of the Treasury).

<sup>126</sup> Rev. Proc. 64-22, 1964-1 C.B. 689; see also Rev. Proc. 2012-18, 2012-10 I.R.B. 455 (citing with approval the portion of Revenue Procedure 64-22 discussing administration); Rev. Proc. 2000-43, 2000-2 C.B. 404 (same).

<sup>127</sup> E.g., 1984-1 C.B. ii. The Cumulative Internal Revenue Bulletin is a compilation of all the Internal Revenue Bulletins issued in each year.

<sup>128</sup> See, e.g., Taxpayer Bill of Rights 2, Pub. L. No. 104-168, 110 Stat. 1452 (1996) (codified as amended in scattered sections of 26 U.S.C.) (listing the rights of taxpayers in dealing with the IRS); Internal Revenue Service Restructuring and Reform Act of 1998, Pub. L. No. 105-206, 112 Stat. 685 (codified as amended in scattered sections of 26 U.S.C.) (reforming the IRS with an eye to improving taxpayer service); *Id.* § 1203, 112 Stat. 721



dollar tax shelters.<sup>129</sup> And the most recent movement, following a series of directives by the Trump administration, has been to cultivate regulations that are “simple, fair, efficient, and pro-growth.”<sup>130</sup>

The IRS continues to juggle each of these concerns in its modern policymaking: simplicity, clarity, fairness, efficiency, and, most of all, its central revenue-raising function. It has been aided by an extensive scholarly literature addressing each of these goals.<sup>131</sup> But these are all *normative* goals, not statutory ones. Whether inspired by judicial deference, by modern political trends, or by some combination thereof, the IRS has moved decisively toward normativity in its rulings.

### B. *The Tax Court Has Maintained the Same Proportion of Statutory and Normative Terms*

But what about the Tax Court? Tax Court judges are likely aware of broader trends, like the controversies surrounding tax shelters from the 2000s. At the same time, Tax Court judges are (at least in theory) impartial arbiters not directly responsible to the executive branch,<sup>132</sup> such that their priorities might vary from the priorities of the current administration.

---

(codified as amended in 26 U.S.C. 7804 note) (requiring, generally, that IRS employees be fired if they engage in one of ten kinds of anti-taxpayer conduct).

<sup>129</sup> See, e.g., Joseph Bankman, *Tax Enforcement: Tax Shelters, the Cash Economy, and Compliance Costs*, 31 OHIO N.U. L. REV. 1, 2 (2005) (describing evidence of huge tax shelters in the early 2000s); Tanina Rostain, *Sheltering Lawyers: The Organized Tax Bar and the Tax Shelter Industry*, 23 YALE J. ON REG. 77, 79 (2006) (describing efforts since the late 1990s to fight tax shelters). For a history of the tax shelter movement of the 2000s, see generally TANINA ROSTAIN & MILTON C. REGAN, JR., *CONFIDENCE GAMES: LAWYERS, ACCOUNTANTS, AND THE TAX SHELTER INDUSTRY* (2014).

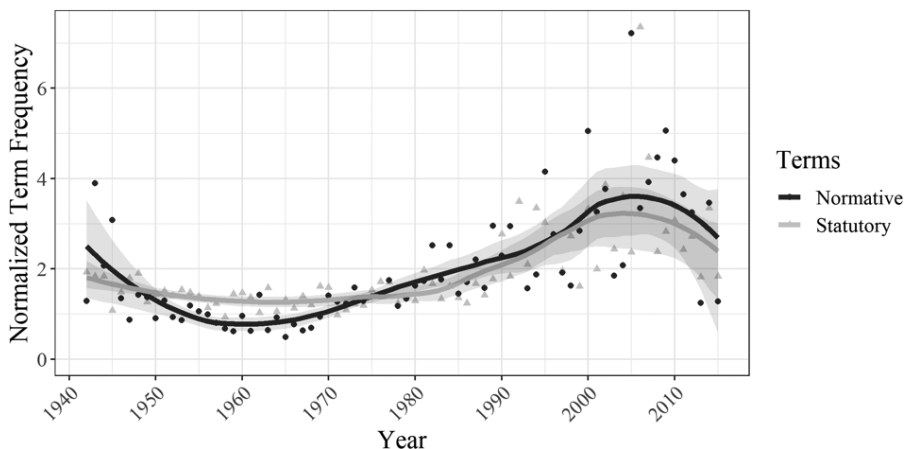
<sup>130</sup> Exec. Order No. 13,789, 82 Fed. Reg. 19,317 (Apr. 21, 2017). The IRS consequently identified and removed 296 regulations that it deemed “no longer necessary because they do not have any current or future applicability.” Eliminating Unnecessary Tax Regulations, 84 Fed. Reg. 9231, 9231 (Mar. 14, 2019). While these executive orders were stated in very general terms, they are nominally binding on the IRS and would have constituted explicit pressure to take normative considerations into account. See Mashaw, *supra* note 2, at 506 (“[B]oth as a practical political and as a normative constitutional matter, we should expect agencies to interpret statutes in the context of presidential direction.”); see also Exec. Order No. 13,777, 82 Fed. Reg. 12,285 (Feb. 24, 2017) (requiring agencies to undertake reforms intended to “lower regulatory burdens on the American people”). See generally Elena Kagan, *Presidential Administration*, 114 HARV. L. REV. 2245 (2001) (discussing presidential direction of agencies).

<sup>131</sup> See, e.g., Lily L. Batchelder, Fred T. Goldberg, Jr. & Peter R. Orszag, *Efficiency and Tax Incentives: The Case for Refundable Tax Credits*, 59 STAN. L. REV. 23 (2006) (discussing efficiency and revenue-raising); John A. Miller, *Indeterminacy, Complexity, and Fairness: Justifying Rule Simplification in the Law of Taxation*, 68 WASH. L. REV. 1 (1993) (discussing simplicity, clarity, and fairness).

<sup>132</sup> Tax Court judges may only be removed for cause. See *supra* note 83 and accompanying text.

It turns out that the Tax Court has remained remarkably steady over the years in its mix between normative and statutory terms. The two have fluctuated within a relatively narrow range for most of the life of the Tax Court. And, importantly, they have generally moved in tandem rather than inversely, in contrast to the IRS.

FIGURE 5. STATUTORY AND NORMATIVE TERMS IN TAX COURT OPINIONS



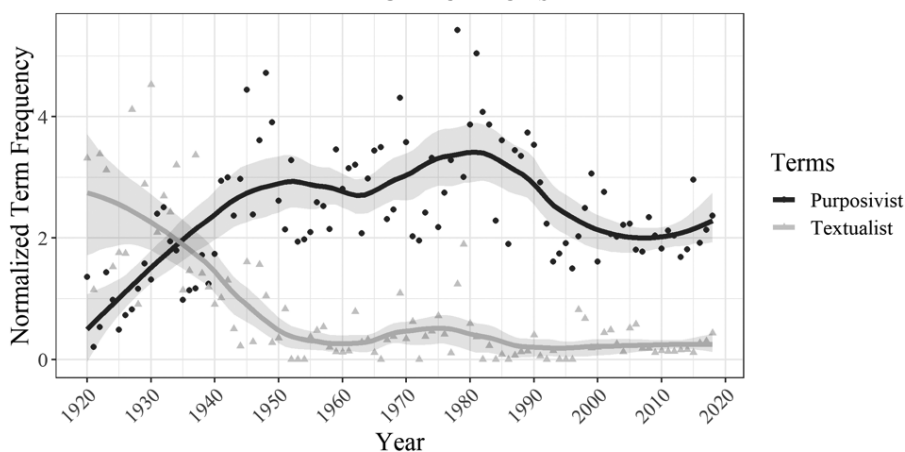
In addition to providing evidence of consistent priorities over time at the Tax Court, Figure 5 also contrasts well with Figure 4, suggesting that the variation shown in Figure 4 is a true effect rather than just noise.

On the other hand, while the Tax Court has remained relatively consistent in the proportion of statutory terms and normative terms it uses in any given year, the frequency of both types of term has changed over time. Both types of term have become more common from a relatively low level during the 1940s through 1970s, to a higher level at present. This could reflect the trend noted in Figure 3, that later Tax Court opinions tend to be longer. If all Tax Court opinions reflect some fixed amount of factual and procedural recitation, longer opinions might cause more space to be allotted to legal analysis. However, this story cannot explain the small dips in the frequency of normative and statutory terms during the 1940s and 2010s. While these dips cannot be attested to with a high degree of confidence—they fall within the confidence intervals, suggesting they may be aberrations rather than true trends—they might be a fruitful subject for future research.

### C. *The IRS Has Become More Purposivist and Less Textualist*

Given that the IRS has significantly reduced the amount of statutory interpretation that it conducts, the next question is whether it has also changed the *type* of statutory interpretation it conducts. As noted above, the Supreme Court and district courts became more purposivist during the 1930s and 1940s but then became less purposivist and more textualist during the 1980s and 1990s.<sup>133</sup> Interestingly, the IRS followed the first shift but not the second, remaining resolutely purposivist despite the rise of the judicial new textualism:

FIGURE 6. PURPOSIVIST AND TEXTUALIST TERMS IN IRS PUBLICATIONS



In fact, in many recent years, the IRS has made almost no references to plain meaning, dictionaries, or the various language canons that I use as proxies for textualism and which the Supreme Court (like the Tax Court, as noted below) has readily adopted. While at first it may appear that the IRS has reduced its use of purposivist terms since 1980, this matches the overall decline in statutory terms discussed in Section III.A. It is worth noting that the IRS's use of textualist terms has declined by at least as much over the same period, such that the relative mix between textualism and purposivism still strongly favors purposivism.

What has prevented the IRS from adopting textualism? My view is that the IRS's close involvement in the legislative process—its role in advising Congress during the drafting of bills<sup>134</sup> and its deep institu-

<sup>133</sup> See *supra* Section I.B.

<sup>134</sup> See, e.g., Parrillo, *supra* note 57, at 266 (“By reason of its unprecedented manpower and its intimacy with Congress (which often meant congressmen depended on agency personnel to help draft bills and write legislative history), the administrative state was the first institution in American history capable of systematically researching and briefing

tional knowledge of the intended meaning of bills<sup>135</sup>—provides it with the means and the motivation to pay special attention to legislative history.<sup>136</sup> This is reflected by the fact that the IRS has chosen to publish legislative history, including relevant reports and hearings, in its Internal Revenue Bulletins since 1941.<sup>137</sup> 1941 marks the original rise of the administrative state as well as purposivism, since specialist agencies like the IRS were able to effectively interpret legislative history in a way that laypeople were not.<sup>138</sup>

This explanation is not specific to tax law—most agencies are involved in the process of drafting statutes and accordingly might have special expertise in interpreting legislative history.<sup>139</sup> Future research could usefully explore whether other agencies have also resisted the modern move toward textualism, like the IRS.<sup>140</sup>

#### *D. The Tax Court Has Become More Textualist and Less Purposivist*

Section I.B observes the movement of federal courts over the past three decades away from purposivism and toward textualism. But the preceding Section indicates that purposivism remains dominant at the IRS. We might ask, as Section I.C does, which is the stronger driver of methodology: cohesion among specialists or cohesion among courts?

---

legislative discourse . . .”); Jarrod Shobe, *Agencies as Legislators: An Empirical Study of the Role of Agencies in the Legislative Process*, 85 GEO. WASH. L. REV. 451, 451 (2017) (finding “that agencies are deeply involved in drafting and reviewing statutory text before enactment, and . . . that Congress often relies heavily on agencies’ significant legislative resources and expertise”); Strauss, *supra* note 31, at 1146 (“The agency may have helped to draft the statutory language, and was likely present and attentive throughout its legislative consideration.”).

<sup>135</sup> See Wallace, *supra* note 76.

<sup>136</sup> See *supra* Section I.B.

<sup>137</sup> See, e.g., 1941-2 C.B. 413–525; 1941-1 C.B. 550–75.

<sup>138</sup> See *supra* notes 57, 134 and accompanying text.

<sup>139</sup> See *supra* notes 134–38 and accompanying text.

<sup>140</sup> An alternative explanation could be that the legislative history of tax statutes might be especially useful due to the work of the JCT. The JCT is a nonpartisan congressional committee that “assists with devising and drafting legislation, and, importantly, produces revenue estimates of every tax provision and prepares explanations of revenue-raising legislative proposals that Congress relies on throughout the legislative process.” Wallace, *supra* note 76, at 183. However, the results in Section III.D *infra* weigh against this explanation. The Tax Court has the same access to JCT publications as the IRS, but it does not participate in the drafting of statutes as the IRS does and did not remain purposivist, unlike the IRS.

FIGURE 7. PURPOSIVIST AND TEXTUALIST TERMS IN TAX COURT OPINIONS

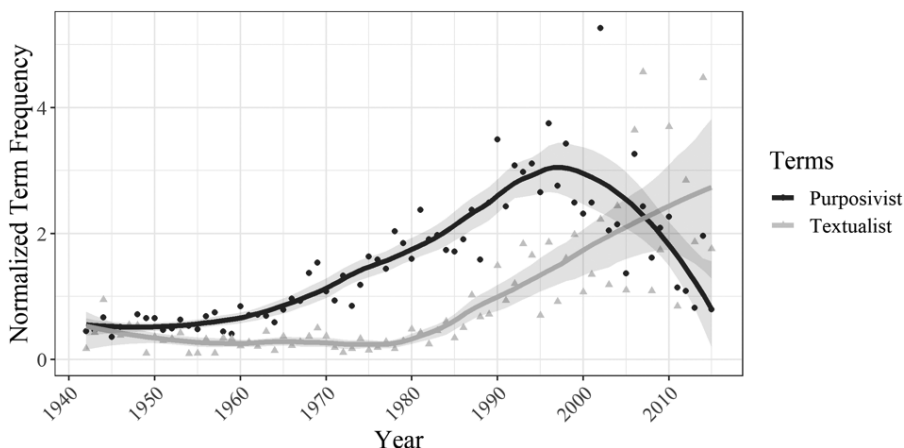


Figure 7 shows that methodological trends in the Tax Court most strongly resemble those in other federal courts. Like district courts, the Tax Court embraced textualist tools in the 1980s and 1990s.<sup>141</sup> The Tax Court also peaked and then declined in its use of legislative history, although it did so approximately a decade later than district courts, in the 1990s rather than the 1980s.<sup>142</sup>

The lag in the Tax Court’s turn away from purposivism is somewhat puzzling. It may be a product of the Tax Court’s continued reliance on certain statutory tools, such as committee reports, in light of their continued use by the IRS and tax experts, or it may reflect reluctance to give up especially useful sorts of legislative history, like the “bluebooks” published by the JCT that summarize legislation in each session of Congress.<sup>143</sup> The following Section addresses this possibility

<sup>141</sup> Cf. Bruhl, *supra* note 56, at 58–61 (using slightly different methodology but finding the same trend).

<sup>142</sup> *Id.* at 57–58. Tax Court data are only available from the court’s founding in 1942, so it is difficult to gauge whether it would have participated in the move toward purposivism around that time. The Tax Court’s predecessor, the Board of Tax Appeals, was an “independent agency in the executive branch” whose “decisions were not final and could be collaterally attacked in federal court,” making it less appropriate for a study of judicial methodology. See Lederman, *supra* note 26, at 1841.

