
Pile of Law: Learning Responsible Data Filtering from the Law and a 256GB Open-Source Legal Dataset

Peter Henderson*, Mark S. Krass*, Lucia Zheng, Neel Guha
Christopher D. Manning, Dan Jurafsky, Daniel E. Ho
Stanford University

Abstract

One concern with the rise of large language models lies with their potential for significant harm, particularly from pretraining on biased, obscene, copyrighted, and private information. Emerging ethical approaches have attempted to filter pretraining material, but such approaches have been ad hoc and failed to take context into account. We offer an approach to filtering grounded in law, which has directly addressed the tradeoffs in filtering material. First, we gather and make available the Pile of Law, a \sim 256GB (and growing) dataset of open-source English-language legal and administrative data, covering court opinions, contracts, administrative rules, and legislative records. Pretraining on the Pile of Law may help with legal tasks that have the promise to improve access to justice. Second, we distill the legal norms that governments have developed to constrain the inclusion of toxic or private content into actionable lessons for researchers and discuss how our dataset reflects these norms. Third, we show how the Pile of Law offers researchers the opportunity to learn such filtering rules directly from the data, providing an exciting new research direction in model-based processing.

Warning: this paper contains quotations that may be offensive or upsetting.

1 Introduction

The presence of private and toxic content in the most popular corpora for pretraining large language models is a well-known problem [11, 58]. But what to do about it is largely a matter of researcher discretion. Some teams implement extensive processes for filtering content deemed toxic or private; others train on data in virtually unmodified form. Resolving all of the difficulties and nuances of content filtering can be challenging, potentially explaining why content filtering has been so uneven.

It is practically difficult to perform reliable and transparent filtering at scale. That is partially because undesirable content is deeply contextual. For example, whether the inclusion of a racial epithet in a dataset is toxic may depend on factors such as the identity of the speaker and the expectations of the readers [43, 114]. Likewise, the existence of privacy violations may depend in part on the extent to which a speaker expected a fact to be widely shared at the time it was expressed [20, 6-7]. And privacy expectations may vary widely across countries [10].

Any filtering process involves complex trade-offs. Filtering for toxicity may have unexpected effects on representation in datasets or the bias of downstream outputs [40, 48, 6]. And filtering too widely for privacy may harm important downstream applications, as when the Census Bureau’s adoption of differential privacy led to errors in redistricting U.S. Congressional districts [102].

Yet researchers are not the first to balance the merits of open-source transparency with potential harms: legal and administrative actors have expended significant resources and process in developing

*Equal contribution.

standards to strike this exact balance. In this work, we suggest that researchers can look to these long-developed (and debated) standards to help ground content filtering mechanisms for large language model training.

This paper makes three contributions. First, we curate and open-source a ~ 256 GB (and growing) dataset of legal and administrative data, which we call Pile of Law, which can be used for assessing norms on data sanitization across legal and administrative settings. This dataset can be an exploratory tool for evaluating different mechanisms for “doing the data work” [104]. And we note that pretraining on the Pile of Law may help with challenging legal tasks that have the potential to improve access to justice [16]. Second, we catalog how government has, through extensive legislation, regulation, and litigation, developed standards for handling the trade-offs between privacy and offensive content on the one hand and transparency, access, and completeness on the other. We suggest actionable insights for researchers based on these legal and administrative norms.² Third, we demonstrate how implicit sanitization rules can be learned from the Pile of Law, providing a path forward for researchers to develop more nuanced filtering mechanisms. We also demonstrate shortcomings in alignment for current sanitization techniques, providing exciting new directions for research.

2 Pile of Law

We curate a ~ 256 GB (and growing) dataset of legal and administrative text.³ The utility of this data is twofold: (1) to aggregate legal and administrative data sources that demonstrate different norms and legal standards for data filtering; (2) to collect a dataset that can be used in the future for pretraining legal-domain language models, a key direction in access-to-justice initiatives [16]. A number of prior works have pretrained smaller models on smaller subsets of legal data, including private data that is subject to restrictive licenses [30, 129]. None of these have conducted an analysis of the legal data itself—and none have curated an open-source, legal-focused pre-training dataset at this scale.

Through extensive efforts, we compile data from 35 data sources, including legal analyses, court opinions and filings, government agency publications, contracts, statutes, regulations, casebooks, and more. Others have aggregated smaller subsets of legal data, such the EuroParl datasets which gather European Parliamentary debates [64, 57]. We have included some of these as subsets of Pile of Law when relevant and plan to continue adding material to the Pile of Law over time, further increasing its utility to the community.

We characterize the dataset in detail in Appendix E. All of the content is already entirely public and mostly available under permissive licenses, but has not previously been compiled at scale for research purposes.⁴ Each of these data sources carries with it an implicit filtering mechanism formed under relevant legal standards of privacy and toxicity, which we discuss throughout subsequent sections and in the Appendix. While the underlying data in Pile of Law is already public record and has implicit filters, we recognize that it may contain sensitive material that has escaped administrative scrutiny. We discuss the ethics of our work and our proposed mechanisms for content removal in Appendix A.

This dataset has obvious utility for pretraining legal-domain foundation models, particularly since, unlike other pretraining data, all material is under open licenses. Though not central to our work, we demonstrate this potential by training an initial BERT-large equivalent model on Pile of Law, yielding comparable results to highly context specific (but smaller) models (see Appendix F for full results). Recent research has shown in legal contexts that pretraining smaller models on highly in-domain data may be better than large models on big data [129, 29]. But in theory, there should be generalizable knowledge and skills that can be learned by training across more diverse sources of data. A well-crafted pretraining procedure that instills analogical reasoning abilities, for example, should transfer across domains. Our dataset is large and diverse enough (covering distinct areas of law like criminal law, contracts, and administrative law) to test this hypothesis in the legal domain, where our initial models can form a baseline.

²Note: while we discuss a number of privacy and toxicity standards, due to the expertise of the authors and the availability of data, this work focuses on the U.S. legal system. We address this and other limitations in Appendix B.

³<https://huggingface.co/datasets/pile-of-law/pile-of-law>.

⁴See Appendix G for a discussion of copyright and licensing in the dataset.

Table 1: Filters Applied in Major Pre-Training Papers

	PSI	Deduplication	Toxic Content	Quality
CCNet [125]	No	MinHash (pages)	No	No
C4 [98]	No	Unknown (3-sentence spans)	Word list	Minimum word counts, presence of curly brackets, ‘lorem ipsum’, etc.
GPT-3 [21]	No	MinHash (pages)	No	Train classifier to distinguish CC from curated high-quality examples
Gopher [97]	No	MinHash (pages)	Google Search	Safe- Min./max. word counts, word-to-symbol ratio, share ellipses, excessive repetition; require stop words
The Pile [44]	No	MinHash (pages)	Ad-hoc source deletion	Train classifier to distinguish CC from curated high-quality examples

3 What Can the Law Teach Us About Content Filtering?

When releasing internal documents concerning individuals, courts and governments have long struggled to balance transparency against the inclusion of private or offensive content. Model creators now face a similar struggle: what content to filter before pretraining a large language model on the data. In this section we survey how governments and courts have handled such content filtering and briefly discuss how Pile of Law implicitly encodes these privacy and toxicity rules. Based on these rules, we provide actionable lessons for researchers training large language models across fields, so that they can adapt similar rules as minimum standards for dataset sanitization. To be clear, we do not take the position that legal rules are optimal nor monolithic. But in many cases they result from a deliberative process that includes judges, legislators, and policymakers in contexts open to public scrutiny, so we think that the machine learning community can at minimum learn from these laws, rules, and norms to improve current ad hoc practice. In short, there is no need to reinvent law.

3.1 Privacy

Despite the growing focus on privacy in NLP [20], Table 1 shows that many major pre-training papers do not explicitly filter potentially sensitive information (PSI).⁵ For example, [44] excludes sources due to concerns over explicit or racist content, but does not assess the prevalence of PSI, despite including web-based sources (e.g. OpenWebText) in which users may have an expectation of anonymity. Instead, pre-training papers have focused their attention on alternatives to filtering, like deduplication [62], federated learning [108, 50], differential privacy [80, 42], and other approaches [32, 74, 65, 127, 82]. But a number of recent papers have demonstrated that large generative models output memorized content [26, 111, 32, 25, 69] even with deduplication [27]. Given that many models are trained without privacy mechanisms, filtering is critical to protecting individuals, which is perhaps why research involving health data still emphasizes that approach [90, 2]. But choosing what to filter is challenging; below, we discuss how governments and courts make such decisions.

How have governments balanced privacy against competing values? First, we examine how several jurisdictions handle privacy filtering. Table 2 provides a brief summary.⁶

Baseline Redactions. Across the jurisdictions we examine, there is a baseline level of filtering. Virtually every jurisdiction in Table 2 protects the identities of minors. At minimum, juveniles must be protected by pseudonyms in public judgments, and outside of some U.S. states, juvenile criminal records are not public. No jurisdiction normally permits the publication of financial account numbers,

⁵We define PSI to mean information that could violate a person’s privacy interests. This could include personally identifiable information, including a person’s name, date of birth, or identification number. Under this definition, a document can contain some PSI (e.g. a name or the facts of a case) while excluding other PSI (e.g. date of birth). But some information that is personally identifiable is not PSI; for example, the name and office contact information of an attorney filing a court brief is identifying but not sensitive.

⁶See Appendix I for a complete version of the table, including citations.

Table 2: Availability of Identifying Information Across Administrative Settings

Jurisdiction	Civil Cases	Criminal Cases	Juvenile Data
U.S. Federal Courts	All case details public unless sealed, except DOBs, ID/account #s.	Def. names public; DOBs, ID/account #s, addresses redacted.	Criminal records confidential. Names redacted from civil cases.
U.S. Admin. Agencies	Most PII omitted from public records.	-	No statute; more protection in practice.
German Courts	Judgments omit all identifying information.	Confidential 3-5 years after sentence completed.	No public access to criminal records.
Chinese Courts	Names/case details public except in certain classes of cases.	Names/case details are public as of 2016.	Juvenile criminal records are categorically exempt from disclosure.
Canadian Courts	Presumption of openness, except specific details and rare sealed cases.	Public; may be sealed after a period of good behavior.	Youth criminal records are always confidential.

dates of birth, or identity numbers like social security numbers.⁷ All of these are bright line rules directly applicable to text corpora.

Value-system contexts. There are also significant points of disagreement corresponding to the role of privacy in different value systems. U.S., Chinese and Canadian courts denote the names of litigants in ordinary civil cases, prioritizing public access and transparency; German courts do not. Likewise, U.S. federal courts virtually never remove criminal cases from the public record [109, p. 1233], a rule also emerging in China [75]. Canada allows most criminal records to be expunged after a period of good behavior. And in Germany, virtually all criminal records are automatically sealed after a set time, and courts have even imposed fines for publicizing a person’s criminal history after expungement [22].

Contextualized privacy. Digging further into these rules highlights how court privacy rules account for context. In the U.S. and Canada, the public disclosure of litigants’ potentially sensitive information (PSI) can be avoided by persuading a court that extenuating circumstances apply [124, 115]. To name one example, courts generally permit pseudonyms when parties allege that they have suffered a sexual assault [124, p. 57]. The chance to *seal* a case, or to file pseudonymously, suggests that even the most open judicial regime allows for censoring in exceptional cases—although the sealing and pseudonymous filing standards suffer from inconsistency and misuse [113].

Likewise, administrative agencies often employ context-aware heuristics when deciding whether to include PSI in public decisions. Although administrative courts are not generally bound by stringent privacy rules like HIPAA [110, 36], the Department of Justice exempts immigration applications from public scrutiny due to privacy concerns [87]; the same is true of Social Security Disability applications. Cases involving veterans’ benefits are released pseudonymously.

Public availability is not a limit. In many cases, the rules for sanitizing PSI and sealing cases do not depend on whether information is already public. For example, the ban on publicly filing documents revealing dates of birth in U.S. federal courts does not depend on whether a litigant’s birth date is otherwise public [122]. In cases where a court does take into account the public availability of information (e.g., sealing standards [115, 121]), contextual countervailing factors can justify keeping a case sealed.

Implications for Pile of Law. All of the above privacy norms mean that each subset of Pile of Law is already filtered for privacy based on legal norms in that jurisdiction. Further filtering could seek to align the whole dataset with the norms in one of the subsets prior to pretraining. Appendix E summarizes the filtering norms present in each subset of the data.

⁷The United States’ Federal Rule of Civil Procedure 5.2 lays out exceptions when these facts are contained in judicial records properly before a federal court and for civil asset forfeiture cases.

Lessons for researchers. First, the law provides a number of useful heuristics that researchers could deploy to sanitize data. Detecting and redacting juvenile names, dates of birth, and account and identity numbers is virtually always appropriate across countries. Legal protections for already-public information show why sanitization may be necessary even for text collected from public-facing web pages. Second, the U.S. system appears to lean most heavily toward transparency. We suggest that researchers can use the U.S. court rules as a floor. Such privacy filtering rules would already go beyond much of current modeling practice. Third, in addition to consensus heuristics, researchers should make contextualized decisions about privacy harms. While this may seem difficult, Section 4 demonstrates how to leverage Pile of Law to learn contextualized standards to mimic legal privacy redaction mechanisms; alternatively, allowing individuals whose information appears in the training corpus to request removal may serve as another stopgap. Last, the U.S. legal rules do not extend as far as some researchers suggest is necessary. For example, [5] suggests that *all* names must be redacted to preserve privacy. This would reflect greater privacy protection than is typically afforded by U.S. law, which prioritizes public openness and transparency about court proceedings, but would be in line with German rules. These pose important value tradeoffs, and we suggest that researchers look as a starting point to the jurisdiction that aligns with such value tradeoffs for filtering other potentially sensitive content.

3.2 Toxicity

How is toxic speech defined in research? The category of ‘toxic speech’ is defined in multiple ways [47, 123, 4]. Some papers define toxicity as “*disrespectful* comments, including . . . identity *attacks*, profanity and *threats*,” thus emphasizing the idea of intentional insult [39, 126, 128]. A broader definition would incorporate *implicit* toxicity, as when a speaker “subtly” or “unconsciously expresses a prejudiced attitude” [18, 70]. [18] cites the example of the question “But where are you from, originally?” Others would take a still broader view, suggesting that any *profanity* is toxic (in addition to hate speech and derogatory content) irrespective of speaker intent [89]. One implication of these divergent choices concerns *mentions* of toxic language, where a speaker refers to something said by another [114]. For example, if a judge writes that “Plaintiff claims that her supervisor called her ‘___’” (where ___ is a profane epithet), an intent-based standard typically would not deem the use of ___ ‘toxic,’ while an approach targeting profanity typically would.

How have governments regulated toxic content? Scholars have documented the role of the law in institutional racism and other forms of oppression, and legal materials from prior eras use words that would by modern standards be considered epithets [31, 13]. Today, the legal profession in most Anglophone countries strongly polices overt discriminatory epithets [38]. Overtly biased speech is prohibited for judges and lawyers in the U.S., Canada, and the U.K. by professional rules [3, 24, 119, 68]; similar norms have been put forward by the U.N. [120]. Judges and lawyers in all countries are routinely sanctioned when they use racist epithets, and most incidents occur verbally or off the bench [38, 60, 118].

Unlike overt, indecorous epithets, legal norms permit the use of speech affected by implicit bias; the incidence of such speech is well-documented [94, 100, 85, 12]. Indeed, it is sometimes encoded in the laws judges and administrative agencies enforce [23]. Furthermore, some lawyers may see themselves as professionally *obligated* to deploy stereotypes when doing so may assist their clients (e.g. immigration [88], defendants in sexual assault cases [34]).

