
Freeing the Law with LOCUS: A Local Ordinance Corpus for the United States

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Progress in legal AI increasingly depends on access to authoritative legal text at
2 scale. Yet one of the most consequential layers of American law remains largely ab-
3 sent from existing machine-readable corpora: local ordinances. Local codes govern
4 zoning, housing, business licensing, public health, noise, animal control, and many
5 other domains of everyday regulation, but they are fragmented across vendor plat-
6 forms designed for human browsing rather than bulk research access. We introduce
7 LOCUS—the Local Ordinance Corpus for the United States—a comprehensive cor-
8 pus and county-harmonized access layer for U.S. municipal and county ordinance
9 codes. The raw corpus, available for release to researchers, represents nearly all
10 publicly available municipal and county ordinance codes. The resulting raw corpus
11 contains codes from 9,239 cities and counties. A smaller county-harmonized LO-
12 CUS access layer provides coverage for the largest 2,309 of 3,144 U.S. counties,
13 accounting for a majority of the population. We use OCR to handle the myriad of
14 document formats that have kept the law from being a public resource. We release
15 the corpus with coverage metadata to support reproducibility, downstream legal AI
16 research, and the incremental expansion of machine-readable access to local law.
17 We train a collection of ModernBERT-based classifiers and scorers to facilitate an-
18 alyzing U.S. local law among several dimensions, such as opacity and paternalism,
19 that have not previously been studied at this scale. LOCUS-v1 and its derivative
20 models are available at: <https://huggingface.co/datasets/LocalLaws/LOCUS-v1>

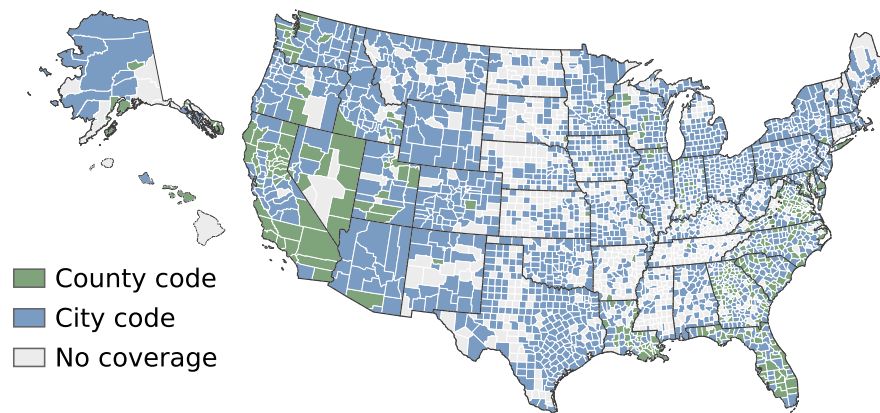


Figure 1: LOCUS represents the longest digitally available code—city or county—for each county.

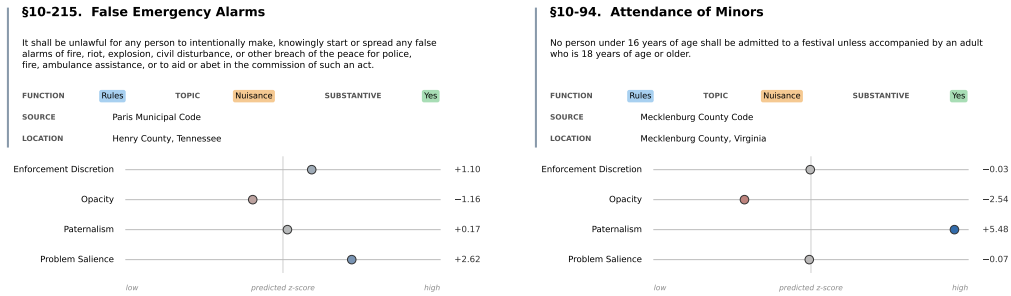


Figure 2: Two example ordinances with predicted scores (in standard units) on four axes (*opacity*, *enforcement discretion*, *paternalism*, and *salience*.) produced by ModernBERT regressors (§5) and function/topic labels produced by ModernBERT classifiers (§4.4). Together they demonstrate the per-ordinance analysis enabled by LOCUS.