<sup>143</sup> See *Joint Committee Bluebooks*, JOINT COMMITTEE ON TAX’N, <https://www.jct.gov/publications.html?func=select&id=9> (last visited Dec. 31, 2019) (describing the bluebooks and providing access to copies since 1969). Note that the bluebooks might not technically be “legislative history,” since they are produced after legislation has already been passed, but they are considered good evidence of contemporaneous understandings about legislation from the last session of Congress. The Supreme Court explicitly disapproved of using a JCT bluebook as legislative history except “to the extent it is persuasive.” *United States v. Woods*, 571 U.S. 31, 48 (2013).

in greater detail, but it remains a question that might be elucidated by future research.

A reader might also wonder how it is possible for the IRS and Tax Court to diverge methodologically in the first place. Why would a textualist Tax Court not simply strike down guidance issued by a purposivist IRS? Judicial deference surely plays a role here, providing agencies latitude to read statutes differently from courts.<sup>144</sup> The Tax Court is also constrained as a practical matter, since fifteen judges can only do so much to police the voluminous guidance that the IRS produces each week. Finally, professional respect may play a role. Formal deference aside, Tax Court judges may informally feel reluctant to repudiate IRS purposivism even if they would have applied more textualist tools when considering the same question *ab initio*. This Article does not draw any conclusion on the precise causal mechanism for the disconnect between the Tax Court and IRS. Likely each of these explanations plays some role, but future research could usefully investigate further.

As Section II.D.2 discusses, this Article attempts to distinguish “doing different things” from “doing things differently” by focusing on changes in relative term frequency over time. But there is a more pointed potential criticism that remains. *Chevron* tends to sort interpretive issues between those within the *Chevron* space and outside of the *Chevron* space. What if issues within the *Chevron* space require purposivist tools (like legislative history) more often, perhaps because they are more policy-focused and less susceptible to resolution using textualist tools? Since agencies have exclusive jurisdiction over issues within the *Chevron* space, this would imply that the IRS appears more purposivist merely because *Chevron* has precluded courts from addressing the most purposivist questions. Likewise, it could be that courts became more textualist merely because the most purposivist issues were removed from their remit.

There are two main reasons to doubt this account. First, it contradicts most of the theoretical and anecdotal literature discussing the rise of the new textualism. This literature generally attributes the modern resurgence in textualism to the intellectual activity of textualists like Justice Scalia,<sup>145</sup> and the reports of judges that lean toward textualism generally reflect theoretical commitments to the primacy of

---

<sup>144</sup> See *supra* Section I.A.

<sup>145</sup> See, e.g., Abbe R. Gluck, *Justice Scalia's Unfinished Business in Statutory Interpretation: Where Textualism's Formalism Gave Up*, 92 NOTRE DAME L. REV. 2053, 2058 (2017).

statutory text.<sup>146</sup> I am not aware of any commentator or judge who has suggested that judges have become more textualist because they face a different set of issues than they did before *Chevron*.

Second, this sorting imperfectly fits the stories told by Figures 6 and 7. In Figure 6, IRS purposivism did not increase after *Chevron*—instead, it remained flat or perhaps even slightly declined. If more purposivist issues were sorted toward the IRS, we would expect the IRS to increase its use of purposivist tools.<sup>147</sup> In Figure 7, the decline of Tax Court purposivism occurred a decade later than the rise of Tax Court textualism, suggesting that the relationship between the two is more complex than a direct tradeoff due to sorting and that the decline in purposivism was not directly attributable to *Chevron*.

### *E. The Tax Court Has Developed a Unique Interpretive Methodology Relative to Other Courts*

Although the Tax Court has generally become more textualist, the specific flavor of its textualism may differ from other trial courts. I use a machine learning classifier to test whether Tax Court opinions may be distinguished based on interpretive methodology alone.<sup>148</sup> I employ two binary classifications: the Tax Court versus generalist district courts, and the Tax Court versus the Court of Federal Claims (CFC). The CFC is another Article I court that handles claims for monetary damages against the federal government.<sup>149</sup>

Both classifications perform moderately well based on the performance measures described above.<sup>150</sup>

---

<sup>146</sup> See generally Gluck & Posner, *supra* note 65 (reporting results from a survey of appellate judges, including judges that lean toward textualism).

<sup>147</sup> Note, however, that the IRS might not grow more purposivist if it were to solely apply normative criteria inside the *Chevron* space, so that interpretive issues within that space are not sorted to the IRS so much as removed from consideration by any authority.

<sup>148</sup> The results in this Section were produced using Tax Court and district court opinions from 2004 to 2018, the modern textualist era of these courts, in order to obtain current results. Because district courts, taken together, produce many more opinions each year than the Tax Court, the sample used for machine learning would tend to be highly imbalanced in favor of district courts. To correct for this, I randomly “undersample” district court opinions by excluding district court cases at random until the two samples are of the same size. See generally Nitesh V. Chawla, *Data Mining for Imbalanced Datasets: An Overview*, in DATA MINING AND KNOWLEDGE DISCOVERY HANDBOOK 875, 875–83 (Oded Maimon & Lior Rokach eds., 2d ed. 2010) (describing undersampling).

<sup>149</sup> See 28 U.S.C. § 1491 (2018).

<sup>150</sup> Rather than relying on the results from a single iteration of the algorithm, these tables reflect the median values from repeated bootstrapping generated in Section IV.B.

TABLE 1. TAX V. DISTRICT COURT CLASSIFIER PERFORMANCE

MCC	0.546
Accuracy	0.772
F <sub>1</sub> Score	0.762

TABLE 2. TAX V. CFC CLASSIFIER PERFORMANCE

MCC	0.446
Accuracy	0.717
F <sub>1</sub> Score	0.707

This suggests that the Tax Court has indeed produced a style of statutory interpretation distinct from the district courts and the CFC, even though all of these courts have taken a broadly textualist turn.

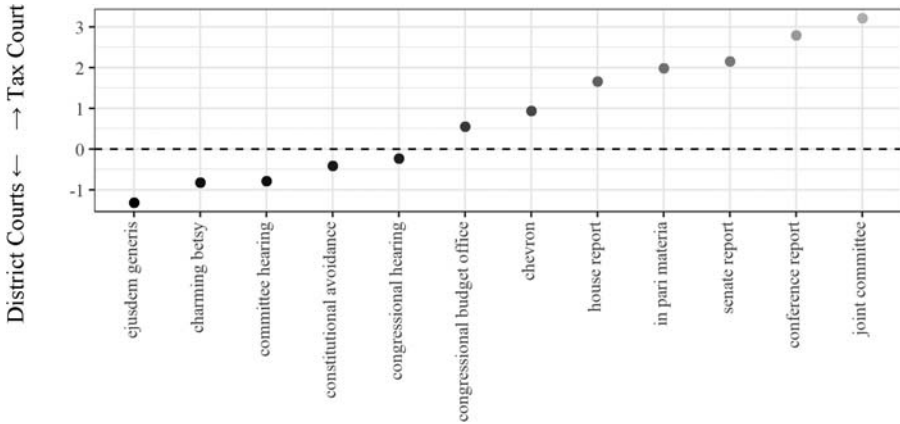
Because the algorithm classifies opinions between the courts by assigning weights to each interpretive term, we can analyze these weights to see which terms are most strongly associated with the Tax Court. Figure 8 presents these weights,<sup>151</sup> evaluating which terms are most predictive of Tax Court opinions (above the dotted line) and which are most predictive of district court opinions (below the dotted line).<sup>152</sup>

<sup>151</sup> These terms were selected because they were statistically significant at the ninety-nine percent confidence level, using bootstrapped confidence intervals, from bootstrapping with one hundred iterations. *See infra* Section IV.B, Appendix Section G (describing bootstrapping to derive confidence intervals in machine learning). No other term was statistically significant above the ninety-five percent level, and consequently they were omitted. As in Tables 1 and 2, the coefficients in Figure 8 are median bootstrapped values.

<sup>152</sup> The listed values are coefficients generated through machine learning, from a logistic regression with log-transformed tf-idf as the independent variables. *See infra* Appendix Section D (discussing tf-idf transformation). Generally speaking, the coefficients should be interpreted as the products of a log-log regression, scaled (by virtue of the tf-idf transformation) so that rarer terms are not disproportionately significant. That is, before scaling, a coefficient of  $\beta$  implies that a  $k$ -fold increase in the frequency of a term is associated with an odds ratio of  $k^\beta$ . *More concretely*, if the coefficient in a log-log regression for a particular term (say, “*in pari materia*”) were 2, then doubling the number of times “*in pari materia*” is used in a case would also increase the probability of a particular case being a Tax Court case to  $2^2 = 4 = 400\%$  of what it would be otherwise. The coefficients in a log-log tf-idf regression can be interpreted in roughly the same manner, but they are scaled to reduce the outsize influence of rare terms. Appendix Section D contains a more specific mathematical description of how the coefficients are calculated.



FIGURE 8. INTERPRETIVE TOOLS, TAX COURT V. DISTRICT COURTS



The results in Figure 8 are intuitively sensible.<sup>153</sup> For legislative history, the Tax Court most heavily prioritizes congressional reports, Congressional Budget Office materials, and materials from the JCT, eschewing congressional hearings. Among textual canons, the Tax Court favors the *in pari materia* canon (requiring that sections of the tax code dealing with similar material “must be construed together”<sup>154</sup>). The prominence of the *in pari materia* canon is not too surprising—many scholars have observed (and approved of) tax authorities’ determination to construe the tax code in a consistent manner,<sup>155</sup> and it is a familiar interpretive move to clarify an ambiguous section of the tax code by reference to other sections.<sup>156</sup> Likewise, other scholars have observed the Tax Court’s reluctance to deploy the *ejusdem generis* canon (requiring that when a general word follows specific words, the general word is assumed to include only words of a similar type<sup>157</sup>—for example, a statute allowing “dogs, cats, and other animals” in a park might not permit tarantulas). Most

<sup>153</sup> See *supra* note 143 and accompanying text.

<sup>154</sup> See, e.g., *Merrill v. Fahs*, 324 U.S. 308, 311, 313 (1945).

<sup>155</sup> See, e.g., *Brudney & Ditslear*, *supra* note 56, at 1298–99 (describing the rule “that when Congress expresses or describes a tax law concept in one part of the Internal Revenue Code, that expression or description should be deemed probative regarding Congress’s treatment of the concept in a separate part of the code”).

<sup>156</sup> See, e.g., *Yates v. Hendon*, 541 U.S. 1, 13–16 (2004); *Drye v. United States*, 528 U.S. 49, 56–57 (1999); *United States v. Reorganized CF&I Fabricators of Utah, Inc.*, 518 U.S. 213, 222–23 (1996); *United States v. Hill*, 506 U.S. 546, 555–56, 556 n.7 (1993); *United States v. Dalm*, 494 U.S. 596, 601–02 (1990); *United States v. Rodgers*, 461 U.S. 677, 695–98 (1983); *United States v. Consumer Life Ins. Co.*, 430 U.S. 725, 745–46 (1977); *Laing v. United States*, 423 U.S. 161, 176–77 (1976).

<sup>157</sup> See *Circuit City Stores, Inc. v. Adams*, 532 U.S. 105, 114–15 (2001) (“[W]here general words follow specific words in a statutory enumeration, the general words are construed to embrace only objects similar in nature to those objects enumerated by the preceding specific words.”).

prominently, courts have expanded the definition of “income” far beyond the initial list of examples provided in the tax code.<sup>158</sup>

For substantive canons, the Tax Court favors *Chevron* deference.<sup>159</sup> (*National Muffler* deference was excluded from the machine learning analysis as a tax-specific standard, which would invariably signal a Tax Court case.) This likely reflects the IRS’s importance as the primary nexus for the administration of federal tax law; deference to its regulations frequently appears in Tax Court cases.<sup>160</sup> Conversely, the Tax Court avoids the constitutional avoidance canon and the *Charming Betsy* canon, which states that “ambiguous congressional statutes should be construed in harmony with international law.”<sup>161</sup> This too makes sense, given that the Tax Court is rarely faced with questions of constitutionality or international law.

#### F. Democratic Judges Are More Purposivist and Republican Judges Are More Textualist at the Tax Court

Past empirical work has frequently asked whether Republican-appointed judges interpret statutes differently than Democrat-appointed judges.<sup>162</sup> Conventional wisdom holds that Republican judges lean textualist, and Democratic judges lean purposivist. This tendency has been observed at the Supreme Court, for example.<sup>163</sup>

<sup>158</sup> See I.R.C. § 61(a) (2018) (listing items that qualify as income); Alice G. Abreu & Richard K. Greenstein, *The Rule of Law as a Law of Standards: Interpreting the Internal Revenue Code*, 64 DUKE L.J. ONLINE 53, 71 (2015) (“[W]hen interpreting the meaning of income, courts often ignore the constraints of *ejusdem generis*.”).

<sup>159</sup> The status of judicial deference regimes as substantive canons is not wholly uncontroversial, but I treat them as such for purposes of this Article without taking a position on that debate. See Raso & Eskridge, *supra* note 43, at 1727 (“As a descriptive matter, we find that deference regimes are more like canons of statutory construction, applied episodically but reflecting deeper judicial commitments, than like binding precedents, faithfully applied, distinguished, or overruled.”). They are, at least, important determinants of how statutes are read, as indicated above.

<sup>160</sup> See, e.g., *N.J. Council of Teaching Hosps. v. Comm’r*, 149 T.C. 466, 478 n.7 (2017) (applying *Skidmore* analysis); *Good Fortune Shipping SA v. Comm’r*, 148 T.C. 262, 275–84 (2017) (applying *Chevron* analysis); *Lindsay Manor Nursing Home, Inc. v. Comm’r*, 148 T.C. 235, 243–61 (2017) (same).

<sup>161</sup> Note, *The Charming Betsy Canon, Separation of Powers, and Customary International Law*, 121 HARV. L. REV. 1215, 1215 (2008).

<sup>162</sup> E.g., Baum & Brudney, *supra* note 65, at 846–47; Brudney & Ditslear, *supra* note 56, at 1301; Krishnakumar, *supra* note 56, at 274–78; David S. Law & David Zaring, *Law Versus Ideology: The Supreme Court and the Use of Legislative History*, 51 WM. & MARY L. REV. 1653, 1671 (2010); Semet, *supra* note 5, at 2289–99, 2316–27.

<sup>163</sup> See, e.g., Law & Zaring, *supra* note 162, at 1654 (“[L]iberal Justices are generally more likely than conservative Justices to cite legislative history.”).