Implications for Pile of Law. The adversarial legal system in many Anglophone countries creates incentives for lawyers to complain about overt racism in written materials, which would violate unambiguous professional rules. Thus, the appearance of epithets in our data is more likely to be confined to quotations, mentions, or to historical legal materials. However, text in our corpus may be toxic according to other definitions; for example, we are unable to quantify the prevalence of implicit biases or offensive stereotyping. Explicit racial, sexual, or offensive terms do appear in modern legal text, but most often in the form of a quotation than direct use. For instance, many cases revolve around evidence documenting racial or gender discrimination, and judges commonly spell out profane or explicit words from the evidentiary record [63, 45]. Finally, elected officials in our legislative transcripts are not bound by the same professional norms as attorneys. Additionally, an interesting future examination may note differences between civil law and common law systems,

examining rates of offensive content between the different legal systems and norms. We provide a per-subset examination of filtering norms in Appendix E.

Lessons for researchers. First, as is true of privacy, the toxicity norms prevalent in many legal systems offer a lower-bound for researchers. Researchers seeking to mimic the standards that apply in courts should sanitize intentional uses of derogatory terms from pretraining data. That said, current filters are not precise enough to handle this standard. Under the rules applicable to lawyers, filters based on simple word lists would be over-inclusive because they would capture *references* to offensive language that may be non-toxic in context. Second, the rules applied in courts suggest that generative models should portray toxic behavior explicitly in some contexts, either to serve the values of ‘accuracy and precision’ or to persuade readers [63, p. 7]; but as [43] argues, this view is contested.

Third, in language model pretraining, there may be reason to exceed minimum judicial standards depending on the length of content needed to contextualize references to offensive speech. Accessible language models like Roberta [78] have a maximum context window of 512 tokens. If a reference to offensive content spans the majority of these tokens, the model will simply uptake the offensive content as if it were being trained for *direct use*. As model contexts grow, it may become more reasonable for researchers to adopt judicial norms.

4 What Can We Learn from Legal Text?

As Section 3 shows, even jurisdictions that impose a strong presumption of transparency on legal documents often allow for contextual decisions that weigh this presumption against the potential harms caused by the inclusion of PSI on the public record. Reducing these rules to tools that can be deployed for filtering may be challenging. But Pile of Law encodes these contextual decisions already, providing a rich opportunity to learn context-aware norms directly. This section demonstrates the promise of Pile of Law for operationalizing legal norms. While not comprehensive, the experiments below demonstrate a path forward for replicating the content-filtering mechanisms of courts and governments by leveraging variation in Pile of Law. In particular, we show that: (1) Pile of Law reflects variation in privacy norms that can be leveraged to learn contextual privacy rules, such as when to redact names in potentially harmful situations; (2) Pile of Law reflects variation in toxicity norms over time and across contexts, toxicity filters fall short of handling these nuances, and researchers can learn much from building toxicity filters that can handle nuances in Pile of Law’s text.

4.1 Learning Contextual Privacy Rules

Case Study 1: Pseudonymity in Immigration Court. The Board of Immigration Appeals (BIA) evaluates petitions appealing immigration decisions and sometimes publishes precedential decisions that affect future cases. Some cases include applicants’ full names, while others replace them with pseudonymous initials. We demonstrate how subsets of the data can be used to learn the value judgements made in making this pseudonymity decision. We split cases into paragraphs and mask terms used to refer to the applicant. We train a distill-BERT base model [105] to predict whether the paragraph should use pseudonymity or not. This model achieves ~80% F1 on the validation set. We then examine what types of content are more likely to trigger a pseudonymity recommendation by conducting a perturbation analysis. We use the Bias in Bios dataset [35], censored for names and pronouns. We prepend an additional sentence to each biography that indicates whether the person: (1) is seeking asylum or is a refugee; (2) experienced torture; (3) committed a non-violent criminal offense; or (4) committed a violent criminal offense. Figure 1 shows that asylum and torture sentences were more likely to trigger pseudonymity while criminal offenses were less likely. This aligns with federal regulations that prevent disclosure of information related to asylum or the Convention Against Torture (8 CFR § 208.6(a)). By contrast, federal regulations allow information disclosure when a criminal proceeding is involved (8 CFR § 208.6(d)(1)(ii)), though no regulation addresses criminal history.

Next, we fit a causal lexicon using the deep residualization method (and associated library) from Pryzant et al. [95]. We control for the year that a case was published since we found that some aspects of privacy standards have shifted year-to-year, which provides a unique opportunity to learn evolving standards of privacy. We select the top 100 most indicative terms for pseudonymity and remove those

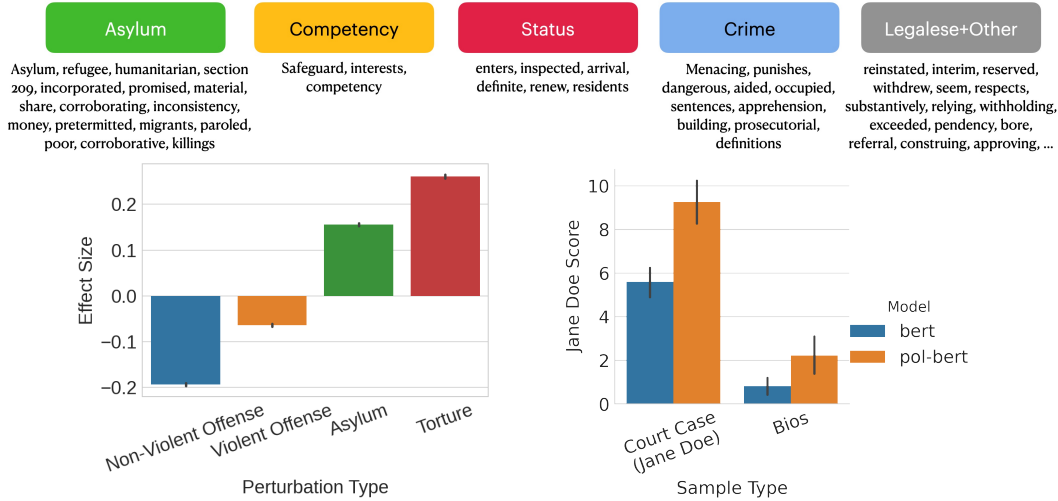


Figure 1: (Top) A causal lexicon learned for the EOIR privacy task, manually sorted by topic with contextual information. Extended version and information in Appendix H.1. (Bottom Left) A distill-bert model is more likely to predict pseudonymity for bios with an asylum or torture perturbation (effect size is difference in pseudonymity likelihood from normal bio and bio with added perturbation). (Bottom Right) Jane Doe Score is the difference in MLM score between a version of the sentence using Jane Doe and a random name. The sample sources are paragraphs using pseudonyms and Bios [35] (no pseudonyms).

where the term only showed up in one case. Then we manually examine contexts and cluster terms into categories. We found that terms most likely to be associated with pseudonymity could be largely clustered into: asylum, mental competency (a legal term used to refer to one’s ability to stand trial), immigration status, and indications of a criminal proceeding. We also find that many terms associated with general legal language were included, suggested some remaining confounding and the need for more research into text-based causal attribution. These causal lexicons are seen in Figure 1.

Case Study 2: Pseudonyms in Civil Litigation. Next, we look to a “zero-shot” version of the experiment above in a broader setting. As noted in Section 3, litigants in U.S. courts can ask to use pseudonyms like “Jane Doe” in court documents, for example in harassment suits. To assess these requests, courts consider contextual factors like “sensitive and highly personal” subject matter, minors, or other extenuating circumstances [124]. We collect ~ 500 paragraphs where a pseudonym (“Jane Doe” or “Jane Roe”) is used from the validation part of the Court Listener Opinions data. For each sentence, we create 100 alternative sentences that replace “Jane Doe” with a name sampled using 1990 Census probabilities (using NAMES). We then compare whether each model is more likely to guess “Jane Doe” using MLM Score [103]. We repeat this process on the Bios dataset [35]. Figure 1 shows that a model trained on Pile of Law (pol-bert) ranks Jane Doe ~ 3 points higher than a standard bert-large-uncased on true pseudonym cases. This suggests that models pre-trained on Pile of Law are more likely to encode appropriate pseudonymity norms. To be sure, pol-bert is slightly more biased for Jane Doe use overall, as is to be expected, but its performance gains persist even after accounting for this bias.

Case Study 3: Privacy Standards in Medical Cases. We examine *inter-source variation* between the Board of Veterans Appeals (BVA) and the Department of Labor’s Employee’s Compensation Appeals Board (DOL). Leading tools for data sanitization remove personal health information as defined by HIPAA [33], including dates or the name of a physician [90]. We ran [90] on all decisions by the BVA and DOL since both adjudicate the extent of applicants’ disabilities, though they are not bound by HIPAA [36]. Showing the difficulty of applying sanitization tools out of domain, virtually *all* decisions included information flagged as HIPAA-protected: 99% included dates; 96% of BVA and 100% of DOL decisions included medical facility names. But the two agencies also differed. About 26% of DOL cases but just 0.36% of BVA cases included a physician name. Physician fraud is more common in worker’s compensation programs like DOL’s [91], but the BVA relies on the

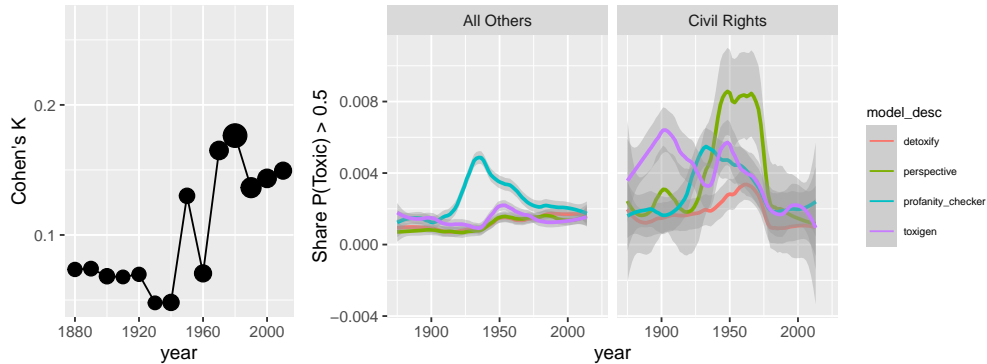


Figure 2: Inter-Model Agreement and Toxicity Over Time. Left: Cohen’s κ , by 10-year bin, calculated for [130] and [59], with dot size proportional to number of examples. Right: Share of sentences assigned $>50\%$ probability of being toxic, by model, time, and topic classification [112].

testimony of VA physicians. The transparency in DOL opinions reflects the higher public interest in physician accountability.

Lessons for researchers. These experiments show that the Pile of Law encodes signals about privacy standards that can be learned to produce more nuanced recommendations about filtering. For example, researchers may consider whether to mimic the EOIR standard to remove names in proceedings related to minors, asylum or safety concerns. Or they may wish to learn and apply the more contextual standard that is used in general U.S. litigation, where a complex set of factors is used to justify the exclusion of names from case texts. Such contextualized filters may help ensure that generative models strike the right balance between accuracy and privacy protection, for example by accurately distinguishing benign releases of names and contact information (e.g., in response to queries about government officials) from harmful ones (sensitive circumstances where harm is plausible).

4.2 Calibration and Value-Alignment in Toxicity Filtering

We also identify three main insights (and challenges) from using toxicity filters on Pile of Law, setting the ground for new research using the dataset: (1) toxicity filters often disagree, creating potential issues for automated filtering; (2) toxicity filters may be value-misaligned when it comes to content that is flagged in Pile of Law; (3) toxicity scores vary highly with the length of the content, making it unclear how to handle long-document filtering.

Case Study: Supreme Court Decisions. Leveraging Pile of Law, we show that there are profound nuances to filtering toxic content. First, toxicity filters encode value judgements and divergent definitions of toxicity. Figure 2 shows Cohen’s κ between profanity-checker and Perspective over time for sentences in Supreme Court cases (Fig. 6 shows the same for all filters). At the sentence level, the tools’ agreement rates are very low, but rise over time, indicating the challenge of handling out-of-domain data far away in time. A vivid example of this challenge is provided in Table 3: *Dred Scott* is the most notoriously racist decision in U.S. history [46], but perhaps due to the archaic language of its holding, *none* of the models is sure that it is toxic.

But civil rights cases illustrate why the disagreement is about conceptual differences, not just domain drift. Figure 2 shows that the period between 1950 and 1970 is associated with a large spike in the share of sentences deemed toxic in U.S. civil rights cases. This period was associated with the end of *de jure* segregation in the United States [117]. Many cases likely *quoted* or *mentioned* racist laws before striking them down. For instance, in Table 3, *Loving* describes a law banning interracial marriage in order to deem it illegal. Quoting this language qualifies as toxic under some but not all definitions, and as Figure 2 shows, that view is encoded in some but not all filters. Accordingly, the filters disagree as to whether *Loving*’s quote is clearly toxic. Document-level filtering could thus easily delete core civil rights cases like *Hunter* and *Loving*—while leaving in *Dred Scott*.

Finally, we note that the context window used to filter out sentences appears to dramatically influence ratings. Perspective segments data into sentences and then labels each sentence, which is the

Table 3: Toxicity Ratings of Quotes From the U.S. Supreme Court, Showing Rating Disagreement

Case	Quote	(1)	(2)	(3)	(4)
Hunter v. Erickson (1969)	“The majority needs no protection against discrimination.”	0.02	0.05	0.00	0.81
Loving v. Virginia (1967)	“[I]f any white person intermarry with a colored person . . . he shall be guilty of a felony and shall be punished by confinement in the penitentiary”	0.52	0.54	0.60	0.94
Dred Scott v. Sandford (1857)	“A free negro of the African race whose ancestors were brought to this country and sold as slaves is not a citizen within the meaning of the Constitution.”	0.29	0.50	0.26	0.54

Note: Model (1) is profanity-check [130]; (2) is Perspective [59]; (3) is Detoxify [49]; and (4) is Toxigen [51].

approach we take above. We find that by using longer span, we can *systematically decrease* the perceived toxicity of a span, even if it is obviously toxic under any definition. We take the top 5k sentences labeled as toxic by Toxigen. We then take 2 sentences before and after the toxic sentence (clamped to the boundaries of the document). We find that the toxicity score drops between **55-57%** (absolute, 95% CI) just by adding this context. While some of this change might be due to correct re-classification of mentions, we provide qualitative examples in which this is clearly untrue in Appendix Table 8.

Lessons for researchers. The experiments above demonstrate that, while toxicity filtering is important to align with the courts’ modern lower bounds banning uses of epithets, it is not clear that existing filters are not consistent and filter out content aligned with different values. Moreover, they can arbitrarily label content as non-toxic in long-document or out-of-distribution settings, which may affect filtering mechanisms. More work is needed to create robust, value-aligned toxicity filters for pretraining and it is unclear if off-the-shelf mechanisms strike the right balance. As we have shown, the Pile of Law provides unique opportunities to develop such methods.

5 Conclusion

In this work we have examined how the law and legal data can inform data filtering practices that are of great importance to responsible large language model training. We provide an extensive legal dataset (the Pile of Law) and illustrate a number of exciting new research directions for future work.

Acknowledgements

We thank SambaNova Systems for generously providing compute resources via the SambaNova Systems Dataflow-as-a-Service™ platform and the Stanford Institute for Human-Centered Artificial Intelligence for computing support. We also thank Jieru Hu for helpful discussions and Krithika Iyer for technical assistance. PH is supported by an Open Philanthropy Project AI Fellowship.

References

- [1] Samira Abnar, Mostafa Dehghani, Behnam Neyshabur, and Hanie Sedghi. Exploring the Limits of Large Scale Pre-training. In *International Conference on Learning Representations*, 2021.
- [2] Allison Adams, Eric Aili, Daniel Aioanei, Rebecca Jonsson, Lina Mickelsson, Dagmar Mikmekova, Fred Roberts, Javier Fernandez Valencia, and Roger Wechsler. Anonymate: A toolkit for anonymizing unstructured chat data. In *Proceedings of the Workshop on NLP and Pseudonymisation*, pages 1–7, 2019.
- [3] American Bar Association. Model Code of Judicial Conduct, 2007. https://www.americanbar.org/groups/professional_responsibility/publications/model_code_of_judicial_conduct.