21 **1 What it means to "free the law"**

22 Legal AI systems increasingly operate over statutes, cases, regulations, contracts, and administrative
 23 materials [Chalkidis et al., 2022, Henderson et al., 2022, Guha et al., 2023]. This expansion has been
 24 accompanied by domain-specific resources for case law [Zheng et al., 2021], contracts [Hendrycks
 25 et al., 2021, Koreeda and Manning, 2021, Tugener et al., 2020], and statutory reasoning [Holzen-
 26 berger et al., 2020]. Despite this, they still lack systematic access to one of the most consequential
 27 layers of American law: local ordinances. These codes govern zoning, housing, building permits,
 28 business licensing, public health, noise, signs, animal control, and other domains of everyday reg-
 29 ulation. For many questions faced by residents, businesses, landlords, and local governments, the
 30 relevant legal text is not only federal or state law, but a municipal or county code.

31 Local law is not merely another collection of statutes. It is a layered system of legal authority. State
 32 statutes, county ordinances, municipal codes, home-rule provisions, charters, preemption doctrines,
 33 and issue-specific delegations can all interact. Whether a state rule, county rule, or municipal rule
 34 controls is often not obvious in the abstract and may depend on the legal domain. This makes local
 35 law a particularly important setting for legal AI: a useful system must not only retrieve text, but
 36 identify the relevant jurisdictional layer and reason about overlap, delegation, and conflict among
 37 sources of authority.

38 We introduce **LOCUS-v1**, a large-scale corpus and county-harmonized access layer for U.S. local
 39 ordinances. The first release of LOCUS adopts a deliberately transparent simplification: for each U.S.
 40 county, we record the most substantial available local code among the county ordinance code and the
 41 ordinance code of the county’s largest municipality, using document length as a reproducible proxy
 42 for local-law coverage. This representation does not purport to decide which local authority controls
 43 every legal question. Rather, it provides a common geographic substrate on which local legal text can
 44 be searched, compared, and connected to population, geographic, Census, and policy data.

45 The need for such a dataset arises because local law is public but not practically available as a
 46 national research corpus. Georgetown Law Library [2026], the most applied-to law school in the
 47 United States comments, “*there is unfortunately no single source where you can find a comprehensive*
 48 *collection of all municipal codes.*” U.S. local codes are fragmented across commercial vendor
 49 platforms designed for in-browser reading rather than bulk research access. Vendors expose different
 50 navigation structures, print workflows, dynamically generated PDFs, and jurisdiction indexes. No
 51 central registry maps every county or municipality to its hosting platform, and no vendor provides
 52 a complete machine-readable index of all jurisdictions it hosts. As a result, constructing a national
 53 corpus requires discovering where each code lives, extracting it through platform-specific workflows,
 54 validating the resulting artifacts, and harmonizing them to a common unit of analysis.

55 We leave full issue-specific hierarchy and conflict modeling to later releases and benchmark tasks.
 56 This staged design reflects both the legal complexity of determining controlling authority and the
 57 need to preserve uncontaminated evaluation settings for future legal-reasoning benchmarks.

58 LOCUS enables a new class of legal AI and empirical legal studies applications. At the retrieval layer,
59 it supports search and question answering over local rules whose terminology varies substantially
60 across jurisdictions. At the representation layer, it enables structured extraction of regulated activities,
61 permits, fees, penalties, effective dates, and cross-references. At the reasoning layer, it creates a
62 foundation for benchmarks that test whether systems can navigate multiple layers of law, identify the
63 relevant jurisdictional authority, and reason about state-local or county-municipal overlap. By making
64 local law observable at national scale, LOCUS turns a fragmented body of public legal authority into
65 infrastructure for legal retrieval, regulatory extraction, comparative policy analysis, and legal-domain
66 language model evaluation.

67 We provide a summary of our corpus (§3), decision points necessary to create it (§4), evaluations of
68 the corpus (§5), and a discussion of how this can improve our understanding of the legal system (§ 6).

69 **2 Related Work**

70 Studying the law has been important in society for centuries [Holmes, 1897]. In the Information
71 Age, the law has become both immediately accessible but increasingly complicated. We are not
72 the first to create corpora for legal NLP [Steinberger et al., 2006, Aletras et al., 2016, Livermore
73 et al., 2017, Harvard Law School Library Innovation Lab, 2018]. Neural network era corpora such as
74 ECHR [Chalkidis et al., 2019] and pile of law [Henderson et al., 2022] contain case law, court and
75 administrative opinions, and legal codes but not the local law. The 162 tasks in LegalBench [Guha
76 et al., 2023] draw heavily from contracts and merger agreements and none involve local ordinances.