I investigate this issue at the Tax Court by dividing opinions by authorship, between Democratic and Republican appointees.<sup>164</sup> A casual survey of interpretive trends among these judges suggests, if anything, the opposite of the conventional story. The two Tax Court judges who have cited legislative history most often (as of 2015, when the court data were assembled)—Judges Morrison and Wright—were both appointed by Republican Presidents.<sup>165</sup> And the three Tax Court judges who have used textualist tools most often—Judges Lauber, Buch, and Nega—were all appointed by President Obama.<sup>166</sup>

However, we should be skeptical of apparent partisan trends as merely collateral effects of larger time trends. Given that the three most textualist judges (again as of 2015) were appointed by President Obama, the question then becomes whether they are textualist because they were appointed by a Democratic President or because they were appointed recently.<sup>167</sup> That is, to what extent is party affiliation a misleading proxy for the year that an opinion was written or the year the author was appointed?

---

<sup>164</sup> Because Tax Court judges serve fifteen-year terms, sometimes they will be reappointed upon the expiration of their terms. Usually, the reappointing President is of the same party as the President originally appointing the judge. For example, Judge Maurice B. Foley was appointed by President Clinton and reappointed by President Obama, see *Chief Judge Maurice B. Foley*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/foley.htm> (last updated Apr. 23, 2019), while Judge Thomas B. Wells was appointed by President Reagan and reappointed by President George W. Bush, see *Judge Thomas B. Wells*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/wells.htm> (last updated Feb. 13, 2013). A few judges were appointed and reappointed by Presidents from different parties—these judges are not clearly either Democratic or Republican and were therefore excluded for purposes of this analysis. For example, Judges Mary Ann Cohen and Joel Gerber were both appointed by President Reagan and reappointed by President Clinton. See *Judge Joel Gerber*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/gerber.htm> (last updated Apr. 8, 2013); *Judge Mary Ann Cohen*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/cohen.htm> (last updated Oct. 2, 2012). Some Tax Court opinions, particularly memorandum opinions, are written by “special trial judges,” who are appointed by the Chief Judge of the Tax Court rather than by the President. See I.R.C. § 7443A(a) (2018); *Wright v. Comm’r*, 105 T.C.M. (CCH) 1440 (2013); *Madison Recycling Assocs. v. Comm’r*, 81 T.C.M. (CCH) 1496 (2001). Since the ideology of a judge appointed by another judge rather than the President will be more attenuated, these opinions are excluded as well.

<sup>165</sup> Press Release, U.S. Tax Court, Death Announcement - Senior Judge Lawrence A. Wright (Mar. 20, 2000), <https://www.ustaxcourt.gov/press/032000.pdf>; *Judge Richard T. Morrison*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/morrison.htm> (last updated July 19, 2019).

<sup>166</sup> *Judge Albert G. Lauber*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/lauber.htm> (last updated Jan. 2, 2020); *Judge Joseph W. Nega*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/nega.htm> (last updated Sept. 12, 2013); *Judge Ronald L. Buch*, U.S. TAX CT., <https://www.ustaxcourt.gov/judges/buch.htm> (last updated Jan. 14, 2013).

<sup>167</sup> See Gluck & Posner, *supra* note 65, at 1300 (“[Y]ounger judges, who attended law school and practiced during the ascendance of textualism, are generally more formalist and accepting of the canons of construction, regardless of political affiliation.”).

Because interpretive methodology could have multiple determinants, visual analysis of time trends and machine learning classifier analysis is potentially unreliable. The better approach is to conduct regression analysis that controls for variables other than party affiliation. The results of this regression analysis are excerpted in Table 3; Section E of the Appendix provides additional detail on methodology and full tables of results.

TABLE 3. TWO-PART REGRESSION RESULTS FOR PARTY AFFILIATION IN TAX COURT OPINIONS, 1942–2015

	<i>Dependent variable: purposivist terms (per million words)</i>		<i>Dependent variable: textualist terms (per million words)</i>	
<b>Democrat</b>	<b>-44.6</b> <b>(65.8)</b>	<b>146.3***</b> <b>(56.9)</b>	<b>-12.4</b> <b>(8.2)</b>	<b>-14.4*</b> <b>(7.6)</b>
Year Judge Appointed		2.6 (2.8)		-0.09 (0.34)
Taxpayer Wins		-0.3 (42.3)		-5.2 (7.5)
Opinion Year Fixed Effects	No	Yes	No	Yes
<i>N</i>	7308	2760	7308	2479

Note: Each column reflects the combined marginal effects from a two-part regression, excerpted from Tables 9 and 10. The fixed effects row indicate whether dummy variables are included for each opinion year, each judge authoring opinions, or both. *N* varies between regressions because some observations lacked determinate values for some variables (for example, some cases lack a clear winner, since the taxpayer won on some issues and lost on others). Standard errors are clustered by judge. \* denotes statistical significance at  $p<0.1$ , \*\* at  $p<0.05$ , and \*\*\* at  $p<0.01$ .

Table 3 presents the results of regressions that test the effect of party affiliation on the use of purposivist and textualist terms, both alone and with full controls. Without controls, Democratic judges appear less likely to use *both* purposivist and textualist terms, albeit not statistically significantly so. However, as noted above, this could simply be the product of confounding omitted variables. With full controls, Democratic judges are statistically significantly more likely to use purposivist terms (at a 99% confidence level) and statistically significantly less likely to use textualist terms (at a 90% confidence level). In percentage terms, Democrats use 26.6% more purposivist

terms and 6.99% fewer textualist terms,<sup>168</sup> suggesting that party affiliation is an important predictor of interpretive methodology.

### *G. Case Outcomes Do Not Statistically Significantly Predict Interpretive Methodology at the Tax Court*

Scholars have previously studied the determinants of taxpayer wins and losses at the Tax Court, with mixed success.<sup>169</sup> None so far have tested the relationship between interpretive methodology and prevailing party in Tax Court cases. To test this question, I coded the winner in each of the Tax Court cases<sup>170</sup> and included the prevailing party in the regressions analyzing determinants of purposivist and textualist term frequencies.

---

<sup>168</sup> These percentages are semi-elasticities calculated by reproducing the two-part regressions described in Table 3, but modified to use log-transformed dependent variables. The coefficients produced by the regressions were then retransformed to estimate linear marginal effects, the same process used to calculate the coefficients in Table 3. A retransformed coefficient  $\beta$  from such a log-linear regression, when associated with a dummy variable, can be used to calculate the percentage semi-elasticity of a change in the associated dummy variable from 0 to 1, using the formula:  $100 \cdot (e^{\beta} - 1)$ . See Eyal Frank, *Log-Linear Regressions: Three Things To Keep In Mind*, EYAL FRANK (Aug. 22, 2015), <http://www.eyalfrank.com/log-linear-regressions-three-things-to-keep-in-mind>.

<sup>169</sup> See, e.g., Robert M. Howard, *Comparing the Decision Making of Specialized Courts and General Courts: An Exploration of Tax Decisions*, 26 JUST. SYS. J. 135 (2005); James Edward Maule, *Instant Replay, Weak Teams, and Disputed Calls: An Empirical Study of Alleged Tax Court Judge Bias*, 66 TENN. L. REV. 351 (1999); Daniel M. Schneider, *Assessing and Predicting Who Wins Federal Tax Trial Decisions*, 37 WAKE FOREST L. REV. 473 (2002). These studies have generally used case outcomes as the dependent variable in regression analysis, attempting to predict case outcomes based on case characteristics. This Article focuses instead on interpretive methodology, using case outcome as one of several independent variables used to attempt to predict methodology.

<sup>170</sup> The coding was conducted algorithmically, exploiting the statement at the end of every Tax Court decision identifying the prevailing party. When a Tax Court case had no clear winner—for example, if the taxpayer prevailed on some issues and the IRS prevailed on others—the case was excluded from the sample. This analysis considers all Tax Court cases from 1942 to 2015.

TABLE 4. TWO-PART REGRESSION RESULTS FOR CASE OUTCOMES  
IN TAX COURT OPINIONS, 1942–2015

	<i>Dependent variable: purposivist terms (per million words)</i>		<i>Dependent variable: textualist terms (per million words)</i>	
Democrat	146.3*** (56.9)		-14.4* (7.6)	
Year Judge Appointed	2.6 (2.8)		-0.09 (0.34)	
<b>Taxpayer Wins</b>	<b>-0.3 (42.3)</b>	<b>11.7 (35.7)</b>	<b>-5.2 (7.5)</b>	<b>-6.8 (6.7)</b>
Opinion Year Fixed Effects	Yes	Yes	Yes	Yes
Judge Fixed Effects	No	Yes	No	Yes
<i>N</i>	2760	4241	2479	4041
Note: Each column reflects the combined marginal effects from a two-part regression, excerpted from Tables 9 and 10. The fixed effects rows indicate whether dummy variables are included for each opinion year, each judge authoring opinions, or both. When judge fixed effects are included, judge characteristics (party and year of appointment) are omitted as multicollinear. <i>N</i> varies between regressions because some observations lacked determinate values for some variables (for example, some cases lack a clear winner, since the taxpayer won on some issues and lost on others). Standard errors are clustered by judge. * denotes statistical significance at $p < 0.1$ , ** at $p < 0.05$ , and *** at $p < 0.01$ .				

Table 4, again excerpted from the full results in Section E of the Appendix, shows that the relationship between case outcomes and methodology is not statistically significant when controls are included, even at the 90% level.

A reader might wonder whether analysis of case outcomes is meaningful given case selection effects, especially in light of George Priest and Benjamin Klein’s famous claim that “the proportion of observed plaintiff recoveries will tend to remain constant over time regardless of changes in the underlying standards applied.”<sup>171</sup> Because the IRS and the taxpayer have the opportunity to settle prior to judgment,<sup>172</sup> the sample of decided cases may be unrepresentative and

<sup>171</sup> George L. Priest & Benjamin Klein, *The Selection of Disputes for Litigation*, 13 J. LEGAL STUD. 1, 31 (1984).

<sup>172</sup> See *Taxpayer Information: During Trial*, U.S. TAX CT., [https://www.ustaxcourt.gov/taxpayer\\_info\\_during.htm](https://www.ustaxcourt.gov/taxpayer_info_during.htm) (last updated Aug. 27, 2019) (noting that Tax Court trials may be settled even after the trial is complete).

biased if litigants tend to settle in clear-cut cases. For example, it could be that more textualist judges are more likely to rule against taxpayers, but that, anticipating this, taxpayers and the IRS tend to settle cases before textualist judges (on terms favorable to the IRS), so that the only cases that go to trial have countervailing unobserved characteristics that make them close cases (for example, facts that favor the taxpayer). If so, the model in this Article might fail to capture the true relationship between textualism and taxpayer victories.

But there is some reason to doubt that Tax Court cases follow the Priest-Klein model of rational settlement. For one, Tax Court cases are unique in that the taxpayer need not pay any litigated taxes until the case is resolved<sup>173</sup>—so there are benefits to the taxpayer (liquidity and deferral) in litigating even a losing case to the end. These benefits may not have offsetting costs to the government, which is not subject to liquidity constraints and whose litigators may not receive sufficient credit for settling quickly.<sup>174</sup> In addition, most (more than eighty percent of<sup>175</sup>) Tax Court cases involve pro se litigants, whose cost of litigation may be lower than those retaining expensive counsel. And because the factual record generally must be assembled in order to respond to the initial IRS audit, Tax Court cases require less additional factfinding than more traditional court cases, again reducing the marginal costs of going to trial. These unusual features may explain why the IRS wins more than seventy-five percent of Tax Court cases,<sup>176</sup> contrary to the Priest-Klein hypothesis, which predicts that trial win rates will follow “a strong bias toward . . . 50 percent.”<sup>177</sup>

Regardless of whether the Priest-Klein model applies, the failure to find a statistically significant relationship is *not* strong evidence that such a relationship does not exist, and this Article does not affirmatively claim that interpretive methodology has no relationship with case outcomes. Moreover, even if there is no consistent predictive relationship between case outcomes and interpretive methodology, methodology could still have an important effect on substantive case outcomes. It could be that every time a dictionary is cited, it decisively

---

<sup>173</sup> See 1 TAXPAYER ADVOCATE SERV., NATIONAL TAXPAYER ADVOCATE ANNUAL REPORT TO CONGRESS 2018, at 295 (2019), [https://taxpayeradvocate.irs.gov/Media/Default/Documents/2018-ARC/ARC18\\_Volume1.pdf](https://taxpayeradvocate.irs.gov/Media/Default/Documents/2018-ARC/ARC18_Volume1.pdf) (“The U.S. Tax Court is the only prepayment judicial forum for taxpayers to resolve their disputes with the IRS.”).

<sup>174</sup> This is essentially a principal-agent problem: Even if government litigators receive credit for avoiding litigation costs of trials, they may not receive credit for bringing in tax revenue earlier than if the trial had not occurred.

<sup>175</sup> 1 TAXPAYER ADVOCATE SERV., *supra* note 173, at 295 (“More than 80 percent of cases in Tax Court are brought by unrepresented taxpayers . . .”).

<sup>176</sup> See *infra* tbl.6.

<sup>177</sup> Priest & Klein, *supra* note 171, at 5.

determines the prevailing party, but the prevailing party is equally likely to be the IRS or the taxpayer. In a well-functioning judicial system, this is in fact desirable—the absence of systemic bias is reassuring rather than a sign that interpretive methodology is superfluous.

## IV

### ROBUSTNESS CHECKS

#### A. *Reading Cases to Confirm Term Frequency Results*

To confirm that term frequency results correspond with conventional notions of textualism, purposivism, statutory interpretation, and normativity, I pulled forty Tax Court opinions and manually evaluated how the terms were used in those opinions. Although I spot-checked each term more informally while producing the list of proxies for each methodology, this Section describes an additional ex post check to ensure the robustness of this Article's methods.

The dataset contained seventy-four years of opinions (1942–2015), which I separated into ten similarly sized time periods. For each methodology, I pulled one opinion at random from each period and reviewed it to confirm that the methodology was used as expected. The full list of these opinions is available online, along with specific details and citations for the methodologies used in each opinion.<sup>178</sup>

#### B. *Bootstrapped Confidence Intervals for Machine Learning*

MCC, Accuracy, and F<sub>1</sub> Score generally tell us about the *magnitude* of the differences between courts that can be captured by a machine learning classifier. But an important measure to determine the robustness of these results is whether they are *statistically significant*, that is, whether the classifier performs better than chance.

To this end, I employ a “bootstrapping” design that repeatedly tests the machine learning algorithm on a resampling of the data. By testing how much the estimates of classifier performance vary between tests, we can calculate the standard error of the test and derive confidence intervals. The algorithm is statistically significantly different from zero—that is, its performance is better than random chance—if the confidence interval for MCC excludes zero, and if the intervals for Accuracy and F<sub>1</sub> Score exclude 0.5 (50%).

Figures 9 and 10 present confidence intervals for each classifier performance metric after bootstrapping with one thousand test itera-

---

<sup>178</sup> *Online Appendix: Spot-Checking Terms*, JONATHAN H. CHOI, <https://www.jonathanhchoi.com/s/Spot-Checking-Terms-10172019.pdf> (last updated Oct. 17, 2019).



tions. For each metric, the median value is represented by the white circle, the 95% confidence interval is represented by the black inner bars, the 99% confidence interval is represented by the grey outer bars, and the null hypothesis (the value that would be generated by a classifier performing no better than chance) is represented by the dashed line.