- [4] Noman Ashraf, Arkaitz Zubiaga, and Alexander Gelbukh. Abusive language detection in youtube comments leveraging replies as conversational context. *PeerJ Computer Science*, 7: e742, 2021.
- [5] Tuomas Aura, Thomas A Kuhn, and Michael Roe. Scanning electronic documents for personally identifiable information. In *Proceedings of the 5th ACM workshop on Privacy in Electronic Society*, pages 41–50, 2006.
- [6] Eugene Bagdasaryan, Omid Poursaeed, and Vitaly Shmatikov. Differential privacy has disparate impact on model accuracy. In *Proceedings of the 33rd International Conference on Neural Information Processing Systems*, pages 15479–15488, 2019.
- [7] Azeem Bande-Ali and Walker Boyle. U.S. Supreme Court annotated transcripts. Github, 2019. https://github.com/walkerdb/supreme_court_transcripts.
- [8] Alexander Baturo, Niheer Dasandi, and Slava J Mikhaylov. Understanding state preferences with text as data: Introducing the UN General Debate corpus. *Research & Politics*, 4(2), 2017.
- [9] Jason Baumgartner, Savvas Zannettou, Brian Keegan, Megan Squire, and Jeremy Blackburn. The Pushshift Reddit Dataset. In *14th International Conference on Web and Social Media*, pages 830–839, 2020.
- [10] Steven Bellman, Eric J. Johnson, Stephen J. Kobrin, and Gerald L. Lohse. International Differences in Information Privacy Concerns: A Global Survey of Consumers. *The Information Society*, 20(5):313–324, 2004.
- [11] Emily M Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. On the Dangers of Stochastic Parrots: Can Language Models Be Too Big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency*, pages 610–623, 2021.
- [12] Janine Benedet. Judicial Misconduct in the Sexual Assault Trial. *U.B.C. L. Rev.*, 52:1, 2019.
- [13] Mary Frances Berry. *The Pig Farmer’s Daughter and Other Tales of American Justice: Episodes of racism and sexism in the courts from 1865 to the present*. Vintage, 2011.
- [14] Andrew Blair-Stanek, Nils Holzenberger, and Benjamin Van Durme. Tax Law NLP Resources, 2020. <https://doi.org/10.7281/T1/N1X6I4>.
- [15] Michael J Bommarito II, Daniel Martin Katz, and Eric M Detterman. LexNLP: Natural language processing and information extraction for legal and regulatory texts. In *Research Handbook on Big Data Law*. Edward Elgar Publishing, 2021.
- [16] Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Koh, Mark Krass, Ranjay Krishna, Rohith Kudithipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avaniika Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E. Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. On the Opportunities and Risks of Foundation Models. *arXiv e-prints*, art. arXiv:2108.07258, August 2021.

- [17] Łukasz Borchmann, Dawid Wisniewski, Andrzej Gretkowski, Izabela Kosmala, Dawid Jurkiewicz, Łukasz Szalkiewicz, Gabriela Pałka, Karol Kaczmarek, Agnieszka Kaliska, and Filip Graliński. Contract Discovery: Dataset and a Few-Shot Semantic Retrieval Challenge with Competitive Baselines. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4254–4268. Association for Computational Linguistics, November 2020.
- [18] Luke Breiffeller, Emily Ahn, David Jurgens, and Yulia Tsvetkov. Finding Microaggressions in the Wild: A Case for Locating Elusive Phenomena in Social Media Posts. In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 1664–1674, Hong Kong, China, November 2019. Association for Computational Linguistics.
- [19] British Columbia Law Society. Code of Professional Conduct for British Columbia, 2021. <https://www.lawsociety.bc.ca/support-and-resources-for-lawyers/act-rules-and-code/code-of-professional-conduct-for-british-columbia/>.
- [20] Hannah Brown, Katherine Lee, Fatemehsadat Mireshghallah, Reza Shokri, and Florian Tramèr. What Does it Mean for a Language Model to Preserve Privacy? In *ACM Conference on Fairness, Accountability, and Transparency (FAcT '22)*, 2022.
- [21] Tom B Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. Language models are few-shot learners. In *Proceedings of the 34th International Conference on Neural Information Processing Systems*, pages 1877–1901, 2020.
- [22] Dana Burchardt. Backlash against the Court of Justice of the EU? The Recent Jurisprudence of the German Constitutional Court on EU Fundamental Rights as a Standard of Review. *German Law Journal*, 21:1–18, 2020. doi: doi:10.1017/glj.2020.16.
- [23] Naomi R Cahn. Looseness of legal language: The reasonable woman standard in theory and in practice. *Cornell L. Rev.*, 77:1398, 1991.
- [24] Canadian Judicial Council. Ethical Principles for Judges, 2004. https://cjc-ccm.ca/cmslib/general/news_pub_judicialconduct_Principles_en.pdf.
- [25] Nicholas Carlini, Chang Liu, Úlfar Erlingsson, Jernej Kos, and Dawn Song. The secret sharer: Evaluating and testing unintended memorization in neural networks. In *28th USENIX Security Symposium (USENIX Security 19)*, pages 267–284, 2019.
- [26] Nicholas Carlini, Florian Tramèr, Eric Wallace, Matthew Jagielski, Ariel Herbert-Voss, Katherine Lee, Adam Roberts, Tom Brown, Dawn Song, Úlfar Erlingsson, et al. Extracting training data from large language models. In *30th USENIX Security Symposium (USENIX Security 21)*, pages 2633–2650, 2021.
- [27] Nicholas Carlini, Daphne Ippolito, Matthew Jagielski, Katherine Lee, Florian Tramèr, and Chiyuan Zhang. Quantifying memorization across neural language models. *arXiv preprint arXiv:2202.07646*, 2022.
- [28] Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, and Ion Androutsopoulos. Large-scale multi-label text classification on EU legislation. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6314–6322, Florence, Italy, 2019. Association for Computational Linguistics.
- [29] Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. LEGAL-BERT: The Muppets straight out of Law School. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2898–2904, 2020.
- [30] Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. Lexglue: A benchmark dataset for legal language understanding in english. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4310–4330, 2022.
- [31] Sumi Cho. Post-racialism. *Iowa L. Rev.*, 94:1589, 2008.

- [32] Maximin Coavoux, Shashi Narayan, and Shay B Cohen. Privacy-preserving Neural Representations of Text. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 1–10, 2018.
- [33] I Glenn Cohen and Michelle M Mello. HIPAA and protecting health information in the 21st century. *JAMA*, 320(3):231–232, 2018.
- [34] Elaine Craig. The ethical obligations of defence counsel in sexual assault cases. *Osgoode Hall L. J.*, 51:427, 2013.
- [35] Maria De-Arteaga, Alexey Romanov, Hanna Wallach, Jennifer Chayes, Christian Borgs, Alexandra Chouldechova, Sahin Geyik, Krishnaram Kenthapadi, and Adam Tauman Kalai. Bias in bios: A case study of semantic representation bias in a high-stakes setting. In *proceedings of the Conference on Fairness, Accountability, and Transparency*, pages 120–128, 2019.
- [36] Department of Veterans Affairs. Health Insurance Portability and Accountability Act Applicability in VBA, 2003. <https://www.va.gov/ogc/docs/ADV3-2003.pdf>.
- [37] Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, 2019.
- [38] Jana DiCosmo. Racism in the Legal Profession: A Racist Lawyer Is an Incompetent Lawyer. *Nat’l Law. Guild Rev.*, 75:82, 2018.
- [39] Lucas Dixon, John Li, Jeffrey Sorensen, Nithum Thain, and Lucy Vasserman. Measuring and mitigating unintended bias in text classification. In *Proceedings of the 2018 AAAI/ACM Conference on AI, Ethics, and Society*, pages 67–73, 2018.
- [40] Jesse Dodge, Maarten Sap, Ana Marasović, William Agnew, Gabriel Ilharco, Dirk Groeneveld, Margaret Mitchell, and Matt Gardner. Documenting Large Webtext Corpora: A Case Study on the Colossal Clean Crawled Corpus. In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 1286–1305, 2021.
- [41] Zachary Elkins, Tom Ginsburg, James Melton, Robert Shaffer, Juan F Sequeda, and Daniel P Miranker. Constitute: The world’s constitutions to read, search, and compare. *Journal of web semantics*, 27:10–18, 2014.
- [42] Oluwaseyi Feyisetan, Abhinav Aggarwal, Zekun Xu, and Nathanael Teissier. Research challenges in designing differentially private text generation mechanisms. *arXiv preprint arXiv:2012.05403*, 2020.
- [43] Richard Thompson Ford. Racial epithets and racial etiquette. *Capital University Law Review*, 49(4):527–534, 2021.
- [44] Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, et al. The Pile: An 800GB dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*, 2020.
- [45] Alan E Garfield. To swear or not to swear: using foul language during a Supreme Court oral argument. *Wash. U. L. Rev.*, 90:279, 2012.
- [46] Jamal Greene. The anticanon. *Harv. L. Rev.*, 125:379, 2011.
- [47] Tommi Gröndahl, Luca Pajola, Mika Juuti, Mauro Conti, and N Asokan. All You Need is “Love” Evading Hate Speech Detection. In *Proceedings of the 11th ACM workshop on artificial intelligence and security*, pages 2–12, 2018.
- [48] Suchin Gururangan, Dallas Card, Sarah K Drier, Emily K Gade, Leroy Z Wang, Zeyu Wang, Luke Zettlemoyer, and Noah A Smith. Whose language counts as high quality? Measuring language ideologies in text data selection. *arXiv preprint arXiv:2201.10474*, 2022.

- [49] Laura Hanu and Unitary team. Detoxify. Github, 2020. <https://github.com/unitaryai/detoxify>.
- [50] Andrew Hard, Kanishka Rao, Rajiv Mathews, Swaroop Ramaswamy, Françoise Beaufays, Sean Augenstein, Hubert Eichner, Chloé Kiddon, and Daniel Ramage. Federated learning for mobile keyboard prediction. *arXiv preprint arXiv:1811.03604*, 2018.
- [51] Thomas Hartvigsen, Saadia Gabriel, Hamid Palangi, Maarten Sap, Dipankar Ray, and Ece Kamar. Toxigen: A large-scale machine-generated dataset for adversarial and implicit hate speech detection. *arXiv preprint arXiv:2203.09509*, 2022.
- [52] David Hausman, Daniel E Ho, Mark Krass, and Anne M McDonough. Executive Control of Agency Adjudication: Capacity, Selection and Precedential Rulemaking. *Journal of Law, Economics & Organization*, 40, 2024.
- [53] Allison Hegel, Marina Shah, Genevieve Peaslee, Brendan Roof, and Emad Elwany. The Law of Large Documents: Understanding the Structure of Legal Contracts Using Visual Cues. *arXiv preprint arXiv:2107.08128*, 2021.
- [54] Peter Henderson, Jieru Hu, Joshua Romoff, Emma Brunskill, Dan Jurafsky, and Joelle Pineau. Towards the systematic reporting of the energy and carbon footprints of machine learning. *Journal of Machine Learning Research*, 21(248):1–43, 2020.
- [55] Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. CUAD: An Expert-Annotated NLP Dataset for Legal Contract Review. *arXiv preprint arXiv:2103.06268*, 2021.
- [56] Zihan Huang, Charles Low, Mengqiu Teng, Hongyi Zhang, Daniel E Ho, Mark S Krass, and Matthias Grabmair. Context-aware legal citation recommendation using deep learning. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, pages 79–88, 2021.
- [57] Javier Iranzo-Sánchez, Joan Albert Silvestre-Cerda, Javier Jorge, Nahuel Roselló, Adria Giménez, Albert Sanchis, Jorge Civera, and Alfons Juan. Europarl-st: A multilingual corpus for speech translation of parliamentary debates. In *ICASSP 2020-2020 IEEE International Conference on Acoustics, Speech and Signal Processing (ICASSP)*, pages 8229–8233. IEEE, 2020.
- [58] Yacine Jernite, Huu Nguyen, Stella Biderman, Anna Rogers, Maraim Masoud, Valentin Danchev, Samson Tan, Alexandra Sasha Luccioni, Nishant Subramani, Isaac Johnson, Somaieh Nikpoor, Jörg Frohberg, Aaron Gokaslan, Peter Henderson, Rishi Bommasani, and Margaret Mitchell. Data governance in the age of large-scale data-driven language technology. In *2022 ACM Conference on Fairness, Accountability, and Transparency*, pages 2206–2222, 2022.
- [59] Jigsaw and Google’s Counter Abuse Technology team. Perspective, 2021. <https://perspectiveapi.com/>.
- [60] Sheri Lynn Johnson, John H Blume, and Patrick M Wilson. Racial Epithets in the Criminal Process. *Mich. St. L. Rev.*, page 755, 2011.
- [61] Lisa LaPlant Jon Quandt, Eric Mill. govinfo. Github. <https://github.com/usgpo/api>, 2018.
- [62] Nikhil Kandpal, Eric Wallace, and Colin Raffel. Deduplicating training data mitigates privacy risks in language models. *arXiv preprint arXiv:2202.06539*, 2022.
- [63] Randall Kennedy and Eugene Volokh. The new taboo: Quoting epithets in the classroom and beyond. *Cap. U. L. Rev.*, 49:1, 2021.
- [64] Philipp Koehn. Europarl: A parallel corpus for statistical machine translation. In *Proceedings of Machine Translation Summit*, pages 79–86, 2005.
- [65] Satyapriya Krishna, Rahul Gupta, and Christophe Dupuy. ADePT: Auto-encoder based Differentially Private Text Transformation. In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 2435–2439, 2021.

- [66] Alexandre Lacoste, Alexandra Luccioni, Victor Schmidt, and Thomas Dandres. Quantifying the carbon emissions of machine learning. *arXiv preprint arXiv:1910.09700*, 2019.
- [67] David S Law and Mila Versteeg. The declining influence of the United States Constitution. *N. Y. U. L. Rev.*, 87:762, 2012.
- [68] Law Society of Ontario. Rules of Professional Conduct, 2022. <https://lso.ca/about-lso/legislation-rules/rules-of-professional-conduct>.
- [69] Jooyoung Lee, Thai Le, Jinghui Chen, and Dongwon Lee. Do language models plagiarize? *arXiv preprint arXiv:2203.07618*, 2022.
- [70] Alyssa Lees, Daniel Borkan, Ian Kivlichan, Jorge Nario, and Tesh Goyal. Capturing covertly toxic speech via crowdsourcing. In *Proceedings of the First Workshop on Bridging Human-Computer Interaction and Natural Language Processing*, pages 14–20, Online, April 2021. Association for Computational Linguistics.
- [71] Spyretta Leivaditi, Julien Rossi, and Evangelos Kanoulas. A benchmark for lease contract review. *arXiv preprint arXiv:2010.10386*, 2020.
- [72] Mark A Lemley and Bryan Casey. Fair learning. *Tex. L. Rev.*, 99:743, 2020.
- [73] Justin D Levinson, Mark W Bennett, and Koichi Hioki. Judging implicit bias: A national empirical study of judicial stereotypes. *Fla. L. Rev.*, 69:63, 2017.
- [74] Yitong Li, Timothy Baldwin, and Trevor Cohn. Towards robust and privacy-preserving text representations. In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 25–30, 2018.
- [75] Benjamin L Liebman, Margaret E Roberts, Rachel E Stern, and Alice Z Wang. Mass Digitization of Chinese Court Decisions. *Journal of Law and Courts*, Fall:176–201, 2020.
- [76] Marco Lippi, Przemysław Pałka, Giuseppe Contissa, Francesca Lagioia, Hans-Wolfgang Micklitz, Giovanni Sartor, and Paolo Torroni. Claudette: an automated detector of potentially unfair clauses in online terms of service. *Artificial Intelligence and Law*, 27(2):117–139, 2019.
- [77] Michael Lissner and Brian W Carver. Courtlistener.com: A platform for researching and staying abreast of the latest law. 2010. <http://www.courtlistener.com>.
- [78] Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. Roberta: A robustly optimized bert pretraining approach. *arXiv preprint arXiv:1907.11692*, 2019.
- [79] Eneldo Loza Mencía and Johannes Fürnkranz. Efficient multilabel classification algorithms for large-scale problems in the legal domain. In *Semantic Processing of Legal Texts*, pages 192–215. Springer, 2010.
- [80] Lingjuan Lyu, Xuanli He, and Yitong Li. Differentially private representation for nlp: Formal guarantee and an empirical study on privacy and fairness. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2355–2365, 2020.
- [81] Sibella Matthews, Vincent Schiraldi, and Lael Chester. Youth justice in Europe: Experience of Germany, the Netherlands, and Croatia in providing developmentally appropriate responses to emerging adults in the criminal justice system. *Justice Evaluation Journal*, 1(1):59–81, 2018.
- [82] H Brendan McMahan, Galen Andrew, Ulfar Erlingsson, Steve Chien, Ilya Mironov, Nicolas Papernot, and Peter Kairouz. A general approach to adding differential privacy to iterative training procedures. *arXiv preprint arXiv:1812.06210*, 2018.
- [83] Microsoft. Microsoft presidio. Github., 2022. <https://github.com/microsoft/presidio/>.
- [84] Eric Mill. Opening up government reports through teamwork and open data. *OpenGov Voices*, 2014.