77 Access to the law is a historical challenge which has been reshaped in part by the internet. Georgia
78 v. Public.Resource.Org, Inc., No. 18-1150 (decided April 27, 2020) [Supreme Court of the United
79 States, 2020] upheld that laws, statutes, and court decisions are public domain, in so far as digital
80 content goes. Since that time the rise of large language models and other modern techniques has
81 enabled intelligent data processing on an unprecedented scale; standardizing over 9,239 one-thousand
82 page documents would not have been feasible several years ago. Local laws have been understudied
83 in part due to data access that we hope LOCUS will resolve.

84 **3 Properties of LOCUS**

85 Our corpus benefits both the technical and social science communities by providing valuable data and
86 insight. We discuss the harmonized LOCUS access layer and additional data provided for researchers.

87 **3.1 A County-Harmonized Access Layer**

88 LOCUS adopts a transparent simplification: for each U.S. county, it identifies the most substantial
89 available local code among the county ordinance code and the ordinance code of the county’s largest
90 municipality. This design does not purport to determine which layer of law controls in every doctrinal
91 context. Instead, it provides a reproducible substrate for retrieval, comparison, and future benchmarks
92 on state–county–municipal legal reasoning.

93 Figure 3 summarizes our publicly released corpus: 2,211,516 chunks of text, out of which the
94 majority are judged to be substantive laws in nature. We define substantive as concerned with rules
95 or enforcement, rather than any text that is purely structural, process-oriented, or purely context;
96 the majority of our annotations are rules. These substantive laws deal with four major categories:
97 buildings, business licensing, zoning, and nuisance. The remainder, roughly a third of the laws are
98 categorized with near 90% precision as other. We investigate the headers of these chunks and find
99 that other constitutes topics such as government, employment matters, and animal regulation (this
100 last category makes Alaska have a disproportionately large share of ‘other’ chunks). Table 1 provides
101 examples that illustrate the diversity of laws.

102 **3.2 Additional Data for Researchers**

103 In addition to the released data, we collect an additional 7,000 documents of other cities and counties.
104 We intend to make this data available to researchers with signed release similar to MIMIC [Johnson

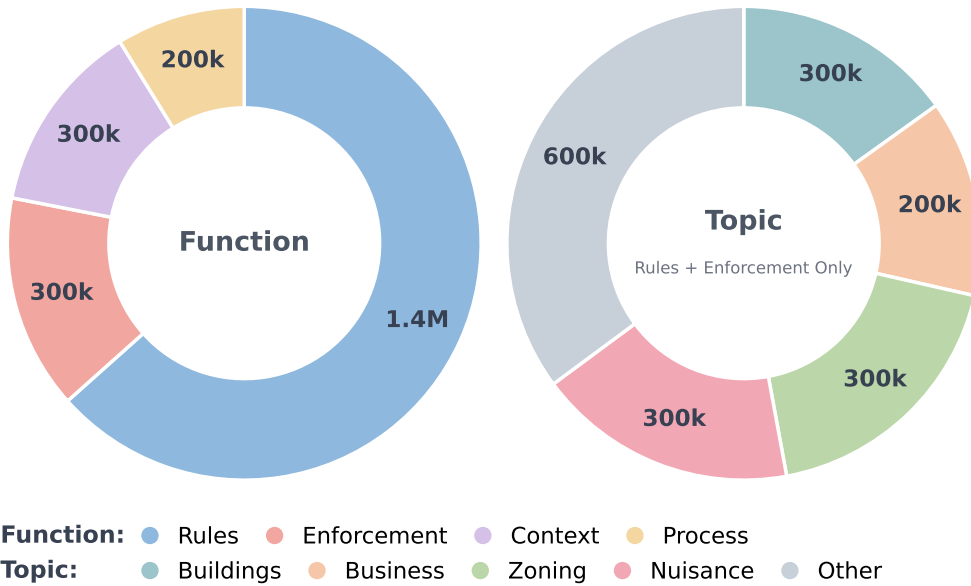


Figure 3: We annotate our corpus at the chunk level along its *Function*, and the substantive laws {Rules and Enforcement} according to the *Topic* referenced. Table 1 provides example texts.