FIGURE 9. BOOTSTRAPPED CONFIDENCE INTERVALS, TAX COURT V. DISTRICT COURTS

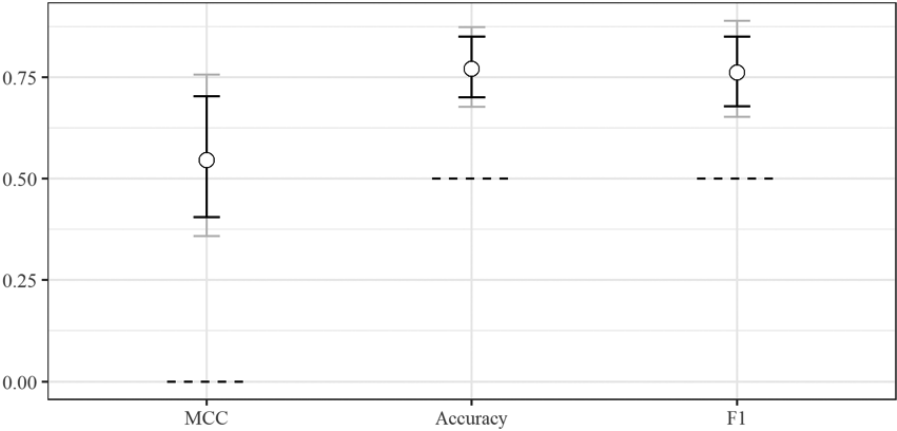
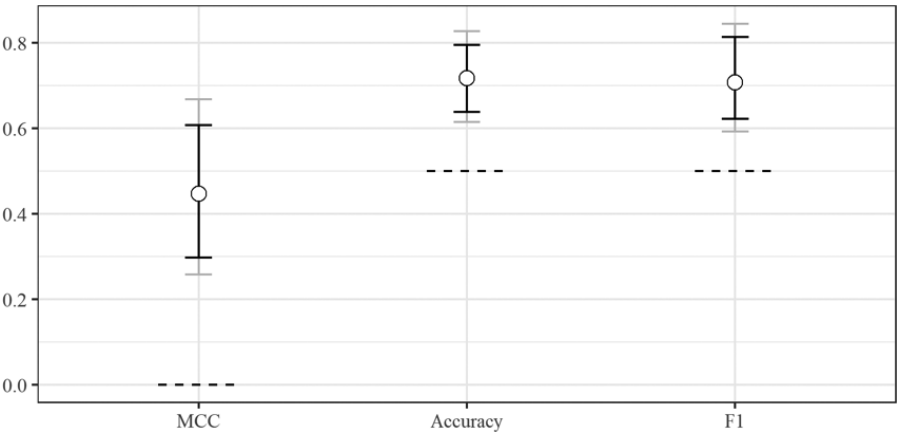


FIGURE 10. BOOTSTRAPPED CONFIDENCE INTERVALS, TAX COURT v. CFC



Figures 9 and 10 demonstrate that for each of the performance metrics—MCC, Accuracy, and F<sub>1</sub> Score—the classifier performs statistically significantly better than chance at a 99% confidence level, providing additional assurance of the results in Section III.E. Section

G of the Appendix contains additional detail on the bootstrapping calculations.

### C. *Validating OCR Quality over Time*

Another potential concern is that apparent trends could be produced merely by variation in the quality of computer OCR over time. This could introduce systematic bias if, for example, older documents were written in text that is more difficult to scan, or the quality of the records degraded over time (due to stains, tears, etc.). If (hypothetically) there were a ten percent chance that any particular word in the 1925 Cumulative Internal Revenue Bulletin were misspelled and therefore not identified, but a zero percent chance in 2018, the matching rate in 2018 would be overstated relative to 1925 by ten percent. A skeptical reader might particularly doubt the dataset of Internal Revenue Bulletins produced specifically for this Article.

I technologically mitigate this issue by using spell-checking to correct obvious errors.<sup>179</sup> But this is not a complete solution—for example, again purely hypothetically, if the OCR rendered the word “the” as “tbo,” the spell-checker would not correct that misspelling.<sup>180</sup>

One way to judge the variation in spelling errors over time, which may be a proxy for OCR quality, is to examine the ratio of terms—purposivist, textualist, normative, and statutory—before and after spell-checking. Figure 11 depicts this ratio, obtained for any year by taking the count of all terms examined in this Article in the Cumulative Internal Revenue Bulletin before conducting spell-checking, divided by the count of such terms after conducting spell-checking.<sup>181</sup>

---

<sup>179</sup> See *supra* Part II.

<sup>180</sup> This is because the spell-checking algorithm used only fixes words that are incorrect by one character (that is, whose Levenshtein distance is one). “Tbo” is different from “the” by two characters—in fact, it would likely be corrected as “to” rather than as “the.”

<sup>181</sup> On some occasions, the ratio will exceed 100%. This could happen if, for example, the misspelled word “mcode” were corrected to “mode.” In this case, the misspelling could be registered as an instance of “code,” which (if used in conjunction with a word like “interpret”) would be registered as a statutory term, while the corrected spelling would not be. For Figure 11, the ratio is capped at 1.00.

FIGURE 11. RATIO OF ALL TERMS BEFORE AND AFTER SPELL-CHECKING

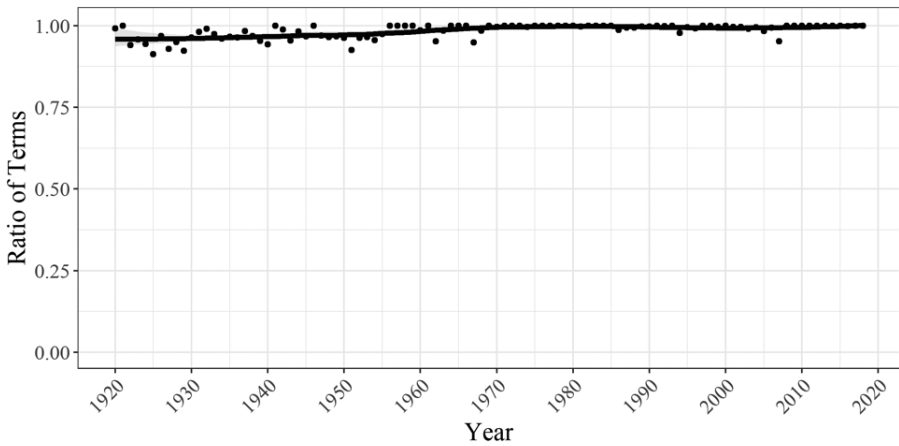


Figure 11 shows that earlier years contain more misspellings, as might be expected. The misspellings are mostly concentrated in the very earliest years: the ratio stabilizes after the late 1960s, with most years afterward at 1.00 (implying that zero spelling errors were found for the entire year).

Figure 11 therefore suggests that the most significant recent trends analyzed by this Article—especially the increase in normative terms during the 1980s—are unlikely to be the product of variation in OCR quality. This is also borne out by the fact that a number of the most prominent trends in this Article are declines in term frequency, such as the decline in statutory terms at the IRS that began in the 1920s.<sup>182</sup> The increase in OCR quality over time suggests that, if anything, these declines might be understated.

#### D. Confirming that Results Are Not Driven by Changes in Terminology

A final potential issue is that interpreters might become more or less familiar with the formal names of interpretive concepts, without necessarily relying more or less on the underlying concept. For example, it could be that judges in the 1930s applied the *ejusdem generis* canon without labeling it as such, whereas judges today are taught *ejusdem generis* and refer to it by name. Thus, many of the apparent trends in this Article—especially those concerning textualist terms with relatively obscure Latin titles—could purely reflect changes in terminology, rather than methodology.

<sup>182</sup> *Supra* fig.4.

One category of textualist term that resists this critique, however, is dictionaries. It would be difficult for a court to cite a dictionary without actually using the word “dictionary” (or some comparable publication covered by this Article, like “World Book” or “Linguae Britannicae”). Consequently, we can compare trends in citations to dictionaries against trends in the use of all other textualist terms (that is, language canons and holistic-textual canons)—if trends in textualist terms are driven by changes in terminology, the two should diverge.

FIGURE 12. DICTIONARIES AND OTHER TEXTUALIST TERMS IN IRS PUBLICATIONS

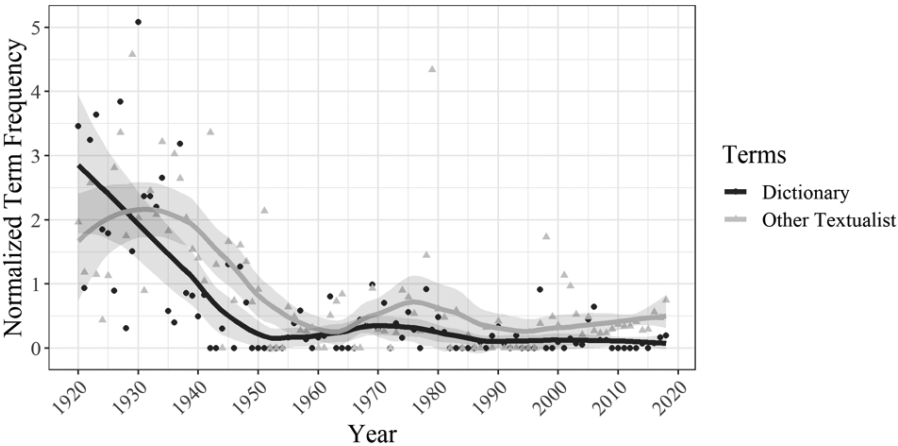
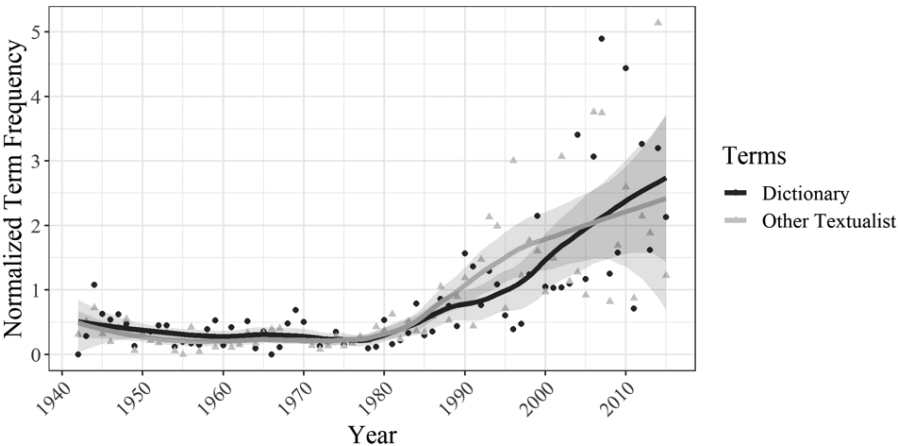


FIGURE 13. DICTIONARIES AND OTHER TEXTUALIST TERMS IN TAX COURT OPINIONS



As Figures 12 and 13 show, the overall trends are similar—consequently, the results in this Article for textualist terms appear to reflect true changes in methodology, rather than just terminology.

### CONCLUSION

Most statutory interpretation occurs at agencies, rather than courts.<sup>183</sup> But little empirical scholarship exists on how agencies interpret statutes,<sup>184</sup> and none has contrasted the methodologies of agencies and courts. This has left jurists and scholars in the dark while grappling with some of the most important questions in modern jurisprudence, including the effect of judicial deference doctrines like *Chevron*.

This Article uses the IRS and the Tax Court as case studies in administrative and judicial statutory interpretation. It concludes that the two differ substantially. First, the data show that the IRS has shifted toward normative policy judgments in its decisionmaking, less often engaging in statutory interpretation at all. Second, the data show that the IRS has become more purposivist over time when interpreting statutes, unlike the now-textualist Tax Court.

This Article also has implications for the study of tax law. It helps taxpayers better to tailor their arguments before the IRS and the Tax Court. Moreover, it provides evidence confirming the “exclusively judicial role” that the Supreme Court has controversially held the Tax Court to play,<sup>185</sup> in that the Tax Court reads statutes more like other courts than like the IRS.

Finally, this Article complicates the standard story of tax exceptionalism. On one hand, the two primary interpreters of federal tax law significantly differ in their methodology, so that tax law is not uniformly more purposivist than other fields, as many scholars have proposed.<sup>186</sup> On the other hand, although the Tax Court has become more textualist in general, it favors different specific interpretive tools than other courts,<sup>187</sup> suggesting that while certain authorities may be purposivist or textualist in broad terms, each may adopt its own flavor of purposivism or textualism.

---

<sup>183</sup> See Mashaw, *supra* note 2, at 502–03 (describing agencies as “the primary official interpreters of federal statutes”).

<sup>184</sup> See *supra* notes 5–6 and accompanying text.

<sup>185</sup> *Freytag v. Comm’r*, 501 U.S. 868, 892 (1991); see also *supra* note 73 and accompanying text. The judicial status of the Tax Court is important because it implies that Tax Court opinions are subject to de novo review, rather than deferential review, as an agency determination would be. In addition, it has implications for the appointments process at the Tax Court. See *supra* note 73.

<sup>186</sup> See *supra* notes 76–78 and accompanying text.

<sup>187</sup> See *supra* Section III.E.

## APPENDIX

A. *Data Sources*

All of the Python code used in this Article is available for reference online.<sup>188</sup> All of the data used in this Article are available upon request, except for court opinions that I am prohibited from sharing under the terms of my researcher license, as described below.<sup>189</sup>

1. *IRS Publications*

The IRS publications used in this Article were extracted from two sources. First, I downloaded all of the Cumulative Internal Revenue Bulletins, published annually by the IRS from 1919 until 2008, from the website of the U.S. Government Publishing Office.<sup>190</sup> Second, I downloaded all of the Internal Revenue Bulletins posted on the IRS's website, which include the years from 2003 until the present.<sup>191</sup> Both sources provide files in .pdf format, which I converted to plain text using Adobe's OCR software. I found alternative OCR software to produce the same or slightly worse results. The OCR was of reasonably high quality, but, to ensure accurate term frequency counts, I also wrote a program to conduct pre-processing (removing whitespace, regularizing capitalization, fixing hyphenation across pages, and conducting spell-checking).<sup>192</sup> Where the documents could not feasibly be processed using algorithms, I edited them manually (for example, to remove irrelevant material such as legislation and legislative history). The beginning and ending years, 1919 and 2019, were omitted as partial years that might be biased if IRS guidance follows an annual cycle.

Internal Revenue Bulletins include all official IRS publications for each year—regulations, revenue rulings, revenue procedures, and other miscellaneous statements. They do not include unpublished guidance on which taxpayers (other than the petitioner) are not generally entitled to rely—for example, private letter rulings issued to particular taxpayers that services such as *Tax Notes* may obtain through FOIA requests. The Internal Revenue Bulletins do contain copies of all tax legislation enacted for the year, along with relevant

---

<sup>188</sup> *Code*, *supra* note 91.

<sup>189</sup> See *infra* Appendix Section A.2.

<sup>190</sup> U.S. GOV'T PUB. OFF., <https://www.govinfo.gov> (last visited Dec. 4, 2019).