- [85] Suzanne J Miller. Judicial language in new jersey sexual violence cases. *Rutgers U. L. Rev.*, 73:141, 2020.
- [86] Margaret Mitchell, Simone Wu, Andrew Zaldivar, Parker Barnes, Lucy Vasserman, Ben Hutchinson, Elena Spitzer, Inioluwa Deborah Raji, and Timnit Gebru. Model cards for model reporting. In *Proceedings of the conference on fairness, accountability, and transparency*, pages 220–229, 2019.
- [87] Nancy Morawetz. A better balance for federal rules governing public access to appeal records in immigration cases. *Hastings LJ*, 69:1271, 2017.
- [88] Deborah A Morgan. Not gay enough for the government: Racial and sexual stereotypes in sexual orientation asylum cases. *Law & Sexuality: Rev. Lesbian, Gay, Bisexual & Transgender Legal Issues*, 15:135, 2006.
- [89] Chikashi Nobata, Joel Tetreault, Achint Thomas, Yashar Mehdad, and Yi Chang. Abusive language detection in online user content. In *Proceedings of the 25th international conference on world wide web*, pages 145–153, 2016.
- [90] Beau Norgeot, Kathleen Muenzen, Thomas A Peterson, Xuancheng Fan, Benjamin S Glicksberg, Gundolf Schenk, Eugenia Rutenberg, Boris Oskotsky, Marina Sirota, Jinoos Yazdany, et al. Protected health information filter (philter): accurately and securely de-identifying free-text clinical notes. *NPJ digital medicine*, 3(1):1–8, 2020.
- [91] Nicholas M Pace and Julia Pollak. Provider fraud in california workers’ compensation: Selected issues, 2017. https://www.rand.org/content/dam/rand/pubs/research_reports/RR1700/RR1703/RAND_RR1703.pdf.
- [92] Jonas Pfeiffer, Naman Goyal, Xi Lin, Xian Li, James Cross, Sebastian Riedel, and Mikel Artetxe. Lifting the curse of multilinguality by pre-training modular transformers. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 3479–3495, Seattle, United States, July 2022. Association for Computational Linguistics.
- [93] Prakash Poudyal, Jaromír Šavelka, Aagje Ieven, Marie Francine Moens, Teresa Gonçalves, and Paulo Quaresma. Echr: legal corpus for argument mining. In *Proceedings of the 7th Workshop on Argument Mining*, pages 67–75, 2020.
- [94] Praatika Prasad. Implicit racial biases in prosecutorial summations: Proposing an integrated response. *Fordham L. Rev.*, 86:3091, 2017.
- [95] Reid Pryzant, Kelly Shen, Dan Jurafsky, and Stefan Wagner. Deconfounded lexicon induction for interpretable social science. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1615–1625, 2018.
- [96] National Historical Publications and Records Commission. Founders online, 2010. <https://founders.archives.gov/>.
- [97] Jack W. Rae, Sebastian Borgeaud, Trevor Cai, Katie Millican, Jordan Hoffmann, Francis Song, John Aslanides, Sarah Henderson, Roman Ring, Susannah Young, Eliza Rutherford, Tom Hennigan, Jacob Menick, Albin Cassirer, Richard Powell, George van den Driessche, Lisa Anne Hendricks, Maribeth Rauh, Po-Sen Huang, Amelia Glaese, Johannes Welbl, Sumanth Dathathri, Saffron Huang, Jonathan Uesato, John Mellor, Irina Higgins, Antonia Creswell, Nat McAleese, Amy Wu, Erich Elsen, Siddhant Jayakumar, Elena Buchatskaya, David Budden, Esme Sutherland, Karen Simonyan, Michela Paganini, Laurent Sifre, Lena Martens, Xiang Lorraine Li, Adhiguna Kuncoro, Aida Nematzadeh, Elena Gribovskaya, Domenic Donato, Angeliki Lazaridou, Arthur Mensch, Jean-Baptiste Lespiau, Maria Tsim-poukelli, Nikolai Grigorev, Doug Fritz, Thibault Sottiaux, Mantas Pajarskas, Toby Pohlen, Zhitao Gong, Daniel Toyama, Cyprien de Masson d’Autume, Yujia Li, Tayfun Terzi, Vladimir Mikulik, Igor Babuschkin, Aidan Clark, Diego de Las Casas, Aurelia Guy, Chris Jones, James Bradbury, Matthew Johnson, Blake Hechtman, Laura Weidinger, Iason Gabriel, William Isaac,

- Ed Lockhart, Simon Osindero, Laura Rimell, Chris Dyer, Oriol Vinyals, Kareem Ayoub, Jeff Stanway, Lorraine Bennett, Demis Hassabis, Koray Kavukcuoglu, and Geoffrey Irving. Scaling language models: Methods, analysis & insights from training gopher. *arXiv preprint arXiv:2112.11446*, 2021.
- [98] Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J Liu. Exploring the limits of transfer learning with a unified text-to-text transformer. *arXiv preprint arXiv:1910.10683*, 2019.
- [99] Abigail M Reecer. The ethical dilemmas of the office of legal counsel in the wake of a whistleblower complaint. *Geo. J. Legal Ethics*, 33:769, 2020.
- [100] Douglas Rice, Jesse H Rhodes, and Tatishe Nteta. Racial bias in legal language. *Research & Politics*, 6(2):2053168019848930, 2019.
- [101] Federico Ruggeri, Francesca Lagioia, Marco Lippi, and Paolo Torroni. Detecting and explaining unfairness in consumer contracts through memory networks. *Artificial Intelligence and Law*, pages 1–34, 2021.
- [102] Steven Ruggles, Catherine Fitch, Diana Magnuson, and Jonathan Schroeder. Differential privacy and census data: Implications for social and economic research. In *AEA papers and proceedings*, volume 109, pages 403–08, 2019.
- [103] Julian Salazar, Davis Liang, Toan Q Nguyen, and Katrin Kirchhoff. Masked language model scoring. *arXiv preprint arXiv:1910.14659*, 2019.
- [104] Nithya Sambasivan, Shivani Kapania, Hannah Highfill, Diana Akrong, Praveen Paritosh, and Lora M Aroyo. “everyone wants to do the model work, not the data work”: Data cascades in high-stakes ai. In *Proceedings of the 2021 CHI Conference on Human Factors in Computing Systems*, pages 1–15, 2021.
- [105] Victor Sanh, Lysandre Debut, Julien Chaumond, and Thomas Wolf. Distilbert, a distilled version of bert: smaller, faster, cheaper and lighter. *arXiv preprint arXiv:1910.01108*, 2019.
- [106] Arman Sarvarian. Common ethical standards for counsel before the european court of justice and european court of human rights. *European Journal of International Law*, 23(4):991–1014, 2012.
- [107] Eran Shalev. Ancient masks, american fathers: Classical pseudonyms during the american revolution and early republic. *Journal of the Early Republic*, 23(2):151–172, 2003.
- [108] Reza Shokri and Vitaly Shmatikov. Privacy-preserving deep learning. In *Proceedings of the 22nd ACM SIGSAC Conference on Computer and Communications Security*, pages 1310–1321, 2015.
- [109] Amy Sholsberg, Evan Mandery, and Valerie West. The expungement myth. *Albany Law Review*, 75:1229–1242, 2011.
- [110] Social Security Administration. HIPAA and the Social Security Disability Programs. <https://www.ssa.gov/disability/professionals/hipaa-cefactsheet.htm>.
- [111] Congzheng Song and Ananth Raghunathan. *Information Leakage in Embedding Models*, page 377–390. Association for Computing Machinery, New York, NY, USA, 2020. ISBN 9781450370899.
- [112] Harold J Spaeth, Lee Epstein, Andrew D Martin, Jeffrey A Segal, Theodore J Ruger, and Sara C Benesh. Supreme court database, version 2013 release 01. *Database at http://supremecourtdatabase.org*, 2013.
- [113] Mike Spector, Jaimi Dowdell, and Benjamin Lesser. How secrecy in U.S. courts hobbles the regulators meant to protect the public. *Reuters*, 2010. <https://www.reuters.com/investigates/special-report/usa-courts-secrecy-regulators>.
- [114] Dan Sperber and Deirdre Wilson. Irony and the use-mention distinction. *Philosophy*, 3: 143–184, 1981.

- [115] Supreme Court of Canada. *Sherman Estate v. Donovan*, 2021. <https://www.canlii.org/en/ca/scc/doc/2021/2021scc25/2021scc25.html>.
- [116] The President and Fellows of Harvard University. Caselaw access project. <https://case.law/api/>.
- [117] Mark V Tushnet. *The Warren Court in historical and political perspective*. University of Virginia Press, 1993.
- [118] U.K. Judicial Conduct Investigations Office. *Annual Report 2015–2016*, 2016. https://jciodev.microsoftcrmpartals.com/_entity/annotation/61785fbb-752a-eb11-a813-000d3a0bacd3.
- [119] United Kingdom Bar Standards Board. *Bar Standards Board Handbook*. <https://www.barstandardsboard.org.uk/the-bsb-handbook.html?part=E3FF76D3-9538-4B97-94C02111664E5709&audience=&q=>.
- [120] United Nations Office on Drugs and Crime. *Commentary on the Bangalore Principles of Judicial Conduct*. 2007. https://www.unodc.org/documents/nigeria/publications/Otherpublications/Commentry_on_the_Bangalore_principles_of_Judicial_Conduct.pdf.
- [121] United States Court of Appeals for the Second Circuit. *Sealed Plaintiff v. Sealed Defendant*, 2008.
- [122] U.S. Election Assistance Commission. Availability of state voter file and confidential information, 2020. https://www.eac.gov/sites/default/files/voters/Available_Voter_File_Information.pdf.
- [123] Betty van Aken, Julian Risch, Ralf Krestel, and Alexander Löser. Challenges for toxic comment classification: An in-depth error analysis. In *Proceedings of the 2nd Workshop on Abusive Language Online (ALW2)*, pages 33–42, 2018.
- [124] Eugene Volokh. Pseudonymous litigation, 2021. <https://www2.law.ucla.edu/volokh/pseudonym.pdf>.
- [125] Guillaume Wenzek, Marie-Anne Lachaux, Alexis Conneau, Vishrav Chaudhary, Francisco Guzmán, Armand Joulin, and Édouard Grave. CCNet: Extracting High Quality Monolingual Datasets from Web Crawl Data. In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 4003–4012, 2020.
- [126] Wenjie Yin and Arkaitz Zubiaga. Towards generalisable hate speech detection: a review on obstacles and solutions. *PeerJ Computer Science*, 7:e598, 2021.
- [127] Xiang Yue, Minxin Du, Tianhao Wang, Yaliang Li, Huan Sun, and Sherman SM Chow. Differential privacy for text analytics via natural text sanitization. In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3853–3866, 2021.
- [128] Susan Zhang, Stephen Roller, Naman Goyal, Mikel Artetxe, Moya Chen, Shuohui Chen, Christopher Dewan, Mona Diab, Xian Li, Xi Victoria Lin, et al. OPT: Open Pre-trained Transformer Language Models. *arXiv preprint arXiv:2205.01068*, 2022.
- [129] Lucia Zheng, Neel Guha, Brandon R Anderson, Peter Henderson, and Daniel E Ho. When does pretraining help? assessing self-supervised learning for law and the CaseHOLD dataset of 53,000+ legal holdings. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, pages 159–168, 2021.
- [130] Victor Zhou. profanity-check. Github. <https://github.com/vzhou842/profanity-check>, 2019.

Checklist

1. For all authors...
 - (a) Do the main claims made in the abstract and introduction accurately reflect the paper’s contributions and scope? [\[Yes\]](#)
 - (b) Did you describe the limitations of your work? [\[Yes\]](#)
 - (c) Did you discuss any potential negative societal impacts of your work? [\[Yes\]](#) Please see our ethics statement.
 - (d) Have you read the ethics review guidelines and ensured that your paper conforms to them? [\[Yes\]](#) Please see Appendix A, where we provide a point-by-point discussion of how our paper conforms to the 2022 NeurIPS ethics guidelines.
2. If you are including theoretical results...
 - (a) Did you state the full set of assumptions of all theoretical results? [\[N/A\]](#)
 - (b) Did you include complete proofs of all theoretical results? [\[N/A\]](#)
3. If you ran experiments...
 - (a) Did you include the code, data, and instructions needed to reproduce the main experimental results (either in the supplemental material or as a URL)? [\[Yes\]](#) See Appendix D.
 - (b) Did you specify all the training details (e.g., data splits, hyperparameters, how they were chosen)? [\[Yes\]](#)
 - (c) Did you report error bars (e.g., with respect to the random seed after running experiments multiple times)? [\[Yes\]](#)
 - (d) Did you include the total amount of compute and the type of resources used (e.g., type of GPUs, internal cluster, or cloud provider)? [\[Yes\]](#) See Appendix C.
4. If you are using existing assets (e.g., code, data, models) or curating/releasing new assets...
 - (a) If your work uses existing assets, did you cite the creators? [\[Yes\]](#)
 - (b) Did you mention the license of the assets? [\[Yes\]](#)
 - (c) Did you include any new assets either in the supplemental material or as a URL? [\[Yes\]](#)
 - (d) Did you discuss whether and how consent was obtained from people whose data you’re using/curating? [\[Yes\]](#)
 - (e) Did you discuss whether the data you are using/curating contains personally identifiable information or offensive content? [\[Yes\]](#) Please see Appendix A, where we provide a point-by-point discussion of how our paper conforms to the 2022 NeurIPS ethics guidelines.
5. If you used crowdsourcing or conducted research with human subjects...
 - (a) Did you include the full text of instructions given to participants and screenshots, if applicable? [\[N/A\]](#)
 - (b) Did you describe any potential participant risks, with links to Institutional Review Board (IRB) approvals, if applicable? [\[N/A\]](#)
 - (c) Did you include the estimated hourly wage paid to participants and the total amount spent on participant compensation? [\[N/A\]](#)

A Ethics Statement

While we recognize that the NeurIPS Ethical Guidelines are not an exhaustive list, much of our paper concerns the broader ethical questions surrounding the privacy and toxicity harms of our data. Accordingly, we limit this Appendix to specifically addressing the nine points explicitly mentioned in the 2022 version of the Guidelines.