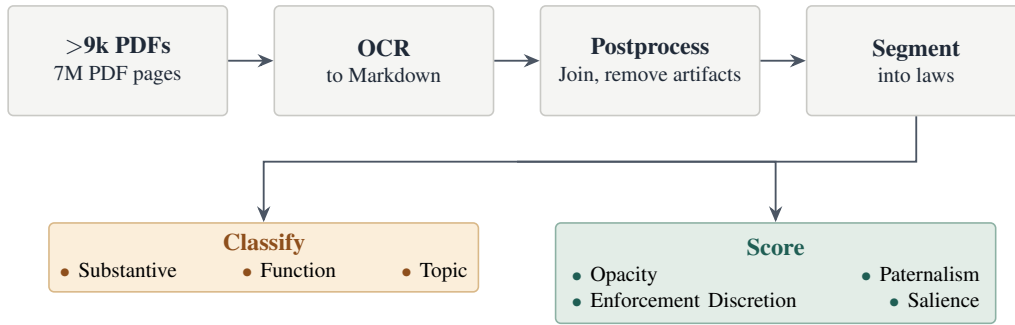


Figure 4: **Processing pipeline.** A corpus of more than 9,000 PDFs (7M pages total) is OCR'd into Markdown, cleaned, and segmented into individual laws. Each segment is then independently *classified* for function, topic, and substance and *scored* on four normative dimensions.

105 et al., 2016, 2023]. Given current LLM ingestion policies, we believe this is necessary for any future
 106 evaluation of local law coverage by foundational models [Dahl et al., 2024].

107 **4 Constructing LOCUS**

108 An overview of the pipeline is shown in Figure 4.

109 **4.1 Collecting the data**

110 The original raw corpus contains 9,239 valid PDFs totaling approximately ~80 GB. Constructing
 111 LOCUS required solving a coupled systems and legal-data problem across thousands of jurisdictions.
 112 Our pipeline uses browser automation and vendor-specific download logic to collect municipal and
 113 county codes from major hosting platforms. The construction process surfaced several nontrivial
 114 failure modes, including server-side PDF assembly limits, filename collisions among non-unique
 115 municipality names, hidden interface thresholds, 15 second crawl delays, anti-bot measures, and
 116 multi-county consolidated cities. Addressing these failures required targeted recovery techniques

Label	Representative example
<i>Function</i>	
RULES	No direct seller shall engage in direct sales within the city without receiving a permit for that purpose as provided herein.
ENFORCEMENT	Any Code Enforcement Officer may issue notices of violation and administrative citations, inspect public and private property, and enforce any available administrative remedies.
STRUCTURAL	Park Regulations 973.01 Adoption and purpose. 973.02 Powers. 973.03 Enforcement. 973.04 Application to concessions.
CONTEXT	The purpose of this notice and review procedure is to notify the public of the permit review process for development proposed in areas having identified significant resources and functional values.
PROCESS	Application for a license required by this division shall be filed with the city clerk on forms provided for that purpose.
<i>Topic</i>	
BUILDINGS	The registered design professional shall submit sufficient technical data to substantiate the proposed alternative engineered design and prove that the performance meets the intent of this code.
BUSINESS	No direct seller shall engage in direct sales within the city without receiving a permit for that purpose as provided herein.
ZONING	Where a change of use of an existing structure requires additional parking or other requirements applicable to a new use, a site plan shall be submitted for review.
NUISANCE	Bike paths may be used only for the operation of bicycles and pedestrian use.
OTHER	Monies appropriated for salaries, wages and related benefits shall not be used for general operations, capital outlay, or other purposes without recommendation from the Mayor and specific approval of a majority of the council.

Table 1: Representative examples for the five *Function* labels and the five merged *Topic* labels in LOCUS. All items in the topic group are annotated as *Rules* or *Enforcement* in their function.

117 rather than a single generic scraper. Furthermore, we manually collect self-hosted or pdf-restricted
118 codes for cities and counties which are not covered by this methodology.

119 4.2 Identifying salient laws

120 Given the huge amount of data, and the diversity of its content and format, we employ a two-level
121 zero-shot approach as the initial labeling approach. Given that our data is being ingested in thousands
122 of different formats after OCR, we need to remove structural content (i.e., stray headers, table of
123 contents) and identify the substantive chunks.