<sup>191</sup> *IRS Online Bulletins*, INTERNAL REVENUE SERV., <https://www.irs.gov/irb> (last visited Dec. 4, 2019).

<sup>192</sup> I used the *pyspellchecker* library in Python, version 0.4.0, with a Levenshtein distance of 1, after excluding any terms analyzed in this Article. See *Pyspellchecker 0.5.3*, PYPI, <https://pypi.org/project/pyspellchecker> (last updated Nov. 25, 2019).

committee reports.<sup>193</sup> Since the tax legislation and legislative history were not original material produced by the IRS, I removed these from the documents for purposes of this Article.

This Article analyzes regulations and subregulatory guidance together. Historically, the line between different types of guidance has sometimes been fuzzy, and the significance of each type of guidance has changed over time. There was little formal distinction between regulations and subregulatory guidance before the Administrative Procedure Act (APA) was passed in 1946<sup>194</sup> and especially before the Federal Register Act was passed in 1935.<sup>195</sup> Even after the APA, most tax regulations are designated “interpretative” by the Treasury and are allegedly not subject to notice-and-comment requirements, again making them hard to distinguish from subregulatory guidance.<sup>196</sup> Moreover, many of the changes in IRS regulatory practice were endogenous to broader political movements that I am trying to capture in this Article—for instance, the passage of the APA was the culmination of years of New Deal politics,<sup>197</sup> the same politics that produced the shift toward purposivism that is a primary finding of this Article.

Given that it would be difficult and perhaps undesirable to disaggregate different types of guidance, I have analyzed all published tax guidance together. The fact that so many of the results discussed in this Article move in opposite directions suggests that this has not

---

<sup>193</sup> The IRS began to publish committee reports in 1939. *See* 1939-1 C.B. pt. 2, at 1. Its decision to publish committee reports may contribute to, or may reflect, the IRS’s general emphasis on committee reports as indicia of legislative history.

<sup>194</sup> Pub. L. No. 79-404, 60 Stat. 237 (1946).

<sup>195</sup> Pub. L. No. 74-220, 49 Stat. 500 (1935).

<sup>196</sup> *See, e.g.,* Kristin E. Hickman, *Coloring Outside the Lines: Examining Treasury’s (Lack of) Compliance with Administrative Procedure Act Rulemaking Requirements*, 82 NOTRE DAME L. REV. 1727, 1729 (2007) (“Treasury also contends, however, that most Treasury regulations are interpretative in character and thus exempt from the public notice and comment requirements by the APA’s own terms.”). Most Treasury regulations, including interpretative regulations, do ultimately go through notice and comment, although often after they have entered into effect as temporary regulations. *Id.* at 1730–31. Many critics have nevertheless alleged that interpretative tax regulations lack force of law. *See, e.g.,* Stanley S. Surrey, *The Scope and Effect of Treasury Regulations Under the Income, Estate, and Gift Taxes*, 88 U. PA. L. REV. 556, 557 (1940) (arguing that what was then section 62 of the Internal Revenue Code “d[id] not invest interpretative regulations with the force of law”); Steve R. Johnson, *Intermountain and the Importance of Administrative Law in Tax Law*, TAX NOTES, Aug. 23, 2010, at 837 (“Interpretive regulations do not have force of law; they merely inform the public of what the agency believes the statute means.”).

<sup>197</sup> *See* George B. Shepherd, *Fierce Compromise: The Administrative Procedure Act Emerges from New Deal Politics*, 90 NW. U. L. REV. 1557, 1560–61 (1996) (“The APA was a cease-fire armistice agreement that ended the New Deal war on terms that favored New Deal proponents.”).

biased the results by, for example, inflating the proportion of strictly procedural matters over time. However, more granular analysis of more specific slices of published guidance—such as the “legislative” regulations that must go through conventional notice and comment<sup>198</sup>—would be an interesting project for future research.

## 2. *Court Opinions*

The court opinions analyzed in this Article were downloaded from the Caselaw Access Project, a joint project of the Harvard Law School Library and Ravel Law.<sup>199</sup> The Project is an extensive and high-quality database that contains “nearly all cases from an American court” between 1658 and 2018.<sup>200</sup> In order to write this Article, I obtained a researcher license from the Caselaw Access Project to download bulk data for the Tax Court and other courts. The terms of the license prohibit sharing bulk data with other researchers, so this is the only dataset used for this Article that I cannot make available upon request.

## 3. *Excluding Non-Substantive Opinions*

Past work has generally measured the percentage of judicial opinions containing a particular interpretive tool (say, dictionaries or legislative history) out of the opinions in which some statutory interpretation occurs. The goal is to exclude opinions that are largely procedural in order to smooth variations in docket composition year over year. To achieve this, these papers have identified a “denominator” of interpretive opinions that divide the number of opinions containing hits for a particular tool.<sup>201</sup>

The Internal Revenue Bulletin contains a few texts in which novel statutory interpretation does not occur—particularly the IRB’s reproduction of the past year’s legislation and legislative history. I removed these from the analysis, which is mathematically equivalent

---

<sup>198</sup> Cf. Hickman, *supra* note 196, at 1730–31 (discussing how most, or all, Treasury regulations ought to be considered “legislative”).

<sup>199</sup> CASELAW ACCESS PROJECT, <https://case.law> (last visited Dec. 4, 2019). Ravel Law was subsequently acquired by LexisNexis. Thanks to Mike Lissner, Executive Director of the Free Law Project, for advice on obtaining these data and for providing the court data for early analyses of Tax Court and Supreme Court decisions. See *Bulk Data*, COURTLISTENER, <https://www.courtlistener.com/api/bulk-info> (last visited Dec. 4, 2019).

<sup>200</sup> Jason Tashea, *Caselaw Access Project Gives Free Access to 360 Years of American Court Cases*, A.B.A. J. (Oct. 30, 2018, 7:10 AM), [http://www.abajournal.com/news/article/caselaw\\_access\\_project\\_gives\\_free\\_access\\_to\\_360\\_years\\_of\\_american\\_court\\_cas](http://www.abajournal.com/news/article/caselaw_access_project_gives_free_access_to_360_years_of_american_court_cas).

<sup>201</sup> See, e.g., Bruhl, *supra* note 56, at 32–33; Calhoun, *supra* note 59, at 495–96.



to the denominator approach used in other articles.<sup>202</sup> The issue of procedural opinions does not arise for the Tax Court, since the dataset for this Article includes only Tax Court “division opinions,” which address novel legal issues.<sup>203</sup> (Tax Court opinions intended only to speak to settled law are called “memorandum opinions” or “oral opinions,” which are unpublished and theoretically lack precedential weight).<sup>204</sup>

### B. Terms Analyzed

The terms used in this Article were drawn from prior empirical work<sup>205</sup> as well as my own reading of relevant sources. All terms are listed in lower case, since the searches I conducted were not case sensitive. All terms were treated as stems for purposes of the counts, meaning that terms with different prefixes or suffixes would also be included. For example, “senate report” below includes “senate reports” as well.

Synonyms for the same concept (for example, “implied repeal” and “implicit repeal”) are all listed for completeness. In order to prevent the machine learning algorithm from overestimating predictive performance based on mere stylistic variation, I group together different terms within a particular category for purposes of the machine learning analysis. For example, the number of citations to Senate reports are aggregated, regardless of whether they are written as “S. Rep.”, “S. Rpt.”, or “Senate report.” Without these groupings, the

---

<sup>202</sup> To illustrate, say that a sample of documents has 150 documents overall, 50 that cite dictionaries and 100 that engage in any statutory interpretation. The denominator approach divides the 50 citing dictionaries by the 100 (the denominator) citing statutory interpretation, yielding 50%. My approach, which is computationally simpler, divides 50 by the 150 minus 50, also yielding 50%.

<sup>203</sup> See I.R.C. §§ 7459–60 (2018) (describing the process to issue division opinions); HAROLD DUBROFF & BRANT J. HELLWIG, *THE UNITED STATES TAX COURT: AN HISTORICAL ANALYSIS* 750–54 (2d ed. 2014) (discussing the historical and modern differences between division opinions and memorandum opinions); Grewal, *supra* note 75, at 2073–79.

<sup>204</sup> In practice, memorandum opinions are often cited and relied upon, and they do sometimes contain original legal decisionmaking. See DUBROFF & HELLWIG, *supra* note 203, at 753; Grewal, *supra* note 75, at 2073–79. However, the point remains that Tax Court division opinions are all intended to contain novel legal interpretation, and it is reassuring that the category is if anything underinclusive.

<sup>205</sup> See Bruhl, *supra* note 56, at 30–31, 38–39, 41, 53 (listing and describing the use of search terms to assess judicial purposivism, textualism, and canon use); Staudt et al., *supra* note 57, at 1932–35, 1940–42, 1949–51, 1955–59 (listing terms associated with textualism, purposivism, judicial deference, and canons of construction); Calhoun, *supra* note 59, at 524–25 app. I (listing dictionaries). I thank Aaron Bruhl for sharing the search terms that he used in his comparative study of judicial statutory interpretation.

algorithm might demonstrate a perfect ability to distinguish one court from another merely based on differences in citation practices.

One hazard attending machine learning is that conducting classification on an entire corpus of text—considering every word in a series of documents and testing whether each word has any predictive value—can produce seemingly strong predictive relationships purely by chance. This practice, known as “data dredging,” is a perennial risk when machine learning is used for social science research.<sup>206</sup> To avoid it, I constrain the vocabulary of words that the classifier may consider in the learning process to the interpretive terms set out in this Section. Importantly, the interpretive vocabulary was selected based on my *ex ante* views on interpretive methodology and draws heavily on the existing vocabularies selected by other authors,<sup>207</sup> rather than being selected *ex post* based on which terms had predictive value after running a machine learning algorithm. In doing so, I reduce the risk that the classifier may appear to successfully predict a result merely by chance or as a result of researcher choices.

## 1. *Purposivist Terms*

### Congressional Reports

<i>conference report</i>	<i>h.r. rept.</i>
<i>conf. rep.</i>	<i>h. r. rept.</i>
<i>conf. rpt.</i>	<i>h.r.rep.</i>
<i>conf. rept.</i>	<i>h.r.rpt.</i>
<i>conf.rep.</i>	<i>h.r.rept.</i>
<i>conf.rpt.</i>	<i>senate report</i>
<i>conf.rept.</i>	<i>s. rep.</i>
<i>house report</i>	<i>s. rpt.</i>
<i>h. rep.</i>	<i>s. rept.</i>
<i>h. rpt.</i>	<i>s.rep.</i>
<i>h. rept.</i>	<i>s.rpt.</i>
<i>h.rep.</i>	<i>s.rept.</i>
<i>h.rpt.</i>	<i>committee report</i>
<i>h.rept.</i>	<i>comm. rep.</i>
<i>h.r. rep.</i>	<i>comm. rpt.</i>
<i>h. r. rep.</i>	<i>comm. rept.</i>
<i>h.r. rpt.</i>	<i>comm.rep.</i>
<i>h. r. rpt.</i>	<i>comm.rpt.</i>
	<i>comm.rept.</i>

<sup>206</sup> See Gregg R. Murray & Anthony Scime, *Data Mining*, in *EMERGING TRENDS IN THE SOCIAL AND BEHAVIORAL SCIENCES* 1, 3–4 (Robert Scott & Stephen Kosslyn eds., 2015).

<sup>207</sup> In contrast, the normative terms were selected specifically for this Article.

## Congressional Hearings

*congressional hearing*  
*congressional record*  
*cong. rec.*  
*cong.rec.*  
*rec. doc.*

*committee hearing*  
*senate hearing*  
*house hearing*  
*conference hearing*

## Miscellaneous Legislative History

*legislative history*  
*history of the legislation*  
*conference committee*  
*joint committee*  
*jct*  
*congressional budget office*  
*cbo*

*senate committee*  
*s. comm.*  
*s. subcomm.*  
*house committee*  
*h.r. comm.*  
*h. subcomm.*  
*h. r. subcomm.*

## 2. Textualist Terms

Some potential synonyms for “plain meaning” were excluded, on the basis that courts have not always used them in a textualist manner. For example, the “literal meaning” of a statute<sup>208</sup> is often cited as a criticism of textualism rather than an endorsement of it.<sup>209</sup> Accordingly, I excluded that term in order to avoid false positives.

## Dictionaries<sup>210</sup>

*dictionary*<sup>211</sup>  
*dictionarium*  
*linguae britannicae*

*world book*  
*funk & wagnalls*

<sup>208</sup> Cf. Stephen C. Mouritsen, *The Dictionary Is Not a Fortress: Definitional Fallacies and a Corpus-Based Approach to Plain Meaning*, 2010 BYU L. REV. 1915, 1973 fig.5 (analyzing the terms “plain meaning,” “ordinary meaning,” “natural meaning,” “literal meaning,” and “common meaning”).

<sup>209</sup> See, e.g., *U.S. Padding Corp. v. Comm’r*, 88 T.C. 177, 184 (1987) (“We may then look to the reason of the enactment and inquire into its antecedent history and give it effect in accordance with its decision and purpose, sacrificing, if necessary, the literal meaning in order that the purpose may not fail.” (quoting *Ozawa v. United States*, 260 U.S. 178, 194 (1922))).

<sup>210</sup> These terms were borrowed in part from John Calhoun’s listing. See Calhoun, *supra* note 59, at 524–25 app. I.

<sup>211</sup> Occurrences of the word “dictionary” in “dictionary act” are excluded.

## Linguistic Canons<sup>212</sup>

*expressio*<sup>213</sup>

*expresio*

*inclusio*<sup>216</sup>

*noscitur a sociis*<sup>217</sup>

*ejusdem generis*<sup>214</sup>

*last antecedent*<sup>215</sup>

*plain meaning*

## Holistic-Textual Canons

*whole act*

*whole-act*

*whole code*

*whole-code*

*in pari materia*<sup>219</sup>

*meaningful variation*

*consistent usage*

*surplusage*<sup>218</sup>

*superfluity*

*superfluities*

## 3. Statutory Terms

Unlike the other terms in this Article, a document's statutory score was determined based on the number of statutory *sentences*. A sentence was designated as statutory if it included at least one word from the column on the left below and one word from the column on the right below.

<sup>212</sup> See Bruhl, *supra* note 56, at 56 ("The category of linguistic canons is composed of four familiar rules of word association and grammar: *ejusdem generis*, *noscitur a sociis*, *expressio unius*, and the rule of the last antecedent. All of these linguistic canons can be captured with good accuracy through electronic searches.").

<sup>213</sup> This phrase and its variants refer to the Latin maxim that *expressio unius est exclusio alterius*, meaning that express listing of certain items in a statute is presumed to exclude any unmentioned comparable items. *Chevron U.S.A. Inc. v. Echazabal*, 536 U.S. 73, 80 (2002) ("[E]xpressing one item of [an] associated group or series excludes another left unmentioned." (quoting *United States v. Vonn*, 535 U.S. 55, 65 (2002))).

<sup>214</sup> See *supra* notes 157–58 and accompanying text.

<sup>215</sup> See Jacob Scott, *Codified Canons and the Common Law of Interpretation*, 98 GEO. L.J. 341, 358 (2010) ("The last antecedent rule is somewhat confusing and hypergrammatician; it limits the operation of qualifying phrases to the last phrase in a sentence (rather than applying that limitation to the entire sentence).").