1. Does the data contain any personally identifiable information or sensitive personally identifiable information?

The data contains PII. As we argue in the main text, each document presumptively follows the privacy norms of the jurisdiction where it was written, which necessarily means that some documents include PII that would violate the norms of a different jurisdiction (e.g., the inclusion of names in U.S. criminal cases would violate German privacy rules). Nonetheless, as we argue in the paper, a reasonable minimum standard for filtering is deference to privacy rules of a particular jurisdiction that already weighs transparency benefits against privacy harms. That is especially true given that our data is largely produced by governmental entities, not merely sanctioned by them. These rules are developed to balance privacy against transparency in a manner specific to cultural context and we respect these standards.

That said, while others have suggested different data governance strategies [58], we follow the approach of CourtListener, which has already grappled with the exact trade-off of compiling legal data [77]. CourtListener runs a filtering mechanism to remove information like SSNs that may have slipped through the courts' first pass redaction layer. CourtListener also provides a mechanism for stakeholders to take down cases. We also run a filter to validate that no Social Security Numbers (SSNs) were present in the data using Microsoft Presidio [83]. Like them, we check for high-risk information like SSNs and restrict indexing of the data by search engines to the best of our ability. We also provide instructions for requesting content removal on the dataset website that reflect the CourtListener mechanism. We also have enabled the HuggingFace community feature to allow requests for dataset changes and removal of content.

We do not redact information further for several reasons. Some of the data subsets, in particular the court subsets, request that reproduction "exercise due diligence in ensuring the accuracy and currency of the materials reproduced" (see, e.g., <https://www.bccourts.ca/Privacy%20Statement.html>). By using existing anonymization filters that might noisily redact factual information or legal citations, we cannot ensure such accuracy and concurrency. As such, we chose to respect the data sources' decisions on the question of balancing privacy against transparency and to reproduce the content in the most accurate way possible. Furthermore, removing names for common law data would break the important legal references and context if done in a non-contextual manner. For example, a quote like "In *Brown*, Justice Warren wrote..." would become "In [MASK], Justice [MASK] wrote...", which would not provide any information about the underlying case law nor the facts of the case to which the law should apply.

Since the goal of the paper is to provide mechanisms and data for "doing the data work," we largely leave the data unfiltered beyond the mechanisms described above. Because we may be bound by upstream licenses prohibiting further restrictions on use, we place the data compilation under a Creative Commons Attribution-NonCommercial-ShareAlike 4.0 International license, but underlying data subsets may be bound by other licenses (see Table 5).

2. Does the data contain information that could be deduced about individuals that they have not consented to share?

It is possible that the case details included in sealed, pseudonymized, or otherwise sanitized case materials could be used to deduce the identity of litigants. Nonetheless, as we note above and in the main text, the degree of protection afforded to litigants in these cases is precisely the protection to which they are entitled by law. Furthermore, litigants have ample notice about the degree of protection afforded to their identities in public cases and mechanisms to challenge decisions on anonymity. We acknowledge that there may be shortcomings to this process, however. Thus, following the mechanism put forth by CourtListener and others, we provide a mechanism via the HuggingFace community feature to request the removal of content. In this process, we ask that stakeholders file a takedown request with the upstream data source and notify us of its modification in that source. Then we will update our archive with the upstream. This ensures that we comply with upstream

licenses and takedown standards. In exceptional cases where the upstream source does not take down especially harmful data, we will consider such situations on a case-by-case basis. We also will prevent search engine indexing of the data to the best of our ability.

3. Does the data encode, contain, or potentially exacerbate bias against people of a certain gender, race, sexuality, or who have other protected characteristics?

As we note in the main text, some of the court opinions we include may rely on precedential cases that relied upon past legal regimes including racial or other discrimination. Further, there is empirical evidence of implicit bias among judges [73]. However, as we argue in the paper, there is reason to believe that explicitly derogatory expressions directed at particular vulnerable groups are absent from recent legal text.

Further, there is a high degree of public interest in the materials included in our dataset, because they reflect the public decisions of governmental institutions and may continue to bind future parties. For example, civil rights cases we reference (e.g., *Brown, Loving*) are crucial civil rights law that legal models must understand despite (or perhaps because of) their references to unjust legal regimes. And since we argue that a large benefit of our data is that it provides a new way to examine and address shortcomings in toxicity filters, we leave content unfiltered.

4. Does the paper contain human subject experimentation and whether it has been reviewed and approved by a relevant oversight board?

No.

5. Does the paper rely on data that have been discredited by the creators?

No.

6. Consent to use or share the data. Explain whether you have asked the data owner’s permission to use or share data and what the outcome was.

Consent for the curated data is implicit in each subset, in particular as indicated by the license it was released, or the legal standard under which it was made public. For example, parties to lawsuits in U.S. federal court have ample notice that their names and case details will become part of the public record, since that is encoded in the rules we reference above.

7. Domain specific considerations when working with high-risk groups.

As we note in the text of the paper, while high-risk groups appear in this data, we rely upon the filtering mechanisms afforded for such circumstances by the government publishers of the text. While our work may make data about such groups easier to access, it is not obvious that this will cause harm given the significant positive implications of improved access to legal information.

8. Filtering of offensive content. For instance, when collecting a dataset, how are the authors filtering offensive content such as racist language or violent imagery?

There may be racist or toxic language in the dataset, but as we mention in this work, many of these instances arise when a court is describing the facts of the case as opposed to a direct use of epithets. To be sure, as we note in the main text, historical precedents and cases relying on them are more likely to include toxic material given the legal system’s historical role in enforcing discriminatory laws. Further, there is empirical evidence of implicit bias among judges [73]. And since we argue that a large benefit of our data is that it provides a new way to examine and address shortcomings in toxicity and privacy filters, we leave content unfiltered so that contextual filters can be learned from the data.

We also note that there may well be violent content particularly in discussion of criminal cases, but for legal language models it is often important to understand the facts of a case and apply the law to these facts. As we argue in this work, one path forward is to learn context-dependent filters based on the court’s protection of victims that would retain content relevant to the law while filtering out identifiable information for protected individuals.

9. Compliance with GDPR and other data-related regulations. For instance, if the authors collect human-derived data, what is the mechanism to guarantee individuals’ right to be forgotten (removed from the dataset)?

As noted in the paper, our data is derived from governmental sources that impliedly comply with the privacy rules applicable in their respective jurisdictions. GDPR applies only to persons “monitoring the behavior of individuals in the EU.” The European data we collect that concerns ‘individuals’ other than Parliamentarians comes from the European Court of Human Rights which must already have complied with relevant provisions. HIPAA does not apply to any of the data we collect. For example, as we note in the main text, although many of the administrative cases we collect concern medical records, the organizations issuing those decisions are not HIPAA custodians and their decisions are not subject to HIPAA protection [36]. Nonetheless, we provide a mechanism to take down content as noted above.

B Limitations

Due to licensing restrictions, our collation of judicial texts from the Anglophone countries is concentrated on U.S. texts (though it does include a number of international sources, including European Parliament proceedings and the decisions of Canadian courts). For example, BAILII, the online repository of judgments issued by British and Irish courts, forbids “storing search results or HTML versions of judgments.”⁸ We encourage future efforts to identify freely-available versions of legal English-language texts in other countries to improve the geographic coverage of the data.

Furthermore, our experiments are meant to be demonstrative and set a research agenda. The learned filters we describe in Section 4 will likely require more work to be applied to significantly out of domain data (though we show that they are somewhat robust through a perturbation analysis). We believe this is an exciting new research pathway.

⁸See <https://www.bailii.org/bailii/copyright.html>

C Carbon Impact Statement

As suggested by Henderson et al. [54], Lacoste et al. [66] and others, we report the energy and carbon impacts of our experiments. While we are unable to calculate precise carbon emissions from hardware counters, we give an estimate of our carbon emissions. For the BERT models, we estimate roughly 8 weeks of RDU usage total at full capacity, including debugging, hyperparameter optimization, iteration on experiments on 8 RDUs. We are unsure of the exact TDP for these chips, but as rough estimate assume that the TDP of each is similar to an A100 GPU or a TPUv2. Assuming minimal CPU time, this is roughly 4704 kWh of energy used and 1121 kg CO_{2eq} at the California yearly average carbon intensity of 238.4 g CO_{2eq} / kWh. Other experiments took relatively smaller amounts of compute using various CPUs and GPUs, mostly at inference time. Fine-tuning of the BERT-Large model on CaseHOLD used 4 A100s for 2 days. Other experiments ran inference on a variety of GPUs/CPUs. To put a conservative estimate on this, we assume roughly 4 weeks of total added compute time at full capacity. Assuming an average TDP of 300W across these machines, this yields 201.6 kWh and 480.6 kg CO_{2eq}. We emphasize that these are extremely rough estimates (and may be overestimates due to their assumptions) and they take into account all experimentation, including iterations and debugging.

D Code and Data Availability

Pile of law is hosted at <https://huggingface.co/datasets/pile-of-law/pile-of-law>. The main pol-bert checkpoint we use in the paper is available at <https://huggingface.co/pile-of-law/legalbert-large-1.7M-2/> and an additional experimental checkpoint using a different random seed yet otherwise identical settings is at <https://huggingface.co/pile-of-law/legalbert-large-1.7M-1/>. We will also make intermediate checkpoints (every 50k timesteps) available on request. The EOIR experiment data and trained model are available at https://huggingface.co/datasets/pile-of-law/eoir_privacy and https://huggingface.co/pile-of-law/distilbert-base-uncased-finetuned-eoir_privacy, respectively.

Code for experiments and data collection will be available at <https://github.com/Breakend/Pileoflaw>. We include pre-processing code for the data used for pre-training, but do not provide the pre-training code itself since we directly use the out-of-the-box SambaNova pre-training tool with no changes of our own (other than at pre-processing time).

E File of Law Data Description

We briefly describe the data curated for each subsection, grouped by logical similarity.

E.1 Legal Case Opinions and Filings

CourtListener Opinions, CourtListener Docket Entries and Court Filings. CourtListener provides a large set of U.S. court case opinions across a number of federal and state courts. Our dataset includes all *judicial opinions* in the CourtListener opinions dataset. This is similar to the FreeLaw portion of [44]. *Judicial opinions* are long-format documents written by judges explaining and justifying a decision about the factual or legal disposition of a case. Judicial opinions are typically written in a more authoritative and less argumentative style.

To supplement the data with argumentative language, we scrape CourtListener’s RECAP Docket Entry Documents API from present day until 2018 (due to rate limits we limit to the past few years; future iterations of the dataset may gather more documents). RECAP includes *briefs*, which are a party’s arguments to the court, as well as interim judicial opinions, exhibits (i.e., the evidence supporting a party’s claim), and other miscellaneous case records. RECAP docket entry documents are any publicly available filings provided on a court docket. Dockets do not generally include evidence that is not referenced by either party.

U.S. Supreme Court Docket Entries and Court Filings. The U.S. Supreme Court docket system contains information about the status of pending and decided cases that have been filed at the Court. A docket is a log containing the history of each case in the form of brief chronological entries summarizing the court proceedings and the court filings, the underlying documents (pleadings, motions, briefs, etc.) filed with the Court in the proceedings of a case. Though dockets and court filings do not have precedential value, the information contained can sometimes provide additional insight into why the Court issued a particular decision or opinion.

U.S. Board of Veterans’ Appeals Decisions. The U.S. Board of Veterans’ Appeals (BVA) is an internal administrative agency of the U.S. government that hears appeals from veterans who were denied disability or other benefits by the Veterans’ Benefits Administration. Every single case involves a veteran arguing for more benefits, opposed by the government arguing for the status quo. The vast majority of BVA opinions concern either (A) the severity of a veteran’s disability or (B) the etiology of the disability and whether it can be traced to their period of service in the U.S. military [56].

U.S. Federal Trade Commission Advisory Opinions. When a business or trade group wishes to engage in a practice that may violate competition or consumer protection laws, it can request an official opinion from the U.S. Federal Trade Commission (FTC), which issues an *advisory opinion*, often written in a similar manner to a judicial opinion.

U.S. National Labor Relations Board Decisions. The U.S. National Labor Relations Board (NLRB) governs the relations between employers and unions, and when one party alleges that there has been a violation of the National Labor Relations Act (e.g., an unfairly held union election), the NLRB is tasked with deciding whether an unfair labor practice has occurred.

U.S. Department of Justice Executive Office for Immigration Review *Immigration & Nationality Decisions*. When a person disagrees with the decision of an immigration judge, they can appeal to the Board of Immigration Appeals (BIA) and/or directly to the Attorney-General (AG), who may also decide to review BIA decisions independently [52]. We download decisions of the BIA and the AG included in the DOJ’s *Immigration & Nationality Decisions*.

U.S. Tax Court PLR Corpus. The U.S. Tax Court resolves disputes between taxpayers and the Internal Revenue Service (IRS). We include the data released by [14], which includes Tax Court memorandum and summary opinions scraped from the Tax Court website and IRS Private Letter Rulings (PLRs), scraped from the IRS website. PLRs are written statements, issued at the request of a taxpayer, that interprets and applies tax laws to the taxpayer’s represented set of facts.

U.S. Department of Labor Employees’ Compensation Appeals Board Decisions. When a person disagrees with the decision of the U.S. Department of Labor Office of Workers’ Compensation Programs (OWCP), they can appeal to the Employees’ Compensation Appeals Board (ECAB). We download decisions of the ECAB from 2006 to April 2022.

European Court of Human Rights Opinions. The European Court of Human Rights (ECHR) hears appeals by individuals or states bound by the European Convention on Human Rights, and issues precedential opinions binding state parties. We include the ECHR corpus, which contains 42 decisions of the ECHR and was created and introduced by [93].

Canadian Court Opinions. The Ontario Court of Appeals (ONCA) and the British Columbia Court of Appeals (BCCA) are intermediate appeals courts that hear all manner of cases originating in their respective jurisdictions. We download the available opinions listed on the ONCA and BCCA websites.

E.2 Legal Analyses

U.S. Office of Legal Counsel Memos. The U.S. Office of Legal Counsel (OLC) advises the U.S. President on the legality of any action the president wishes to take. OLC issues memos that provide in-depth legal analyses on various topics that have confronted presidents in the past. These should not be considered ground truth for the law, but provide good examples of in-depth legal analysis and are similar to judicial opinions in tone.

U.S. Department of Justice Inspector General Reports. Almost every federal agency in the U.S. federal government has an Office of the Inspector General, which is responsible for independent oversight of the agency. This includes regular audits of the agency's spending, monitoring of active government contractors, and investigations into wasteful or corrupt agency practices, which are compiled and published in public reports. We download the reports of every U.S. IG published online, which were originally scraped by [84].

E.3 Laws

U.S. Code of Federal Regulations, U.S. State Codes, U.S. Code, U.S. Federal Rules of Evidence, U.S. Federal Rules of Civil Procedure. We scrape the U.S. Code of Federal Regulations, U.S. State Codes, U.S. Code, Federal Rules of Evidence, and U.S. Federal Rules of Civil Procedure. These are statutes and rules that govern much of the U.S. legal system and are important knowledge for machine learning algorithms.

U.S. Bills, U.S. Federal Register. We scrape proposed U.S. Bills and the U.S. Federal Register, which includes proposed and final rules along with other regulatory actions, to imbue understanding of statutory construction (Bills) and regulatory construction (Federal Register). We use the Government Publishing Office's govinfo API and Bulk Data Repository [61] to scrape proposed U.S. Bills and the U.S. Federal Register. The Bulk Data Repository provides digitized proposed U.S. Bills from the 113th Congress (2013-2014) to the 117th Congress (2021-2022) and proposed Federal Register rules from 2000-2022. We scrape from the beginning of the available time range to October 2021.