124 After preliminary investigation of Anthropic and Gemini, we settle on OpenAI’s GPT-5.4 as a fast
125 and reliable annotator for this data [OpenAI, 2026]. After comparing a 500 sample of 5.4 mini and
126 nano, we select nano as a cost-effective and only marginally worse option for large-scale annotation.
127 Inspired by LLM-as-a-Judge [Zheng et al., 2023], we evaluate the 5.5% of annotations deemed most
128 challenging with a much more expensive GPT-5.4 model. The model agrees on 64,977 out of 108,889
129 predictions. The more advanced model often decreases its predictions of rules in favor of process and
130 enforcement. Crucially, no models hesitated in identifying structural content, which was ultimately
131 removed from our release. We intend to maintain this dataset and hope to get support from the
132 LLM and law communities in improving these annotations as we update the corpus. Ideally, direct
133 evaluation by lawyers and judges would enable us to exceed the limitations of LLM-as-a-Judge.

134 4.3 OCR and Processing

135 The ordinances are stored in diverse layouts and formats, including single- and double-column
136 layouts, born-digital, exported, and scanned documents, etc. To best handle this diversity, the pipeline

137 for building LOCUS starts by running optical character recognition (OCR) to convert every image of
138 a page to Markdown.

139 We accomplish this with LightOnOCR-2-1B [Taghadouini et al., 2026], an open 1B parameter vision-
140 language model (VLM) based on Qwen-3 [Bai et al., 2025] finetuned on 16MM PDF pages that scores
141 highly on a standard OCR benchmark, OlmOCR-Bench [Poznanski et al., 2025]. LightOnOCR-2-1B
142 generates Markdown text from a page image. We find that this model is robust to the diversity of the
143 raw ordinances, consistently generating correct text in natural reading order.

144 The rest of our post-processing pipeline consumes the unified Markdown output to stitch together
145 laws across pages. We strip artifacts such as repeated headers, footers, and page numbers, and merge
146 content that crosses pages such as paragraphs and tables. The next stage of this post-processing
147 pipeline is to segment the joined content into individual laws, identifying section and subsection
148 headers.

149 The final step of our post-processing pipeline is to classify the substantivity, function, and topic of
150 each extracted law. We discuss the construction of these classifiers in 4.4. Each is trained on the
151 roughly 100M parameter ModernBERT-base [Warner et al., 2025] encoder, which enables us to
152 efficiently run inference on every law. Segments that are classified as purely structural, rather than
153 containing any laws, are omitted from the dataset.

154 The raw ordinances are contained in roughly 7M pages. We are able to scale our OCR pipeline on
155 Modal¹. Given the relatively small size of LightOnOCR-2-1B and Modal’s batch inference support,
156 we were able to efficiently run the entire pipeline and process documents across all formats at roughly
157 \$0.30 per 1,000 pages.

158 4.4 Annotating the Law

159 To organize the ordinances, we develop three classifiers: **substantivity**, **function**, and **topic**. A
160 breakdown of the label space and selected examples are shown in Table 1.

161 We build these classifiers by sampling 100,000 laws from the pipeline discussed in the previous
162 section and using GPT-5.4-nano to annotate each of them. The resulting labels are used to train a
163 ModernBERT classifier [Warner et al., 2025], which can be efficiently used for inference across the
164 rest of the dataset. The classifiers are trained using 80,000 samples for training, 10,000 for parameter
165 sweeps, and finally evaluated on a 10,000 instance subset.

166 From this collection, LOCUS-v1 derives a county-harmonized release that records a representative
167 local-law artifact for each covered county, together with the structured metadata from the classifiers.

168 4.5 Creating a Harmonized Access Layer

169 Our access layer illustrated in Figure 1 is built by a simple algorithm run on all the codes: for every
170 county in the United States, is there an existing county code and an existing city code, ideally from
171 the largest city in that county? If both exist, pick the longest by page length. This is an imperfect
172 process but length of code and population of jurisdiction were correlated.² By doing this, we are
173 able to provide a code for counties *representing* 94% of the United States by population. Since for
174 example the second order city or the population living in the county outside the city are not captured
175 by this, this access layer applies to a smaller literal population than the full data.