<sup>216</sup> "*Inclusio unius*" is a relatively rare variant whose effect is identical to the *expressio unius* canon. See *LawProse Lesson #227: Part 2: "Including but Not Limited to,"* LAWPROSE: BLOG, [www.lawprose.org/lawprose-lesson-227-part-2-including-but-not-limited-to](http://www.lawprose.org/lawprose-lesson-227-part-2-including-but-not-limited-to) (last visited Dec. 4, 2019) ("In legal literature, *expressio unius* is more than 15 times as common as *inclusio unius*.").

<sup>217</sup> See Staudt et al., *supra* note 57, at 1933 ("[T]he meaning of one term is 'known by its associates' (i.e., understood in the context of other words in the list).").

<sup>218</sup> See, e.g., *Corley v. United States*, 556 U.S. 303, 314 (2009) ("[O]ne of the most basic interpretive canons . . . [is] that '[a] statute should be construed so that effect is given to all its provisions, so that no part will be inoperative or superfluous, void or insignificant.'" (quoting *Hibbs v. Winn*, 542 U.S. 88, 101 (2004))).

<sup>219</sup> See *supra* notes 154–56 and accompanying text.

Includes: AND Includes:

construe  
construing  
construction  
interpret  
reading

statute  
statutory  
legislation  
congress  
code  
section

In addition, the following terms were included in the vocabulary used for machine learning analysis.

plain language  
legislative intent  
statutory purpose  
vagueness  
vague

*ambiguity*  
*ambiguities*  
*ambiguous*  
*unambiguous*

#### 4. Normative Terms

As noted above,<sup>220</sup> the phrase “effective tax administration” is excluded from counts using the following terms, as are the phrases “treasury inspector general for tax administration” and “small business regulatory enforcement fairness act.” In addition, any occurrences of normative terms in sentences that also contained purposivist terms, textualist terms, or substantive canons were excluded, in order to avoid policy judgments that occur in the interpretive process (for example, legislative history that discusses fairness).

good public policy  
public policy goal  
public policy grounds  
tax administration  
efficient administration  
efficient tax collection  
efficient enforcement  
compliance burden  
financial burden  
administrative burden

regulatory burden  
burdensome  
compliance cost  
complexity  
intrusive  
fairness  
unfair  
injustice  
unjust  
clarity

## 5. Substantive Canons

Deference regimes, such as those under *Chevron* and *Skidmore*, have sometimes been considered precedents and sometimes consid-

<sup>220</sup> See *supra* note 93 and accompanying text.

ered canons of construction.<sup>221</sup> I classify them as substantive canons for purposes of this Article but do not otherwise take a position on which categorization is more accurate.

### General Substantive Canons

<i>charming betsy</i>	<i>repeal by implication</i>
<i>rule of lenity</i>	<i>implied repeal</i>
<i>absurd result</i>	<i>implicit repeal</i>
<i>avoidance canon</i>	<i>implicitly repeal</i>
<i>canon of avoidance</i>	<i>presumption against preemption</i>
<i>constitutional avoidance</i>	<i>presumption against pre-emption</i>

### Deference Canons

<i>chevron</i>	<i>seminole rock</i>
<i>skidmore</i>	<i>auer</i>

### C. *Non-Normal Distribution of Term Frequencies in Tax Court Opinions*

Term frequencies in Tax Court opinions have several important distributional features that demand special attention in statistical analysis (including machine learning analysis). First, they are “semicontinuous”<sup>222</sup>: They vary continuously (they are not limited to whole numbers) but cannot be less than zero (since no opinion can use any term less than zero times). Second, they are “zero-inflated”<sup>223</sup>: Many courts use *no* terms of any particular type, such that the median number of purposivist, textualist, statutory, and normative terms used in Tax Court opinions is zero in each case. Third, they are “log-normal”: Even excluding zero values, the distributions exponentially decrease, with long right tails (i.e., most cases use few terms, but some cases use a large number of terms), requiring log-transformation to turn them into normal distributions.

Each of these features violates the conventional assumption of normal distribution that underlies conventional statistical analysis, including standard ordinary least squares (OLS) regression and machine learning on raw term frequencies. Log-normality also casts doubt on visual analysis of the term frequency charts in this Article.

<sup>221</sup> See *supra* note 159.

<sup>222</sup> See Yongyi Min & Alan Agresti, *Modeling Nonnegative Data with Clumping at Zero: A Survey*, 1 J. IRANIAN STAT. SOC’Y 7, 7–8 (2002) (“We refer to a variable as *semicontinuous* when it has a continuous distribution except for a probability mass at 0.”).

<sup>223</sup> See *id.* at 7 (2002) (“Applications in which data take nonnegative values but have a substantial proportion of values at zero occur in many disciplines. The modeling of such ‘clumped-at-zero’ or ‘zero-inflated’ data is challenging.”).

There is a risk that any analysis of data following a log-normal distribution will be driven by outliers and therefore be less robust. Consequently, this dataset requires additional transformation to confirm the robustness of the results in this Article and should not be interpreted using OLS regression or raw term frequencies alone.

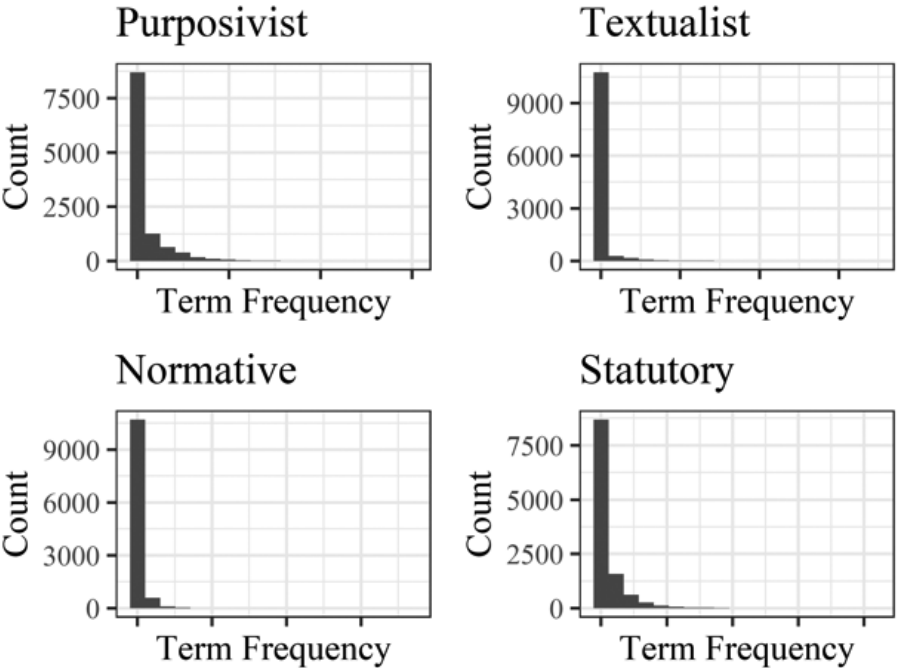
Table 5 illustrates the problem of zero-inflation in the data:

TABLE 5. PERCENTAGE OF TAX COURT OPINIONS WITH ZERO TERMS, 1942–2015

<i>Type of Term</i>	
Purposivist	69.89%
Textualist	93.46%
Statutory	70.31%
Normative	88.27%

Figure 14 illustrates all three issues: semicontinuity, zero-inflation, and the log-normal distribution:

FIGURE 14. HISTOGRAM OF PURPOSIVIST, TEXTUALIST, NORMATIVE, AND STATUTORY TERMS IN TAX COURT OPINIONS, 1942–2015

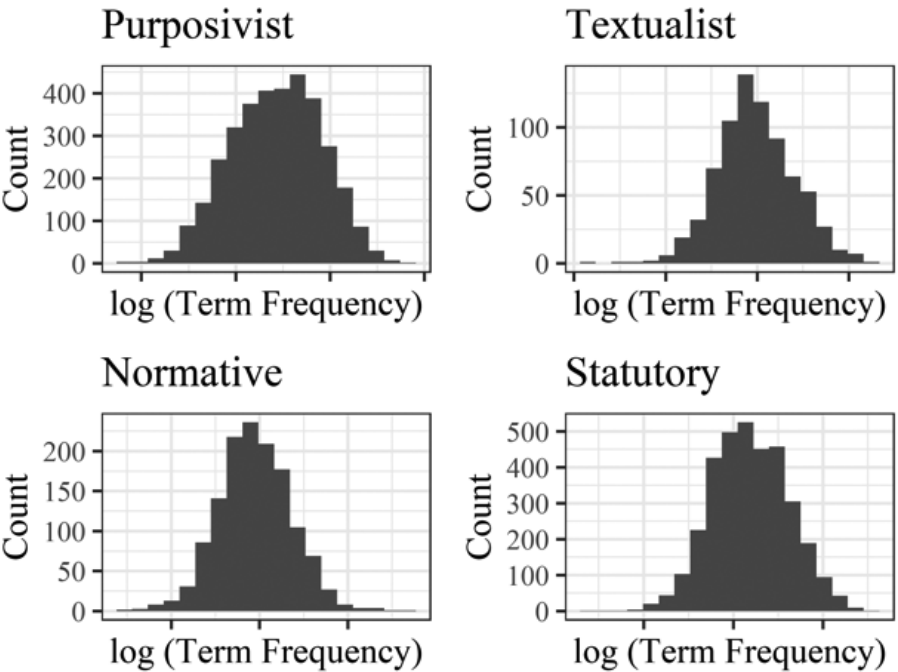


Fortunately, a log-normal distribution can be easily addressed by logarithmically transforming the data. One method is to log-transform the data as follows:

$$\tilde{y} = \log(y) \tag{1}$$

When the data are log-transformed in this way, they take the shape of the normal distribution (when excluding zeros—zero-inflation is a separate problem that I address in Section E.2 of the Appendix). Section F reproduces each term frequency chart in this Article after log-transformation. Figure 15 illustrates that the log-transformation produces approximately normal distributions.

FIGURE 15. LOG-TRANSFORMED HISTOGRAM OF PURPOSIVIST, TEXTUALIST, NORMATIVE, AND STATUTORY TERMS IN TAX COURT OPINIONS, 1942–2015



From the log-transformed histograms, it is evident that the distribution of data points is approximately log-normal when considering opinions with more than zero terms. This confirms that the data can



be described as semicontinuous, with zero-inflation and a log-normal distribution.<sup>224</sup>

This Article employs three methodologies, each of which must appropriately account for these distributional features. To ensure that the charts presented above are valid, Section F of the Appendix presents log-transformed versions of each of them. To ensure that the machine learning methodology is valid, Section D of the Appendix describes how the transformer used in the machine learning analysis normalizes the data prior to the operation of the classifier. Finally, to ensure that regression analysis is valid, Section E.2 of the Appendix employs a two-part regression model specifically designed to address semicontinuity, zero-inflation, and log-normality, which are common issues in natural datasets.

#### *D. Tf-idf Transformation and Classification in Machine Learning*

This Section provides additional detail on the methodology used for the machine learning analysis in this Article, especially in light of the log-normal distribution of term frequencies discussed in the previous Section. Section II.B discussed how Tax Court opinions are first vectorized by obtaining term frequencies for each term of interest, and ultimately classified by an algorithm (in this case, a logistic regression) that improves by iterating over a training set.

Between vectorizing and classification, however, the term frequencies are also *transformed* in order to make the classification statistically valid. The transformation converts raw term frequency to term frequency-inverse document frequency (tf-idf) and normalizes the data in the process. Mathematically, given term frequency  $tf_{t,d}$  with respect to term  $t$  and document  $d$ , term frequency is log-transformed so that:

$$\widetilde{tf}_{t,d} = \log(1 + tf_{t,d}) \quad (2)$$

Notice that this log-transformation is the same one used in Section F of the Appendix. Next, inverse document frequency is calcu-

---

<sup>224</sup> I was not able to separate the dataset of IRS publications cleanly into discrete individual publications (which in any case are much more heterogeneous than court opinions; many IRS publications are merely administrative and only a few lines long). Consequently, I could not conduct the histogram analysis above for IRS publications. However, since it is plausible that IRS publications would also follow the same problematic distribution—anecdotally, I noticed outliers while cleaning the dataset, where the IRS heavily utilized certain interpretive tools in explaining particularly knotty guidance—out of caution, I apply the same logarithmic corrections in Section F of the Appendix for IRS publications as I do for Tax Court opinions.

lated as a function of  $N$ , the number of documents in the corpus, and  $df_t$ , the number of documents in the corpus for which  $tf_{t,d} > 0$ :

$$idf_t = \log(N / df_t) \quad (3)$$

Finally, tf-idf is calculated as a function of log-transformed term frequency and inverse document frequency:

$$tfidf_{t,d} = \widetilde{tf}_{t,d} \cdot idf_t \quad (4)$$

Conceptually, the use of tf-idf rather than raw term frequency prevents certain terms from having an outsize influence on the regression merely because they are rarer. The inclusion of log-transformation in the tf-idf transformation addresses the log-normality of the term frequency distribution.

Because the tf-idf statistic is then used in a classifier modeled as a logistic regression, the use of tf-idf rather than raw term frequency merely multiplies each coefficient in the regression by a scalar and therefore does not affect statistics such as MCC, Accuracy, or  $F_1$  score, nor does it affect the statistical significance of each term. This can be seen by considering the regression that the classifier conducts, where  $p / (1 - p)$  is the odds ratio with respect to the classification category (e.g., a Tax Court opinion), and  $n$  is the number of terms.

$$\log\left(\frac{p}{1-p}\right) = \beta_0 + \beta_1 \cdot tfidf_{1,d} + \beta_2 \cdot tfidf_{2,d} + \cdots + \beta_n \cdot tfidf_{n,d} + \epsilon_d \quad (5)$$

Through Equations 4 and 5, we find that, for any term  $t$ :

$$\beta_t \cdot tfidf_{t,d} = \beta_t \cdot \widetilde{tf}_{t,d} \cdot idf_t \quad (6)$$

$idf_t$  varies with respect to each term and not with respect to each document. This means that it is a scalar multiplier against coefficient  $\beta_t$ . In other words, the relationship between  $\beta_t$  in this regression and  $\widehat{\beta}_t$  from a different regression conducted only on log-transformed term frequency is that:

$$\widehat{\beta}_t = \beta_t \cdot idf_t \quad (7)$$

### *E. Regression Analysis of Tax Court Opinions*

This Section employs regressions to more closely analyze the relationship between interpretive methodology in Tax Court opinions,

on the one hand, and either case outcomes or party affiliation of judges, on the other hand.<sup>225</sup>

All of the regressions in this Section use clustered standard errors with clustering by judge, a variant of robust standard errors that accounts for heteroskedasticity across the “clusters” of opinions written by different judges. The regressions take each Tax Court opinion as a single observation, using term frequency (either purposivist or textualist) as the dependent variable, and taking as the independent variables: (1) party affiliation of the judge authoring the opinion, (2) case outcome, (3) the year that the judge writing the opinion was appointed, (4) fixed effects for the year the opinion was written, and/or (5) fixed effects for the judge that wrote the opinion. In regressions where judge fixed effects are used, party affiliation and the judge’s year of appointment are dropped as multicollinear with the fixed effects. Fixed effects introduce dummy variables for each year or judge, which control for variation in methodology over time and between judges,<sup>226</sup> isolating differences within a particular year and within a particular judge’s docket.