U.S. Founders Letters. The U.S. founders letters describe the birth of the American Republic and its democratic and legal institutions. We scrape ~185,000 documents of correspondences and other writings by U.S. founders, George Washington, Benjamin Franklin, John Adams (and family), Thomas Jefferson, Alexander Hamilton, John Jay, and James Madison, from Founders Online [96].

World Constitutions. Constitutional construction is an important task for state-building – as well as for fundamental legal analysis. As such, we add the world's constitutions [41] to the data.

EUR-Lex. EUR-Lex provides access to European Union (EU) legal documents, including treaties, legal acts from EU institutions, preparatory documents, EU case law, EU international agreements, European Free Trade Association (EFTA) documents, national transposition measures, and national case law related to EU law. EUR-Lex provides greater context on legal systems in the European Union and international law. We add EUR-Lex data from [28], an expanded version of the Eur-Lex dataset released by [79].

E.4 Contracts / Business Documents

Credit Card Agreements, Terms of Service, Edgar Contracts, Atticus Contracts. Credit Card Agreements provided by the U.S. Consumer Financial Protection Bureau,⁹ Terms of Service [76, 101],

⁹<https://www.consumerfinance.gov/credit-cards/agreements/>

Edgar Contracts [17], Atticus Contracts [55] were all added to the data, which may be useful pretraining data for a large number of contract based tasks that are beginning to employ machine learning methods, such as contract review [71, 53, 55].

E.5 Conversations

U.S. Congressional Hearings. A congressional hearing is a meeting or a session of a Senate, House, joint, or special committee of Congress, to obtain information and opinions on proposed legislation, conduct an investigation, evaluate the activities of a government department or the implementation of a federal law, or provide testimony and data about topics of current interest. These hearings provide details on the research and drafting process for proposed legislation. The transcripts of congressional hearings from the 89th Congress (1965-67) to the 117th Congress (2021-2022) were scraped using the Government Publishing Office’s govinfo API and Bulk Data Repository [61].

European Parliament Proceedings Parallel Corpus. The European Parliament proceedings capture debate about proposed EU legislation. We include only the English data in the European Parliament Proceedings Parallel Corpus [64], though the original data provides parallel translations in 11 languages.

U.S. Supreme Court Oral Argument Transcripts. The U.S. Supreme Court holds oral arguments in about 70-80 cases each year. Oral arguments give the Justices an opportunity to ask questions to the attorneys representing the parties to the case and the attorneys to highlight important arguments of the case. We extract oral argument transcript data using raw data made available through [7], from the 2021-08-14 release.

U.N. General Debate Corpus. The General Debate is held at the beginning of each session of the United Nations (UN) General Assembly, which has convened annually since 1946. The General Debate is a forum for world leaders and other senior officials representing UN member states to deliver statements that present their government’s perspective on the major issues in world politics, analogous to legislative state-of-the-union addresses in domestic politics. We include the General Debate statements from the UN General Debate Corpus (UNGDC) [8] in the Pile of Law. The UNGDC contains General Debate statements from 1970 (Session 25) to 2020 (Session 75).

Reddit r/legaladvice & r/legaladviceofftopic. Because most legal language is often difficult to understand for the lay person and does not encode clear answers to simple legal questions, we sought to find a dataset that yields a “plain English” Q&A format. We settled on the two subreddits r/legaladvice and r/legaladviceofftopic. Because of the risk of encoding incorrect legal advice, we heavily filtered the data. We filtered out any posts with profanity using profanity-check [130]. We also only included posts with at least one answer with a score of over 8 net upvotes. We then restructured the data as:

Title: [Post Title]
Question: [Post Content]
Topic: [Post Flair]
Answer #[N]: [Top Answers]...

We used the PushShift API to scrape the entirety of the each subreddit [9].

E.6 Study Materials

Bar Exam Outlines. The bar exam outlines provide key definitions and descriptions of concepts relevant to various subject areas in American law that are tested on the bar examination (e.g., Constitutional Law, Contracts, Criminal Law, etc.). In every U.S. jurisdiction (except in certain cases in Wisconsin), all applicants seeking admission to the bar to practice law in the jurisdiction must pass a bar examination.

Open Source Casebooks. Legal knowledge is hard to parse from individual cases. Casebooks and textbooks have been created to teach students with edits cases and commentary to focus learning on the most important legal topics. As such we gather ~60 casebooks that were released under a Creative Commons license. The casebooks range across a broad range of topics corresponding to core doctrinal material. All licenses and author credits remain self-contained in the individual documents.

Table 4: Description of the Pile of Law by Data Source

Data Source	Data Size	Word Count	Document Count
Court Listener Opinions	59.29GB/19.76GB	7.65B/2.55B	3.39M/1.12M
Court Listener Docket Entries and Court Filings	52.13GB/17.38GB	5.36B/1.79B	1.49M/496K
U.S. Supreme Docket Entries and Court Filings	1.51GB/0.50GB	151.05M/51.73M	48K/16K
U.S. Board of Veterans' Appeals Decisions	13.21GB/4.40GB	1.74B/580.98M	630K/210K
U.S. Federal Trade Commission Advisory Opinions	1.55MB/0.52MB	157K/53K	112/33
U.S. National Labor Relations Board Decisions	994.83MB/331.61MB	120.33M/39.20M	24K/8K
U.S. Department of Justice Executive Office for Immigration Review <i>Immigration & Nationality Decisions</i>	22.89MB/7.63MB	3.05M/1.01M	1671/558
U.S. Department of Labor Employees' Compensation Appeals Board	353.25MB/117.75MB	48.20M/16.01M	21K/7K
European Court of Human Rights Opinions [93]	111.53MB/37.18MB	16.71M/3.47M	7K/1K
Canadian Court Opinions (ON, BC)	182.09MB/60.70MB	23.45M/7.66M	8K/3K
U.S. Office of Legal Counsel Memos	36.98MB/12.33MB	4.36M/1.31M	1038/346
U.S. Office of Inspector General Reports	1.90GB/0.63GB	167.71M/54.18M	29K/10K
U.S. Code of Federal Regulations	670.87MB/223.62MB	79.06M/25.41M	182/61
U.S. Supreme Court Oral Argument Transcripts	1.51GB/0.50GB	151.05M/51.73M	47K/16K
U.S. State Codes	6.77GB/2.26GB	829.62M/441.38M	157/60
U.S. Code	268.40MB/89.47MB	30.54M/18.20M	43/15
U.S. Federal Rules of Evidence	670KB/223KB	77K/36K	51/17
U.S. Federal Rules of Civil Procedure	1.59MB/0.53MB	237K/40K	69/23
U.S. Bills	1.27GB/0.42GB	156.06M/49.4M	84K/28K
U.S. Federal Register	159.29MB/53.10MB	6.61M/53.27M	4060/1354
U.S. Founders Letters	419.33MB/139.78MB	53.27M/17.69M	138K/46K
World Constitutions [41]	24.43MB/8.14MB	3.43M/1.06M	139/48
EUR-Lex [28]	1.31GB/0.44GB	191.65M/65.31M	106K/36K
Credit Card Agreements	70.19MB/23.40MB	10.73M/3.09M	2023/615
Terms of Service [76, 101]	1.57MB/0.52MB	213K/62K	37/13
Edgar Contracts [17]	10.76GB/3.59GB	1.44B/473.50M	741K/247K
Atticus Contracts [55]	31.2GB/10.4GB	3.96B/1.31B	488K/163K
U.S. Congressional Hearings	6.17GB/2.06GB	761.12M/250.04M	24K/8K
U.S. Tax Court PLR Corpus [14]	639.03MB/213.01MB	84.25M/27.62M	41K/14K
European Parliament Proceedings Parallel Corpus [64]	302.71MB/100.90MB	41.55M/13.70M	7K/2K
U.N. General Debate Corpus [8]	134.90MB/44.97MB	17.68M/5.81M	6K/2K
Reddit r/legaladvice & r/legaladviceofftopic	299.04MB/99.68MB	40.42M/13.56M	110K/37K
Bar Exam Outlines	1.18MB/0.39MB	123K/43K	44/15
Open Source Casebooks	87.09MB/29.03MB	9.20M/3.91M	52/14
Total	~ 256GB	~ 30B	~ 10M

Table 5: Filtering Norms by Data Source in the Pile of Law

Data Source	Examples of Filtering Norms
Court Listener Opinions	FRCP 49.1 (requiring partial redactions for social-security number and taxpayer-identification number, date of birth, minor's names, financial account numbers; governing sealing and redaction standards for other information that parties may wish to keep private); State Rules for filing pseudonymously ¹⁰ . Judicial code of ethics govern conduct of judges; American Bar Association Model Rules of Professional Conduct govern attorney conduct.
Court Listener Docket Entries and Court Filings	<i>Id.</i>
U.S. Supreme Docket Entries and Court Filings	<i>Id.</i>
U.S. Board of Veterans' Appeals Decisions	38 CFR 20.1301(c) ("Appeals on or after January 1, 1992, are electronically available for public inspection and copying on the Board's website. All personal identifiers are redacted from the decisions prior to publication.")

¹⁰<https://withoutmyconsent.org/50state/filing-pseudonymously/federal/>

U.S. Federal Trade Commission Advisory Opinions	16 CFR 1.4 (“Written advice rendered pursuant to this section and requests therefor, including names and details, will be placed in the Commission’s public record immediately after the requesting party has received the advice, subject to any limitations on public disclosure arising from statutory restrictions, the Commission’s rules, and the public interest. A request for confidential treatment of information submitted in connection with the questions should be made separately.”)
U.S. National Labor Relations Board Decisions	The U.S. National Labor Relations Board (NLRB) protects information in accordance with the Privacy Act of 1974, the E-Government Act of 2002, P.L. 107-347, and the Federal Records Act, 44 U.S.C. § 3301 et seq. Section 208 of the E-Government Act of 2002 requires all federal agencies to conduct a privacy impact assessment (PIA) for all new or substantially changed technology that collects, maintains, or disseminates personally identifiable information (PII). The goals of a PIA are to ensure conformance with applicable legal, regulatory, and policy requirements for privacy, determine privacy risks, and evaluate processes to mitigate potential privacy risks.
U.S. Department of Justice Executive Office for Immigration Review <i>Immigration & Nationality Decisions</i>	8 CFR 208.6 (describing restrictions on disclosure of information to third parties in relation to asylum claims); 8 CFR 103.10(d) (the Attorney General may select cases to publish as precedential decisions)
U.S. Department of Labor Employees’ Compensation Appeals Board	20 CFR 501.8(c) (decisions shall be publicly available); Agency Policy not governed by law ¹¹ (“To limit a claimant’s exposure on the Internet, the Department of Labor (DOL) avoids referring directly to the claimant’s name in decisions posted on the DOL web site on or after August 1, 2006. . . This policy is intended to protect FECA claimants from unnecessary publicity; it is not based upon a legal requirement. Neither FOIA, nor the Privacy Act, nor any other law compels DOL to take this action. Rather, this policy is based on a desire to address in a responsible way concerns raised by some claimants about the ease of access to their case decisions on the Internet. The policy seeks to comply with legal requirements to make agency decisions available on the Internet, but to do so in a way that limits a claimant’s exposure to Internet users.”)
European Court of Human Rights Opinions [93]	Although EHCR does not publish a formal set of PII practices, it is bound by its own interpretations of the European Convention on Human Rights, which includes among others the right to be forgotten (see <i>Hurbain v. Belgium</i>). ¹² Ordinarily, the Court appeals to publish the names of plaintiffs but to anonymize all other details. Rules 33 and 47 of the Rules of Court also allow for anonymity and takedown requests to be considered in the case of private data. There are no professional standards governing lawyers before the EHCR except for those imposed by the lawyer’s home country [106].

¹¹<https://www.dol.gov/agencies/ecab/decisions-info>

¹²<https://hudoc.echr.coe.int/fre?i=001-210884> (interpreting Articles 8 and 10 of the Convention to permit the censorship of documents that infringe on the right to be forgotten).

Canadian Court Opinions (ON, BC)	The names of litigants and all other case details are ordinarily public in Canadian court decisions, although in certain criminal cases pertaining to sexual offenses they may be sealed by the presiding judge (example). Under the Youth Criminal Justice Act, the records of criminal cases against juveniles are generally confidential. ¹³ Sealing standards in Canada, as in the U.S., permit restrictions on public access where necessary to protect important public interests [115]. Canadian lawyers and judges are bound by ethical norms that prohibit discourteous speech, which impliedly includes derogatory speech [68, 19].
U.S. Office of Legal Counsel Memos	Because the Office of Legal Counsel handles legal questions raised by government agencies about policy decisions, its opinions rarely pass on the particular details of a specific person's interaction with the government. ¹⁴
U.S. Office of Inspector General Reports	OIG offices are effectively the ethics investigators that examine the conduct of U.S. agencies (including compliance with privacy regulations), so it would be extremely unusual for OIG offices to undertake actions that would create privacy or toxicity effects. 5 CFR §2638.106.
U.S. Code of Federal Regulations	Like other prospective rules of general application, the CFR does not ordinarily address individuals and thus is highly unlikely to contain private information. It is possible that in the extremely unlikely event that PII was revealed in the CFR or Federal Register, the Privacy Act of 1974 might apply to provide a remedy. A very small number of regulations may contain archaic terms now considered toxic (see e.g. 10 C.F.R. §§ 800.003); otherwise no racial epithets are used.
U.S. State Codes	Codified statutes (i.e. laws passed by a legislature) generally do not contain private information since they announce prospective rules that do not apply to any individual. However, state codes may include archaic laws that feature speech that would currently be classified as toxic (e.g., Miss. Code Ann. § 37-113-31).
U.S. Code	Codified statutes (i.e. laws passed by a legislature) generally do not contain private information since they announce prospective rules that do not apply to any individual. Further, most archaic references to minority communities have now been expurgated and replaced (see e.g. H.R. 4238).
U.S. Federal Rules of Evidence	Codified rules generally do not contain private information since they announce prospective rules that do not apply to any individual. The FRE contains no overt epithets.
U.S. Federal Rules of Civil Procedure	Codified rules generally do not contain private information since they announce prospective rules that do not apply to any individual. The FRCP contains no overt epithets.

¹³<https://www.justice.gc.ca/eng/cj-jp/yj-jj/tools-outils/sheets-feuillets/recor-dossi.html>

¹⁴There are exceptions, see e.g. [99] (referring to a whistleblower complaint that was referred to OLC) but these are generally not made public.

U.S. Bills	As is true of the U.S. Code, proposed Bills generally do not contain private information, except in the rare cases where a bill is passed for the benefit of one individual (e.g., to grant a person citizenship); in such cases, a person has no expectation of privacy. Proposed bills from past eras may contain archaic derogatory language but recent legislation, which is edited and produced by professionalized offices, are unlikely to do so (House ; Senate).
U.S. Federal Register	The Federal Register contains only official communications about prospective or final rules, as well as comments submitted about those rules, other official agency action, and official actions by the President. The names and contact information that appear most frequently are individuals designated to receive public comments, whose information is personally identifiable but not sensitive in virtue of the offices they occupy. Certain executive orders may also name individuals (e.g. pardons for criminal cases). Agencies are extremely unlikely to reveal personal information outside of these contexts, though in such a highly unusual context the Privacy Act of 1974 might apply to provide a remedy.
U.S. Founders Letters	Although many letters tended to use pseudonyms [107], no formal rules applied to the inclusion of information in these documents.
World Constitutions [41]	No filtering standards, but many constitutions are heavily influenced by the Universal Declaration of Human Rights [67] and recent constitutions are unlikely to have offensive or private content.
EUR-Lex [28]	It is unclear if there are any standards governing European Law. However, as with U.S. legal sources it is unlikely that modern legal text has any toxic or private information since these sources promulgate laws. That being said older laws may have direct epithets and there may be indirectly toxic content in the laws.
Credit Card Agreements	Consumer Financial Protection Bureau Privacy Policy ¹⁵ ; 12 CFR 1070.13(d) (private information should be redacted by CFPB)
Terms of Service [76, 101]	Since Terms of Service are not personalized, there would be no reason for PII to appear in this content. Professional norms and rules govern the drafters (attorneys) of ToS agreements (see ABA Model Rules of Professional Conduct). Reputational harms and anti-discrimination laws may also constrain overtly toxic content.
Edgar Contracts [17]	SEC Policy ¹⁶ (“My name appears in an old enforcement order or release. Is it possible to remove the document so web searches on my name don’t return the sec.gov document at the top of the results list? We don’t remove historical enforcement materials at public request or attempt to influence search result rankings. Enforcement documents from the beginning of sec.gov in 1995 remain available.”)