Table 6 presents summary statistics for Tax Court opinions to facilitate interpretation of the regression results in this Section.

TABLE 6. SUMMARY STATISTICS FOR TAX COURT OPINIONS, 1942–2015

	<i>N</i>	<i>Minimum</i>	<i>Mean</i>	<i>Median</i>	<i>Maximum</i>	<i>Standard Deviation</i>
Democrat <sup>227</sup>	7308	0	0.561	1	1	0.496
Taxpayer Wins <sup>228</sup>	4261	0	0.224	0	1	0.417
Textualist Term Frequency <sup>229</sup>	11,451	0	30.0	0	3427.6	162.5
Purposivist Term Frequency <sup>230</sup>	11,451	0	365.3	0	11,869.4	967.3

225

*See supra* Sections III.F, III.G.

226

*See generally* PAUL D. ALLISON, FIXED EFFECTS REGRESSION MODELS (2009) (describing fixed effects regression models).

227

This variable equals 1 if the judge authoring an opinion is a Democrat and 0 otherwise. Thus, this row indicates that 56.1% of opinions were authored by Democrats.

228

This variable equals 1 if the taxpayer won and 0 otherwise. Thus, this row indicates that the taxpayer won in 22.4% of cases.

229

Terms per million words.

230

Terms per million words.

It should be noted that regressions measure fundamentally different things than classifier accuracy results. Here, classifier accuracy measures how well opinions can be categorized into one of two categories based on interpretive methodology alone. This is roughly analogous to a regression with a binary category dummy (say, Democratic or Republican) as the dependent variable and each specific interpretive term (say, “dictionary”) as the independent variables. (This description glosses over some additional nuance, of course, such as the transformation discussed in the previous Section and the fact that only some classifier techniques are analogous to regression.)<sup>231</sup> Classifier accuracy therefore measures the extent to which methodology alone can explain the variation between the two categories.

In contrast, the regression analysis in this Section more narrowly asks whether specific variables have a statistically significant relationship with methodology. The regressions do not analyze a vector consisting of many different interpretive tools, but rather a single summary statistic reflecting the term frequency of all textualist or purposivist tools, respectively, in each opinion. Most importantly, the experimental hypothesis is completely different. Tests of statistical significance in regression analysis ask only whether a particular variable has *any* effect (i.e., whether we can reject the null hypothesis that the variable has no effect); classifier accuracy gauges the *magnitude* of the effect by asking how much that variable drives outcomes. Classifier accuracy is therefore a cousin of  $R^2$ , the measure of the portion of variation in the dependent variable that can be explained by all of the independent variables. A result might be very statistically significant but still have a low  $R^2$ .

### 1. Ordinary Least Squares Regression Model

Because term frequencies are not normally distributed, as described in Section C of the Appendix, OLS is not an appropriate regression model for these data. Nevertheless, I present OLS results for comparison with the results of the two-part regression model.<sup>232</sup> For document  $d$ , number of years  $y$ , and number of judges  $j$ , the models for each regression in Tables 7 and 8 are (in order, left to right):

$$tf_d = \beta_0 + \beta_1 \cdot Democrat_d + \epsilon_d \quad (8)$$

<sup>231</sup> See *supra* notes 98–102 and accompanying text.

<sup>232</sup> I used OLS regression in Stata, version 16, using robust variance estimates. *Robust Variance Estimates*, STATA, [https://www.stata.com/manuals13/p\\_robust.pdf](https://www.stata.com/manuals13/p_robust.pdf) (last visited Dec. 4, 2019).

$$tf_d = \beta_0 + \beta_1 \cdot Year\ Judge\ Appointed_d + \epsilon_d$$

(9)

$$tf_d = \beta_0 + \beta_1 \cdot Taxpayer\ Wins_d + \epsilon_d$$

(10)

$$tf_d = \beta_0 + \beta_1 \cdot Democrat_d + \beta_2 \cdot Year\ Judge\ Appointed_d + \beta_3 \cdot Taxpayer\ Wins_d + \sum_{i=1}^y \beta_{4,i} \cdot Year_{d,i} + \epsilon_d$$

(11)

$$tf_d = \beta_0 + \beta_1 \cdot Taxpayer\ Wins_d + \sum_{i=1}^y \beta_{2,i} \cdot Year_{d,i} + \sum_{k=1}^j \beta_{3,k} \cdot Judge_{d,k} + \epsilon_d$$

(12)

TABLE 7. OLS REGRESSION RESULTS FOR TAX COURT  
PURPOSIVISM

<i>Dependent variable: purposivist terms (per million words)</i>					
Democrat	-44.6 (80.1)		159.7** (64.6)		
Year Judge Appointed		12.42*** (1.06)	2.8 (3.6)		
Taxpayer Wins			81.9* (46.26)	-13.3 (50.0)	6.6 (39.8)
Opinion Year Fixed Effects	No	No	No	Yes	Yes
Judge Fixed Effects	No	No	No	No	Yes
<i>R</i> <sup>2</sup>	0.0006	0.0535	0.0012	0.1348	0.1540
<i>N</i>	7308	11,451	4261	2763	4255
Note: Each column of this table represents the results of a separate regression. The dependent variable in each regression is the frequency of textualist terms, in words per million. The fixed effects rows indicate whether dummy variables are included for each opinion year, each judge authoring opinions, or both. When judge fixed effects are included, judge characteristics (party and year of appointment) are omitted as multicollinear. <i>N</i> varies between regressions because some observations lacked determinate values for some variables (for example, some cases lack a clear winner, since the taxpayer won on some issues and lost on others). Standard errors are clustered by judge. * denotes statistical significance at p<0.1, ** at p<0.05, and *** at p<0.01.					

TABLE 8. OLS REGRESSION RESULTS FOR TAX COURT  
TEXTUALISM

<i>Dependent variable: textualist terms (per million words)</i>					
Democrat	-12.6 (8.4)			-7.3 (6.8)	
Year Judge Appointed		1.15*** (0.19)		-0.20 (0.53)	
Taxpayer Wins			1.2 (5.9)	-3.4 (6.1)	-3.2 (4.7)
Opinion Year Fixed Effects	No	No	No	Yes	Yes
Judge Fixed Effects	No	No	No	No	Yes
$R^2$	0.0013	0.0130	0.0000	0.0623	0.0994
$N$	7308	11,451	4261	2763	4255
Note: Each column of this table represents the results of a separate regression. The dependent variable in each regression is the frequency of textualist terms, in words per million. The fixed effects rows indicate whether dummy variables are included for each opinion year, each judge authoring opinions, or both. When judge fixed effects are included, judge characteristics (party and year of appointment) are omitted as multicollinear. $N$ varies between regressions because some observations lacked determinate values for some variables (for example, some cases lack a clear winner, since the taxpayer won on some issues and lost on others). Standard errors are clustered by judge. * denotes statistical significance at $p<0.1$ , ** at $p<0.05$ , and *** at $p<0.01$ .					

2. Two-Part Regression Model

Although OLS regression may be useful to set a baseline, it is a poor fit for the Tax Court data analyzed in this Article. As described in Section C of the Appendix, any regression method must specially adjust for the fact that term frequencies in this dataset are semicontinuous,<sup>233</sup> zero-inflated,<sup>234</sup> and log-normal. Each of these features violates the assumption of normal distribution that underlies OLS regression.

However, these features frequently appear in natural datasets, and econometricians have developed alternative regression methods to address them.<sup>235</sup> In this Section, I will use the two-part regression

<sup>233</sup> See Min & Agresti, *supra* note 222, at 7–9.  
<sup>234</sup> See *id.*  
<sup>235</sup> See generally J.A. Cole & J.D.F. Sherriff, *Some Single- and Multi-Site Models of Rainfall Within Discrete Time Increments*, 17 J. HYDROLOGY 97 (1972) (applying an early

model first developed by Naihua Duan et al.<sup>236</sup> and implemented by Federico Belotti et al.<sup>237</sup> Conceptually, the model is divided between a first part to determine whether the dependent variable has a zero or positive value, and a second part to determine the positive value, conditional on the value being positive. This models, for example, a situation in which a judge makes an initial decision on whether to use any textualist terms and, if she does so, a second decision on how many textualist terms to use.

Mathematically (and assuming a single independent variable for simplicity), and for our purposes applying a logistic regression, the first step may be represented as<sup>238</sup>:

$$\text{logit}[P(Y_i = 0)] = x'_{1i} \cdot \beta_1 + \epsilon_i \quad (13)$$

The second step is a regression of the value of  $y_i$  conditional on  $y_i$  being positive, for our purposes assuming a log-normal distribution<sup>239</sup>:

$$\log[y_i | y_i > 0] = x'_{2i} \cdot \beta_2 + \epsilon_i \quad (14)$$

The model separately estimates the marginal effect of each independent variable with respect both to the first part and the second part. But the two parts can also be combined to estimate the overall marginal effect of each independent variable with respect to the dependent variable. That is, the combined marginal effect of  $x_i$  both in changing the likelihood that  $y_i$  will be positive, as well as the marginal predictive effect of  $x_i$  on  $y_i$  in case  $y_i$  is positive. Mathematically, this is represented as<sup>240</sup>:

$$y_i = \hat{y}_i | x_i = (\hat{p}_i | x_i) \cdot (\hat{y}_i | y_i > 0, x_i) \quad (15)$$

Equations 8, 11, and 12 are modified in order to reflect Equations 13 through 15. Results from the two-part regression are presented in Tables 9 and 10. Each table contains three regressions, and each regression is in turn separated between the first part, second part, and combined marginal effect. Note that the coefficients in each of the three columns represent the results of very different regressions and are not directly comparable except in sign.

---

version of a two-part regression model to estimate rainfall); Naihua Duan et al., *Choosing Between the Sample-Selection Model and the Multi-Part Model*, 2 J. BUS. & ECON. STAT. 283 (1984) (applying a two-part model to estimate healthcare expenditures).

<sup>236</sup> See Duan et al., *supra* note 235.

<sup>237</sup> See Federico Belotti et al., *Twopm: Two-Part Models*, 15 STATA J. 3 (2015).

<sup>238</sup> Min & Agresti, *supra* note 222, at 11. This particular example assumes that a logit model is used for the first part, which is the model I use in this Article. A probit model may also be used but would not have been appropriate for these data.

<sup>239</sup> *Id.*

<sup>240</sup> Belotti et al., *supra* note 237, at 7.

$N$  changes between the tables, even for regressions with the same dependent and independent variables, because the first part of the regression drops any observations if the zero-positive dichotomy can be perfectly predicted based on any independent variable, including a dummy variable—for example, if any judge never uses a textualist term, or if no textualist terms were used in any opinion for a given year.

The first-step regression, as noted above, is a logit model. The second-step regression is a generalized linear model (GLM), which is a generalization of the OLS model with some assumptions relaxed. Specifically, I use a GLM model with a log-link function and a Poisson distribution, in order to account for the distribution of term frequencies.<sup>241</sup> The coefficients from the first and second parts are retransformed in order to calculate combined marginal effects on a raw scale, because they are both calculated on non-linear scales. The regressions are presented with McFadden's  $R^2$  statistics for both steps. McFadden's is an alternative measure of goodness-of-fit that is appropriate to logit and log-linked GLM regressions.

To confirm that Poisson is the appropriate distribution family for the GLM model, I conduct a modified Park test. The modified Park test evaluates the relationship between the square of the residuals from the GLM regression (the variance) and the natural logarithm of the regression's predicted values (the mean). OLS regression assumes that there is no relationship (this is the assumption of homoskedasticity).<sup>242</sup> Applying the modified Park test to the second step of the two-part model (with full controls), I find a coefficient of 1.058 with respect to purposivist terms, and 0.968 with respect to textualist terms. Chi-squared tests fail to reject the null hypothesis that these coefficients are equal to 1, implying that a Poisson distribution is the appropriate distribution family in each case.

One interesting supplemental finding to those in Section III.G is that the first-part coefficients for case outcomes are positive, but the second-part coefficients are negative. This suggests that cases in which the taxpayer wins are more likely to include at least one purposivist or textualist term, but cases in which the taxpayer loses are more likely to include more than one (conditional on including at least one). This result is not statistically significant but possibly warrants additional research.

---

<sup>241</sup> For an example of this model in a two-part regression, see *id.* at 10–13.

<sup>242</sup> See Partha Deb & Edward C. Norton, *Modeling Health Care Expenditures and Use*, 39 ANNU. REV. PUB. HEALTH 489, 497 (2018) (“One should use the Gaussian distribution in the GLM when the coefficient on the expected value is close to 0.0 . . . . One should use a Poisson-type distribution . . . when the coefficient is close to 1.0 . . . .”).



TABLE 9. TWO-PART REGRESSION RESULTS FOR TAX COURT PURPOSIVISM

	<i>Dependent variable: purpoivist terms (per million words)</i>					
	Logit (1 <sup>st</sup> )	GLM (2 <sup>nd</sup> )	Combined	Logit (1 <sup>st</sup> )	GLM (2 <sup>nd</sup> )	Combined
Democrat	-0.408* (0.212)	0.161* (0.091)	-44.6 (65.8)	0.127 (0.108)	0.255*** (0.099)	146.3*** (56.9)
Year Judge Appointed				0.0125 (0.0088)	0.0009 (0.0039)	2.6 (2.8)
Taxpayer Wins				0.096 (0.113)	-0.05 (0.071)	-0.3 (42.3)
Opinion Year Fixed Effects	No	No	No	Yes	Yes	Yes
Judge Fixed Effects	No	No	No	No	No	No
McFadden's $R^2$	0.0072	0.0129		0.1326	0.0166	
$N$	7308	7308	7308	2760	2760	2760
				4241	4241	4241

Note: Each "Logit" and "GLM" column reflects the first and second part of a two-part regression, respectively, with the "Combined" column reflecting the marginal effect calculated by combining both columns. The fixed effects rows indicate whether dummy variables are included for each opinion year, each judge authoring opinions, or both. When judge fixed effects are included, judge characteristics (party and year of appointment) are omitted as multicollinear.  $N$  varies between regressions because some observations lacked determinate values for some variables (for example, some cases lack a clear winner, since the taxpayer won on some issues and lost on others). Standard errors are clustered by judge. \* denotes statistical significance at  $p < 0.1$ , \*\* at  $p < 0.05$ , and \*\*\* at  $p < 0.01$ .