¹⁵<https://www.consumerfinance.gov/privacy/>

¹⁶<https://www.sec.gov/os/webmaster-faq#reuse2>

Atticus Contracts [55]	Professional norms and rules of professional conduct would govern the drafters of contracts. However, dated, offensive, and discriminatory clauses still remain in some contracts, ¹⁷ though they don't seem to appear in this data.
U.S. Congressional Hearings	House of Representatives 117th Congress Rule VII, clause 3 (b)(2) (“An investigative record that contains personal data relating to a specific living person (the disclosure of which would be an unwarranted invasion of personal privacy), an administrative record relating to personnel, or a record relating to a hearing that was closed under clause 2(g)(2) of rule XI shall be made available if it has been in existence for 50 years.”); Rule X, clause 11 (f) (“The select committee shall formulate and carry out such rules and procedures as it considers necessary to prevent the disclosure, without the consent of each person concerned, of information in the possession of the select committee that unduly infringes on the privacy or that violates the constitutional rights of such person. Nothing herein shall be construed to prevent the select committee from publicly disclosing classified information in a case in which it determines that national interest in the disclosure of classified information clearly outweighs any infringement on the privacy of a person.”); Rule XVII, clause 4 and 8 (unparliamentary words may be preserved for the record and may be removed only by permission or order of the house) ¹⁸
U.S. Tax Court PLR Corpus [14]	Tax Court’s Rules of Practice and Procedure Rule 27 and 103 (governing privacy redactions and sealing); Internal Revenue Code Section 7461(b) (court can take action “which is necessary to prevent the disclosure of trade secrets or other confidential information, including [placing items] under seal to be opened only as directed by the court.”). Despite these rules, hearings in the past have had offensive or toxic content appear.
European Parliament Proceedings Parallel Corpus [64]	Rules of Parliament Title IX Rule 226.13 (“the petitioner, a co-petitioner or a supporter may request that his, her or its name be withheld in order to protect his, her or its privacy, in which case Parliament shall comply with the request.”)
U.S. Supreme Court Oral Argument Transcripts	Supreme Court Justices are not bound by a code of conduct ¹⁹ , but professional norms likely restrict their speech.
U.N. General Debate Corpus [8]	The corpus files were cleaned to remove typos and OCR conversion errors, but were not otherwise altered. But similar to other corpora of debates in political forums, the speech of U.N. General Debate speakers is likely restricted by professional norms.

¹⁷<https://www.mercurynews.com/2019/02/26/for-whites-only-shocking-language-found-in-property-docs-throughout-bay-area/>

¹⁸<https://sgp.fas.org/crs/misc/R45866.pdf>

¹⁹<https://sgp.fas.org/crs/misc/LSB10255.pdf>

Reddit r/legaladvice & r/legaladviceofftopic	Content Moderation Policies r/legaladvice ²⁰ (no identifying information, no illegal advice, etc.); Content Moderation Policies r/legaladviceofftopic ²¹ (“personal attacks and harassing comments will be removed”; “While it is okay to post published situations, disclosing the names or information of otherwise-anonymous parties, users, etc., is strictly forbidden.”)
Bar Exam Outlines	No restrictions are made on third party content, but given that attorneys wrote the outlines they would be governed by the same professional rules as any attorney.
Open Source Casebooks	No restrictions are made on third party content, but given that attorneys wrote the outlines they would be governed by the same professional rules as any attorney (especially because case books are meant to be educational content).

²⁰https://www.reddit.com/r/legaladvice/wiki/index#wiki_general_rules

²¹<https://www.reddit.com/r/legaladviceofftopic/>

Data Source	License
Court Listener Opinions	Underlying content is Public Domain.
Court Listener Docket Entries and Court Filings	Public Domain.
U.S. Supreme Docket Entries and Court Filings	Public Domain.
U.S. Board of Veterans' Appeals Decisions	Public Domain.
U.S. Federal Trade Commission Advisory Opinions	Public Domain.
U.S. National Labor Relations Board Decisions	Public Domain.
U.S. Department of Justice Executive Office for Immigration Review <i>Immigration & Nationality Decisions</i>	Public Domain.
U.S. Department of Labor Employees' Compensation Appeals Board	Public Domain.
European Court of Human Rights Opinions [93]	Non-commercial OK. Commercial use requires written permission.
Canadian Court Opinions (ON, BC)	All reproduction OK (ON , BC). We acknowledge that these documents are sourced from the Court of Appeal for Ontario and the Court of Appeal for British Columbia, respectively. Note that these are not official versions.
U.S. Office of Legal Counsel Memos	Public Domain.
U.S. Office of Inspector General Reports	Underlying content is Public Domain. Complication is CC0 1.0 Universal public domain dedication .
U.S. Code of Federal Regulations	Public Domain.
U.S. Supreme Court Oral Argument Transcripts	Public Domain.
U.S. State Codes	Public Domain.
U.S. Code	Public Domain.
U.S. Federal Rules of Evidence	Public Domain.
U.S. Federal Rules of Civil Procedure	Public Domain.
U.S. Bills	Public Domain.
U.S. Federal Register	Public Domain.
U.S. Founders Letters	Creative Commons CC0 1.0 Universal license
World Constitutions [41]	CC BY-NC 3.0
EUR-Lex [28]	Creative Commons Attribution 4.0 International
Credit Card Agreements	Provide by Consumer Financial Protection Bureau in the Public Domain, but original copyright and license is unknown. We assume publication is governed by fair use standards.
Terms of Service [76, 101]	Publicly Available, but unknown license. We assume publication is governed by fair use standards.
Edgar Contracts [17]	Unknown license. We assume publication is governed by fair use standards.
Atticus Contracts [55]	CC BY 4.0
U.S. Congressional Hearings	Public Domain
U.S. Tax Court PLR Corpus [14]	Underlying Content is Public Domain, Complication License is CC BY-NC 4.0
European Parliament Proceedings Parallel Corpus [64]	No copyright restrictions on compliation. Non-commercial for underlying data.
U.N. General Debate Corpus [8]	Public Domain
Reddit r/legaladvice & r/legaladviceofftopic	Creative Commons Attribution 4.0 International. Reddit also grants a license to copy and display the underlying data.
Bar Exam Outlines	Publicly Available, but unknown license. We assume publication is governed by fair use standards.
Open Source Casebooks	All CC, varying on exact restrictiveness. Most restrictive: CC BY-NC-SA 4.0 . All licensing information preserved in individual documents.

Table 6: Content licenses with links where applicable. All government-generated content in the United States is public domain.

E.7 Distribution of Topics

We train a TF-IDF SVM on the LexGlue Supreme Court topic modeling task [30, 112]. The task predicts 13 topics from case data. We then run this on our data to get a sense for topic distribution. This simple method aligns with expectations, as seen in Figure 3. NLRB opinions, for example, are mainly classified as “Unions” (as one would expect), constitutions are mainly concerned with civil rights, the federal rules of civil procedure are mostly related to judicial power, and SEC Edgar filings are mostly related to economic activity. However, there are some unintuitive categorizations that can be explained by the coding scheme of Spaeth et al. [112].

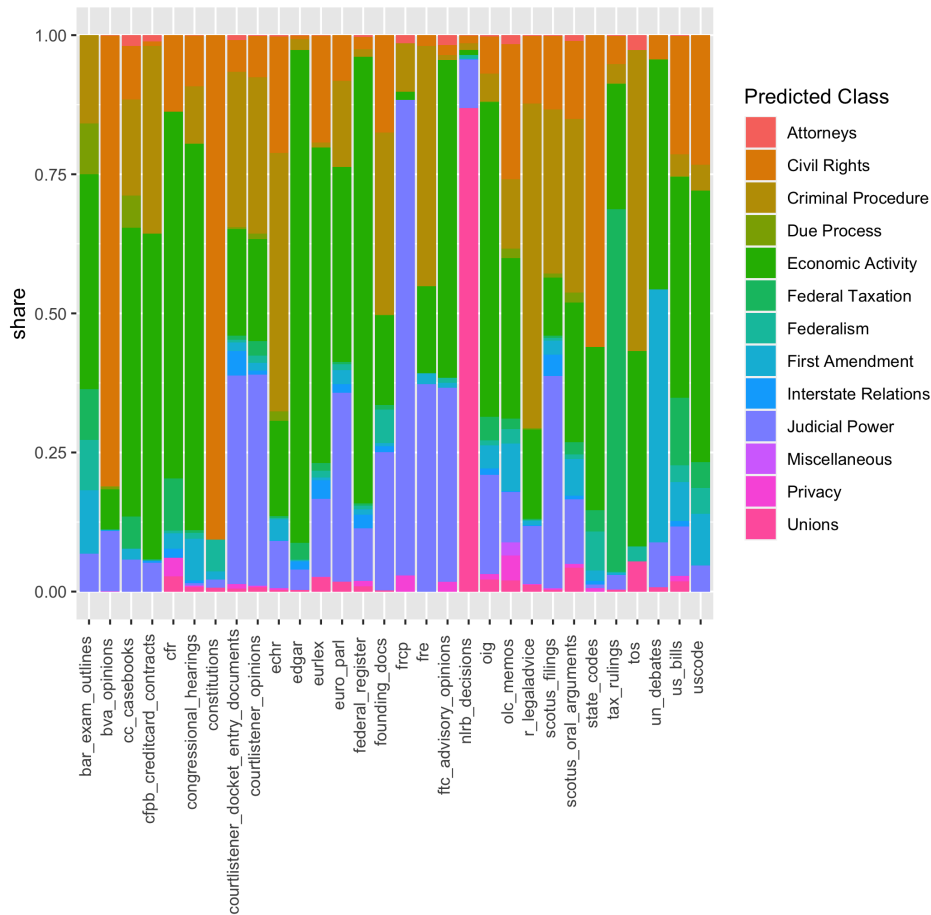


Figure 3: The distribution of topics in Pile of Law.

Model	CaseHOLD (F1)
CaseLaw-BERT (tuned)	78.5
CaseLaw-BERT [30, 129]	75.4
PoL-BERT-Large (tuned)	75.0
Bert-Large-Uncased (tuned)	71.3

Table 7: PoL-BERT-Large compared against results on the CaseHOLD [129] variant provided by [30].

F Models

We use the [SambaNova Systems Dataflow-as-a-ServiceTM](#) platform for training an initial experimental BERT-large baseline model from scratch. The API interface abstracts away training details for training transformer models and uses SambaNova RDU chips which are optimized for data-intensive workflows, like those required for training a BERT-large model. We implement a customized pre-processing method similar to Roberta pretraining (described below). Despite using a training method consistent with prior BERT-like model training protocols and similarly consistent pre-processing, we found training a model on such varied data surprisingly challenging, and detail some discussion below that may be a path forward for interesting pretraining research on Pile of Law.

First, we fit a custom word-piece vocabulary to the training split of Pile of Law consisting of 29k tokens using the HuggingFace WordPiece tokenizer.²² We supplement this with a random set of legal terms from Black’s Law Dictionary to make a total vocabulary size of 32k tokens.

We do not use the NSP task of BERT [37], but rather use the method of Roberta [78]. We train on 512 length sequences for the entirety of training and use the 80-10-10 masking, corruption, leave split of BERT [37]. We use a replication rate of 20 to create different masks for each context. To generate sequences we use the LexNLP sentence segmenter (which handles legal citations which are often falsely mistaken for sentences) [15]. We then fill sentences until they comprise 256 tokens. We add a [SEP] token after this, and then fill sentences such that the entire span is under 512 tokens. If the next sentence in the series is too large we do not add it and fill the context with padding tokens.

For pretraining, we create a smaller training set use a randomly generated ~ 30 GB sample of the Pile of Law sampled evenly across subsets of data. This means that some subsets of data were included in their entirety while others only included a small portion of the total training data.

At first, we randomly shuffled data, but found that training this way was quite unstable and we were unable to get any model to converge to a reasonable optimum. Instead, we had to distribute data sources evenly across shards such that each device saw the same mixture of data.

Even with this, we found that a large batch size and learning rate created large instabilities, potentially due to the diversity of the data. As a result, we used a very small learning rate (5e-6) and batch size (128) which yielded stable training. We ran this for roughly two weeks until 1.7M timesteps. Since we had leftover compute due to the small batch size, we ran two parallel model training runs with different random seeds and select the lowest log likelihood model for fine-tuning. Since we use such a low learning rate and batch size to keep training stable, given limited compute availability, we believe the model may be undertrained.

While we make both models freely available at <https://huggingface.co/pile-of-law/legalbert-large-1.7M-1/> and <https://huggingface.co/pile-of-law/legalbert-large-1.7M-2/>, we only select the lowest perplexity model for fine-tuning and call this PoL-BERT-Large. We will make intermediate checkpoints, taken at 50k step intervals, available on request. Model Cards [86] are hosted on the HuggingFace website along with the models. We evaluate the checkpoint on the CaseHOLD legal reasoning task [129] and use the train/validation/test split from [30]. Table 7 shows the results. We run fine-tuning using a small hyperparameter search for those marked as (tuned). All other results are those reported by [30].

The Pile of Law model does not significantly outperform a BERT-base model trained exclusively on case law data and using a highly in-domain vocabulary for the CaseHOLD task. This is consistent

²²<https://github.com/huggingface/tokenizers>

with a number of recent results that suggest that masked-language model pre-training efficacy may saturate [1], especially as more diverse or distinct data sources are added [92], or may only give significant gains for highly in-domain data [129]. Because Pile of Law has an extremely diverse set of data, it may require more complicated techniques than MLM to boost performance. Moreover, due to the low learning rate we had to use, it may be possible that training for longer is essential for improved gains. We were unable to complete a full epoch over the data. Nonetheless, we see roughly the same performance as reported in [30] and improve over a Bert Large model.

We suspect that these results mean that the choice of vocabulary, pre-train time, and data selection procedure may all play an important role in adaptation. We hope that the BERT-large model provides an initial baseline for Pile of Law pretraining.

G Copyright

Briefly, copyright is of particular concern as a legal risk for pretraining models [16]. While most data in Pile of Law is public domain or under permissive licenses, copyrighted material may appear in the data. While most subsets of the data are extremely unlikely to contain copyrighted data, one source which may contain such material are CourtListener docket entries. In this case attorneys may have submitted exhibits to the court containing such material. For example in *Dr. Seuss Enters., L.P. v. ComicMix LLC*, 983 F.3d 443 (9th Cir. 2020), entire books are available as exhibits (Ex. 4 and Ex. 5). While the use and release of this data as part of court proceedings and legal investigations is protected under fair use in the United States [72], if the data is then used to generate books or other content that competes with the copyright material it may create some risk of infringement. Or if used in a country without a fair use doctrine, there may be infringement concerns. As such, we suggest that the docket entries subset be omitted for training in these particular scenarios.

Finally, while `r/legaladvice`, `r/legaladviceofftopic`, Edgar contracts, and Credit Card Agreements data were published by third parties under a CC license, it may be possible that the underlying data is under a different license. Similarly, bar exam outlines and Terms of Service contracts are under an unknown license, but were made generously publicly available by the creators (links are in the data). For these data sources, the risk of infringement is low and fair use standards should govern in the United States. However, if users wish to be extremely cautious, they may omit these subsets of data as well.

H Experiment Details

H.1 EOIR Privacy

To construct the EOIR dataset from which we train, we first split EOIR opinions into paragraphs. For each paragraph we mask any references to the respondent, including terms like “respondent” and “appellant.” We then extract using regex expressions whether the court used a pseudonym. This is extracted from the header of the file. For example, “Matter of A-B-C.” would be marked as using a pseudonym. From this we set the label to be either pseudonym or no pseudonym. We then train a distillbert model to take the masked paragraph as input and output a prediction of whether the text should be pseudonymized. We provide the masked and labeled data in https://huggingface.co/datasets/pile-of-law/eoir_privacy. As well as the trained distillbert model here: https://huggingface.co/pile-of-law/distilbert-base-uncased-finetuned-eoir_privacy. Figure 4 is an extended figure showing all causal terms learned from this data based on hand-crafted clustering. We describe some reasons for the cluster groupings below.

Reasoning for unobvious cases. While some connections from vocabulary to topic might be obvious, we provide some reason for classification of non-obvious connections in the following Table:

Word	Group Reasoning
section 209	INA §209 deals with asylum claims
inspected	There is a standard that immigrants be inspected and admitted or paroled on entry into the U.S.
migrants	all contextualized references referred to Cuban and Haitian migrants seeking asylum
killings	all references were to extrajudicial or mass killings commonly referred to in asylum claims
poor	all references referred to the poor conditions of the country of origin in asylum claims
pretermitted	under 8 CFR § 1208.13, Immigration Judges can pretermite an asylum application without a hearing
money	all references were to threats of extortion for money, under threat of violent punishment, directed at asylum claimants by figures in their country of origin
inconsistency	immigration judges look for inconsistencies in credible fear interviews when evaluating asylum claims
promised	all references were to promises made to extortionary figures by asylee claims (e.g., the asylum seeker promised to join a gang if their life was spared, and subsequently left the country to seek refuge)
material share	“a material element” of the asylum application was false when evaluating asylum claims, one determination is whether the applicant is part of a group “composed of members who share a common immutable characteristic.” Matter of M-E-V-G-, 26 I&N Dec. 227 (BIA 2014).
safeguard	all references were to procedural safeguards put in place when a person’s competency to stand trial was in question
building,occupied interest	all references were to burglary
corroborating,corroborative	all references were to competing interests during competency evaluations asylum claims are sometimes evaluated based on corroborating/corroborative evidence for the contents of the application

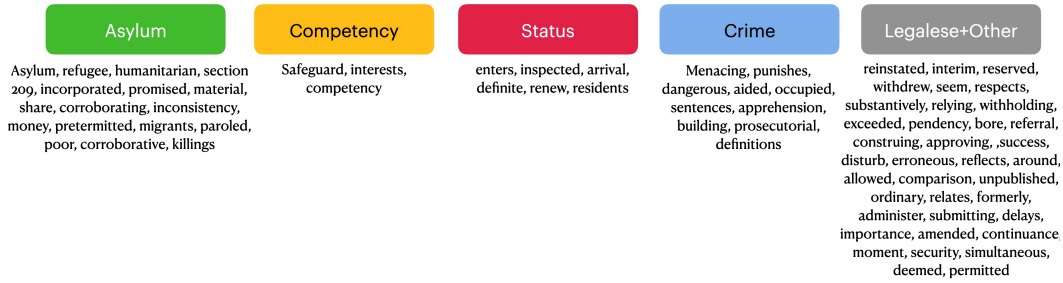


Figure 4: A causal lexicon learned for the EOIR privacy task, manually sorted by topic with contextual information. Year was used as a control variable.

H.2 Mean Toxicity Scores, by Supreme Court Issue Area

Note for toxicity experiments we use the Harvard Case Law Access Project meta-data to associate opinion text to a given date [116].

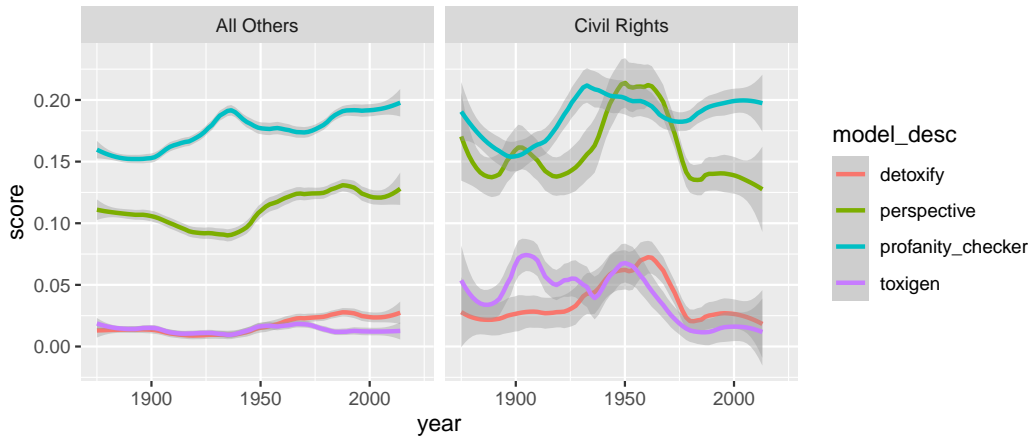


Figure 5: Mean toxicity of ‘highly toxic’ sentence in the average Supreme Court opinion, by issue type and year, 1875-present.

In Figure 5, we provide more detail on the discussion in Section 4.2 regarding the vulnerability of toxicity scores to factual circumstances. For each Supreme Court opinion beginning in 1875, we obtain the sentence in the 97.5th percentile of toxicity as indicated by each of the four filters we study. These sentences represent “highly toxic” sentences within each opinion, while leaving out the 2.5% most extreme cases to reduce sampling variability. We then average the toxicity scores of these highly toxic sentences by year and by whether or not they are categorized as pertaining to ‘civil rights’ by [112].

The right panel suggests that the onset of the Civil Rights Era in the late 1940s corresponds to a rise in the assessed toxicity of Supreme Court cases, despite the fact that many of the cases containing highly toxic sentences were instrumental in dismantling official segregation. For example, the 97.5th percentile of toxicity in *Brown v. Board of Education*, which ended the official segregation of schools, is .535 as rated by the Perspective API.

H.3 Cohen's Kappa for All Models

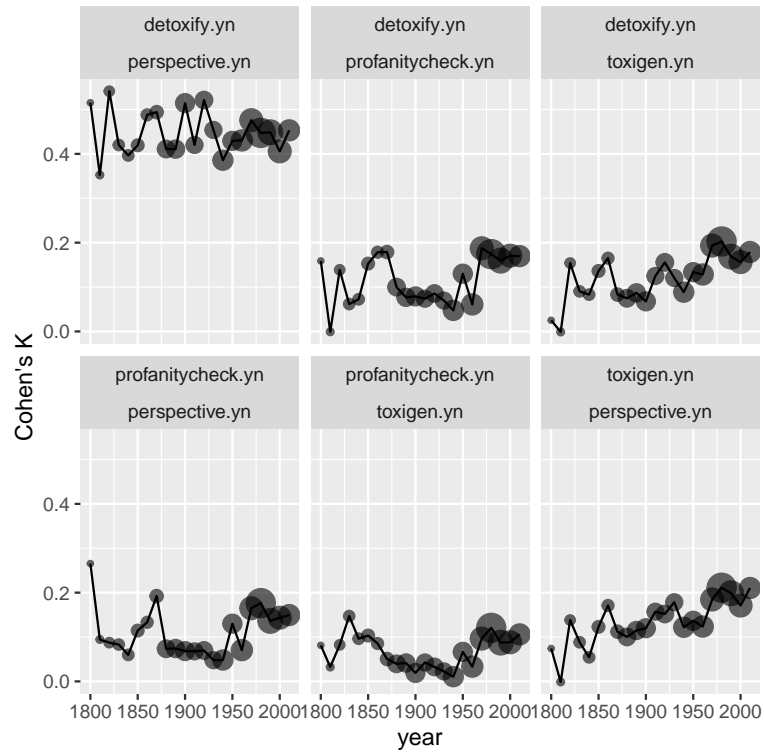


Figure 6: Cohen's κ for all pairs of models studied. Each bin corresponds to the sentences in Supreme Court cases decided over the nearest ten years. Note: detoxify and perspective may have higher agreement since they appear to be trained on similar data sources.

H.4 Qualitative Examples

Content Warning: This section contains quotations that may be offensive or upsetting.

Table 8: Context length perturbation. Adding additional context can make deeply offensive content be labeled as non-toxic using Toxigen [51]. Red indicates the toxic sentence with score paired in red. Black font indicates the appended sentence and toxicity score for the entire paragraph. We redact identifying information and epithets for presentation when applicable. Top two are hand-identified, bottom are randomly samples toxic sentences according to Toxigen.

Category	Quote	Score
True Use	<p>“There is a race so different from our own that we do not permit those belonging to it to become citizens of the United States. Persons belonging to it are, with few exceptions, absolutely excluded from our country. I allude to the Chinese race. But, by the statute in question, a Chinaman can ride in the same passenger coach with white citizens of the United States, while citizens of the black race in Louisiana, many of whom, perhaps, risked their lives for the preservation of the Union, who are entitled, by law, to participate in the political control of the State and nation, who are not excluded, by law or by reason of their race, from public stations of any kind, and who have all the legal rights that belong to white citizens, are yet declared to be criminals, liable to imprisonment, if they ride in a public coach occupied by citizens of the white race. It is scarcely just to say that a colored citizen should not object to occupying a public coach assigned to his own race.”</p>	0.78/0.03
True Use	<p>“The appellant, [NAME REDACTED], age [AGE REDACTED] years, was a homosexual p[****]t. The evening of [DATE REDACTED], the decedent was in a tavern of a hotel in [LOCATION REDACTED]. [NAME REDACTED], who had a room in this hotel, also came into the tavern, and in a short time took a seat by the decedent and engaged him in conversation.”</p>	0.98/0.07
Mention (quoted)	<p>“It was buzzards of the [NAME REDACTED] stripe who shot down my buddy. It was Barnett gunmen who staged the assault on [NAME REDACTED]. It is Jew thugs who have ‘pulled’ practically every robbery in this city. It was a member of the Barnett gang who shot down [NAME REDACTED] while he stood in the shelter of [NAME REDACTED]’s ham-cavern on [NAME REDACTED] avenue.”</p>	0.99/0.00
Mention (quoted)	<p>“ ‘ Other signs would most naturally have been understood as suggesting — falsely—that [NAME REDACTED] was gay. Homosexuality was the theme of many of the signs. There were signs reading “God Hates F[***],” “Semper Fi F[***],” “F[***] Doom Nations,” and “F[**] Troops.” Id., at 3781-3787.’ ”</p>	0.99/0.01
True Use	<p>“And one thing that we very often say and talk about is the three classes of liars. There is the plain liar, the damn liar, and the expert witness. And of all of them, the expert witness is the worst. “There were a few of them here.””</p>	0.91/0.00
Mention (quote)	<p>“ ‘The two statutes under which appellants were convicted and sentenced are part of a comprehensive statutory scheme aimed at prohibiting and punishing interracial marriages. The Lovings were convicted of violating §20-58 of the Virginia Code: “Leaving State to evade law. – If any white person and colored person shall go out of this State, for the purpose of being married, and with the intention of returning, and be married out of it, and afterwards return to and reside in it, cohabiting as man and wife, they shall be punished as provided in § 20-59, and the marriage shall be governed by the same law as if it had been solemnized in this State. The fact of their cohabitation here as man and wife shall be evidence of their marriage.” ’ ”</p>	0.82/0.40
Ambiguous	<p>“Respondents inexplicably make no effort to address Chapter 2 under the Agostini test. Instead, dismissing Agostini as factually distinguishable, they offer two rules that they contend should govern our determination of whether Chapter 2 has the effect of advancing religion. They argue first, and chiefly, that “direct, nonincidental” aid to the primary educational mission of religious schools is always impermissible. Second, they argue that provision to religious schools of aid that is divertible to religious use is similarly impermissible.’ ”</p>	0.77/0.01

I Details on Comparative Law

Table 9: Availability of Identifying Information Across Administrative Settings

Jurisdiction	Civil Cases	Criminal Cases	Juvenile Data
U.S. Federal Courts	Generally available, <i>except</i> dates of birth, financial account numbers, or social security numbers [Fed. R. Civ. Pr. 5.2(a)]. Cases are sealed under exceptional circumstances.	Generally available once filed and, for federal crimes, are never sealed. Filings cannot include dates of birth, SSNs, or residential addresses. [Fed. R. Crim. Pr. 9037(a).]	Juvenile criminal records are generally confidential [18 U.S.C. §5038]. Juvenile names must be partially redacted from civil records [Fed. R. Civ. Pr. 5.2(a)(3)].
U.S. Administrative Agencies	Identifying information, including names, generally omitted from public records. [5 U.S.C. 552(a)]	N/A	Disclosure not explicitly forbidden, but likely more stringent than adult disclosure. [See text]
German Courts	Public judgments exclude all identifying information. [GVG §174]	Most criminal records are automatically expunged five years after the completion of the sentence or three years after death.	Juvenile criminal records are excluded entirely from federal data unless the conviction results in a sentence of over 1 year; even then access is restricted. [81]
Chinese Courts	The names of litigants and case details are public for most cases except if they fall into an excluded category, such as disputes resolved via mediation [75].	The names of litigants and case details are public for criminal cases as of 2016 [75].	Juvenile criminal records are categorically exempt from disclosure [75].
Canadian Courts	The names of litigants and case details are public for most cases unless they meet the standard for sealing.	Criminal records are generally public, but all criminal records are eligible for ‘suspension’ (a form of sealing) after a certain period of good behavior. [Criminal Records Act]	Youth criminal records are always confidential and are automatically sealed after the conclusion of the sentence [Youth Criminal Justice Act].

J Qualitative Examples

We include some qualitative examples from the data to better understand its nature. For example, an OLC Memo might look like the following:

Constitutionality of the Qui Tam Provisions of the False Claims Act Qui tam suits brought by private parties to enforce the claims of the United States violate the Appointments Clause of the Constitution because qui tam relators are “Officers of the United States” but are not appointed in accordance with the requirements of the Appointments Clause. Private qui tam actions violate the doctrine of Article III standing because the relator has suffered no personal “injury in fact.” The qui tam provisions of the False Claims Act violate the separation of powers doctrine because they impermissibly infringe on two aspects of the President’s authority to execute the laws: the discretion whether to prosecute a claim and the authority to control the conduct of litigation brought to enforce the Government’s interests. Given qui tam’s clear conflict with constitutional principles, any argument²³

An excerpt from the US Bills looks like:

113 S2875 IS: National Guard Investigations Transparency and Improvement Act of 2014

U.S. Senate

2014-09-18

Pursuant to Title 17 Section 105 of the United States Code, this file is not subject to copyright protection and is in the public domain.

II

113th CONGRESS

2d Session

S. 2875

IN THE SENATE OF THE UNITED STATES September 18, 2014

Mr. Begich introduced the following bill; which was read twice and referred to the Committee on Armed Services A BILL

To codify in law the establishment and duties of the Office of Complex Administrative Investigations in the National Guard Bureau, and for other purposes.

Be it enacted by the Senate and House of Representatives of the United States of America in Congress assembled,

1. Short title.

This Act may be cited as the National Guard Investigations Transparency and Improvement Act of 2014.

2. Codification in law of establishment and duties of the Office of Complex Administrative Investigations in the National Guard Bureau.

(a) In general.—There is in the Office of the Chief of the National Guard Bureau the Office of Complex Administrative Investigations (in this section referred to as the “Office”)²⁴

²³https://www.justice.gov/sites/default/files/olc/opinions/1989/07/31/op-olc-v013-p0207_0.pdf

²⁴<https://www.congress.gov/bill/113th-congress/senate-bill/2875/text?r=9&s=1>