TABLE 10. TWO-PART REGRESSION RESULTS FOR TAX COURT TEXTUALISM

Dependent variable: textualist terms (per million words)									
	Logit (1 <sup>st</sup> )	GLM (2 <sup>nd</sup> )	Combined	Logit (1 <sup>st</sup> )	GLM (2 <sup>nd</sup> )	Combined	Logit (1 <sup>st</sup> )	GLM (2 <sup>nd</sup> )	Combined
Democrat	-0.555** (0.224)	0.147 (0.115)	-12.4 (8.2)	-0.14 (0.16)	-0.26* (0.15)	-14.4* (7.6)			
Year Judge Appointed				0.0014 (0.0074)	-0.0036 (0.0063)	-0.09 (0.34)			
Taxpayer Wins				0.20 (0.18)	-0.31** (0.12)	-5.2 (7.5)	0.17 (0.15)	-0.33** (0.140)	-6.8 (6.7)
Opinion Year Fixed Effects	No	No	No	Yes	Yes	Yes	Yes	Yes	Yes
Judge Fixed Effects	No	No	No	No	No	No	Yes	Yes	Yes
McFadden's $R^2$	0.0010	0.0137		0.0988	0.0391		0.1501	0.0550	
$N$	7308	7308	7308	2479	2479	2479	4041	4041	4041
Note: Each "Logit" and "GLM" column reflects the first and second part of a two-part regression, respectively, with the "Combined" column reflecting the marginal effect calculated by combining both columns. The fixed effects rows indicate whether dummy variables are included for each opinion year, each judge authoring opinions, or both. When judge fixed effects are included, judge characteristics (party and year of appointment) are omitted as multicollinear. $N$ varies between regressions because some observations lacked determinate values for some variables (for example, some cases lack a clear winner, since the taxpayer won on some issues and lost on others). Standard errors are clustered by judge. * denotes statistical significance at $p<0.1$ , ** at $p<0.05$ , and *** at $p<0.01$ .									

F. Log-Transformed Charts

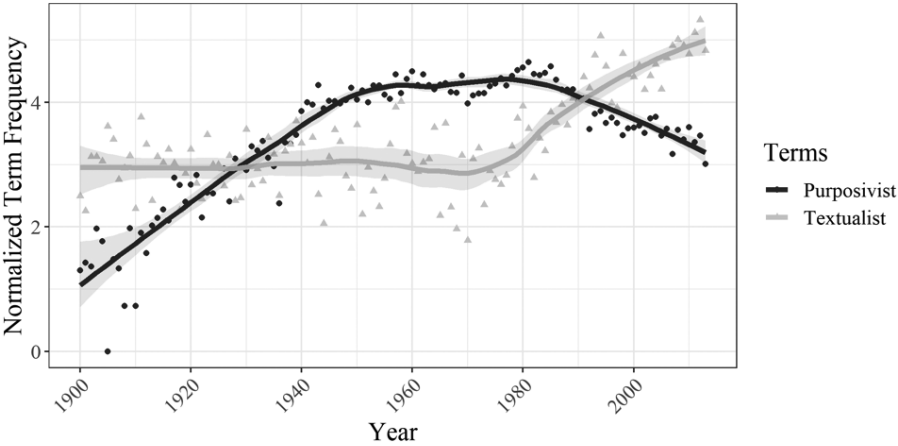
As noted in Section C of the Appendix, LOESS regression analysis of a long-right-tailed non-normal distribution may inadvertently place outsize importance on outliers. In order to visually ensure that the figures used in this Article are robust and not merely driven by outliers, this Section recreates each term frequency chart using the log-transformation specified in Equation 1<sup>243</sup>:

$$\hat{y} = \log(1 + y)$$

(16)

Visual examination of the log-transformed charts suggests approximately the same results as presented earlier in this Article.

FIGURE 16. PURPOSIVIST AND TEXTUALIST TERMS IN SUPREME COURT OPINIONS



<sup>243</sup> Note that in each case, the term frequency subjected to the log-transform is expressed in terms per million words. This makes the left scale of the graph more readable but does not affect the shape of the curve.

FIGURE 17. PURPOSIVIST AND TEXTUALIST TERMS IN DISTRICT COURT OPINIONS

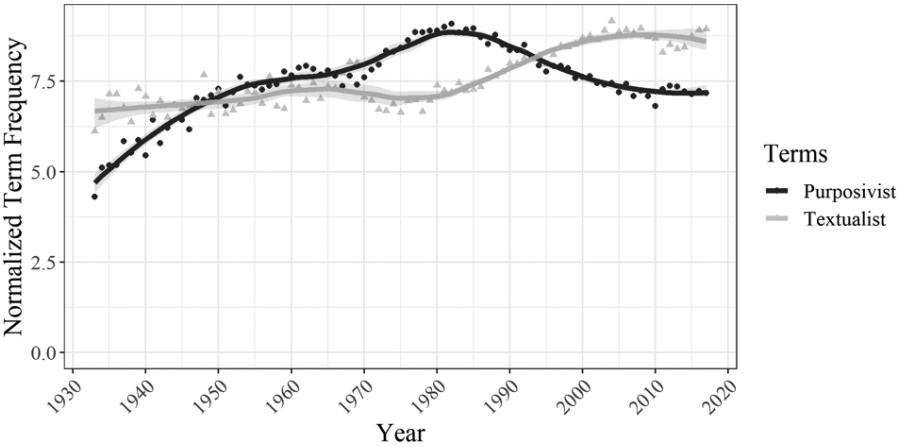
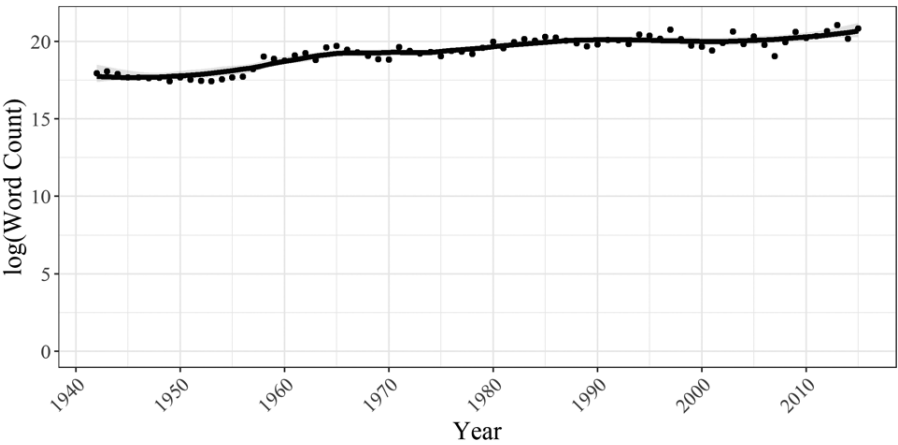


FIGURE 18. AVERAGE WORD COUNT OF TAX COURT OPINIONS



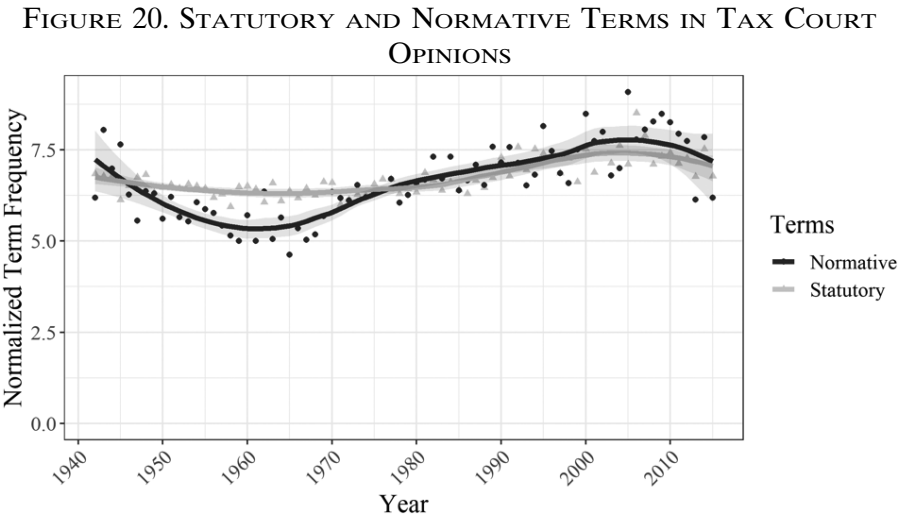
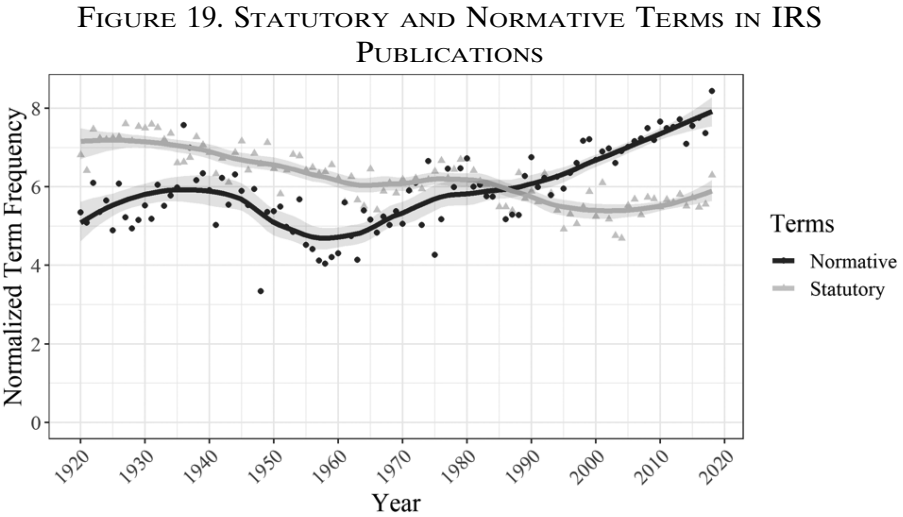


FIGURE 21. PURPOSIVIST AND TEXTUALIST TERMS IN IRS PUBLICATIONS

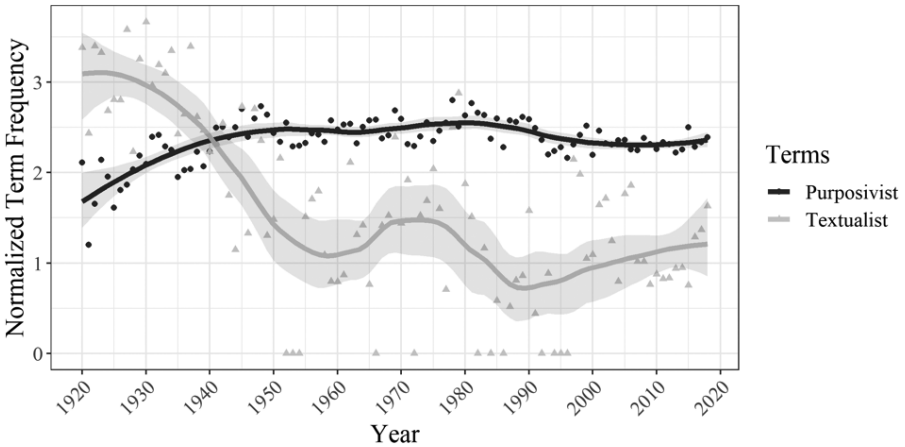
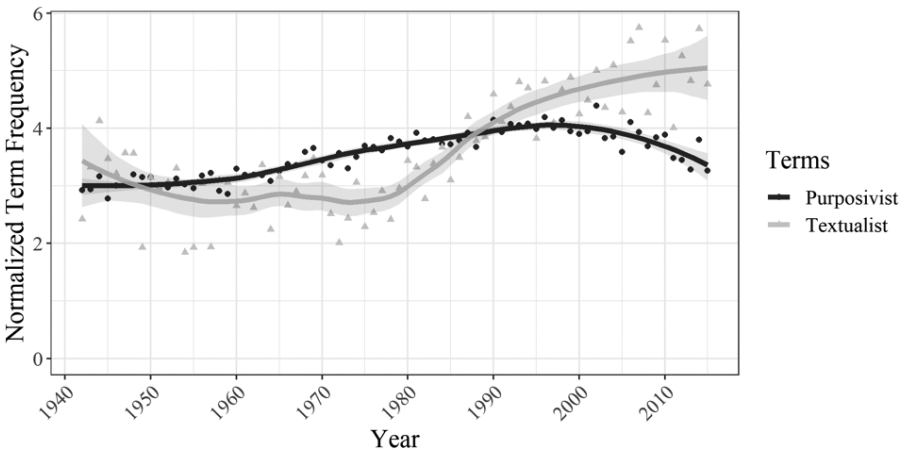


FIGURE 22. PURPOSIVIST AND TEXTUALIST TERMS IN TAX COURT OPINIONS



### G. Bootstrapped Confidence Intervals for Machine Learning

The bootstrapped confidence intervals in Section IV.B were calculated as basic bootstrap confidence intervals, the same form of bootstrapping used to calculate confidence intervals for the longitudinal figures in this Article.<sup>244</sup> These are sometimes known as “empirical confidence intervals” and avoid making certain assumptions about the functional form of standard errors. Consequently, they are better suited to bootstrapping than conventional confidence intervals. The Python code used to conduct the bootstrapping and to calculate the

<sup>244</sup> See *supra* note 64.

confidence intervals is available online.<sup>245</sup> I conducted bootstrapping with one thousand tests.

The histograms generated from the bootstrapping, Figures 23 and 24 below, suggest that each of the performance statistics used (MCC, Accuracy, and  $F_1$  score) was approximately normally distributed over the bootstrapping tests.

FIGURE 23. HISTOGRAM OF MCC,  $F_1$ , AND ACCURACY RESULTS FROM BOOTSTRAPPING, TAX COURT V. DISTRICT COURTS

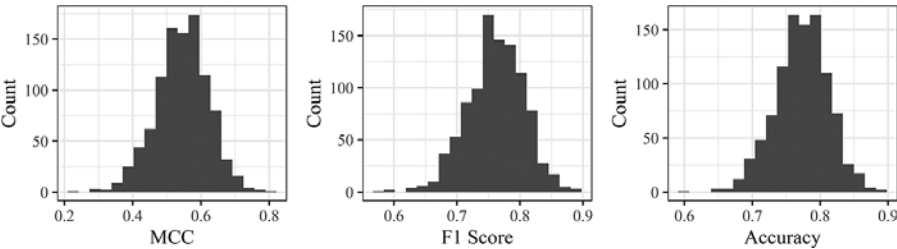
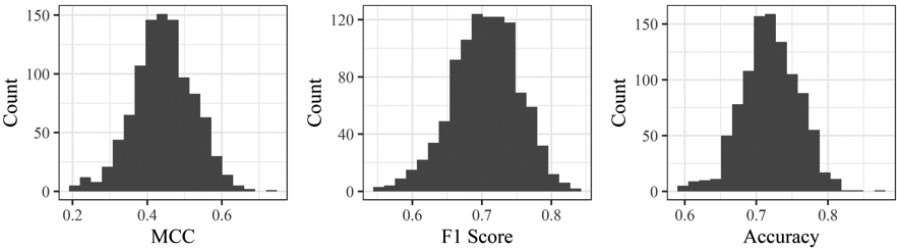


FIGURE 24. HISTOGRAM OF MCC,  $F_1$ , AND ACCURACY RESULTS FROM BOOTSTRAPPING, TAX COURT V. CFC



<sup>245</sup> See *Code*, *supra* note 91.