176 5 A Dimensional Analysis of Local Laws

177 By linking the text of the laws to the locales in which they apply, LOCUS-v1 opens the door for new
178 types of analysis. In addition to the function and topic metadata, we annotate each ordinance in
179 LOCUS-v1 with dimensional data. We consider four dimensions:

180 1. **Enforcement Discretion** (*highly discretionary to non-discretionary*) — how much selective
181 judgment does the law leave to officials?

¹<https://modal.com>

²Counties run on average slightly shorter than cities, but we opted for an easily interpretable selection algorithm rather than introducing weights; this did not dramatically impact the final selection as certain states, such as Maryland, have much more powerful counties than cities.

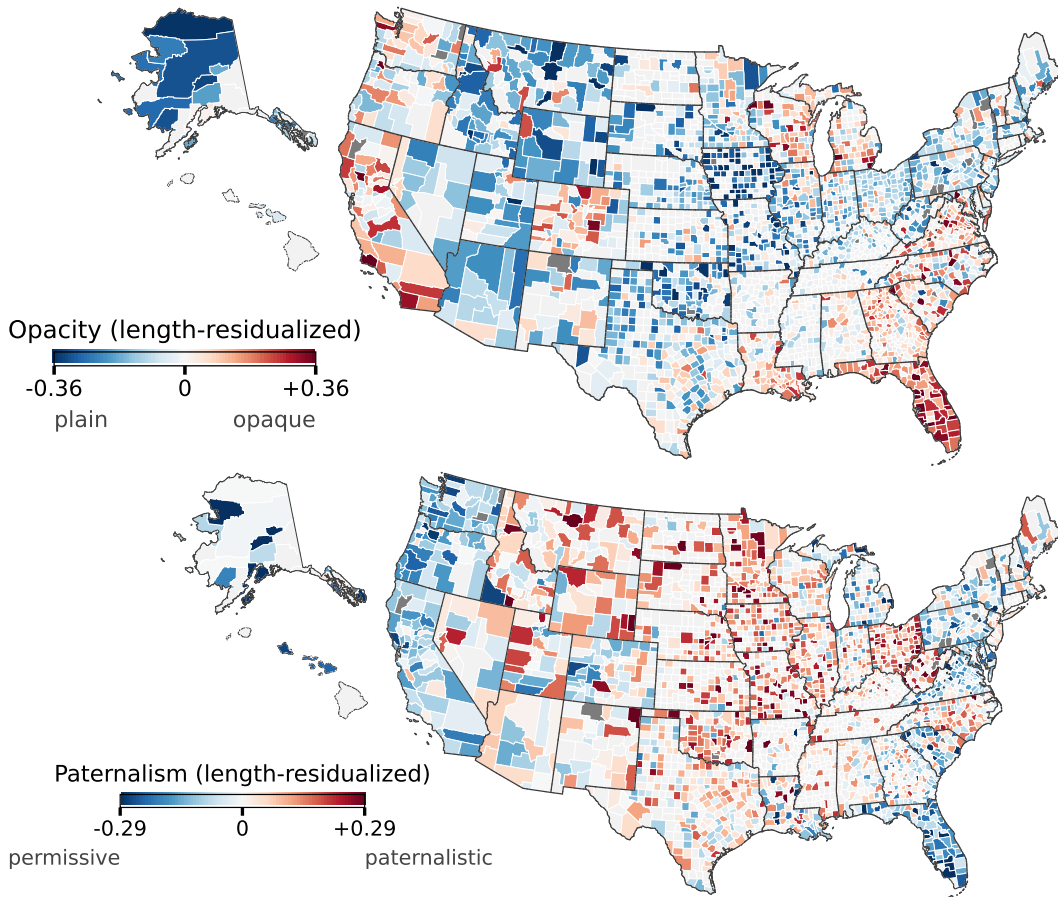


Figure 5: The opacity and paternalism of laws varies across the country. LOCUS facilitates studying the laws for macro trends such as discovering that Florida law is opaque but not paternalistic.

- 182 2. **Opacity** (*opaque to intelligible*) — how hard is it for an ordinary person to know what is
 183 required?
 184 3. **Paternalism** (*paternalistic to externality oriented*) — is it protecting the actor from themself
 185 or protecting others/the public?
 186 4. **Problem Salience** (*highly salient to unimportant*) —how strongly does it represent the issue
 187 as important, urgent, or threatening?

188 Examples of laws occupying these dimensions are given in Figure 2. For instance, the law preventing
 189 minors under the age of 16 from attending a festival unless accompanied by an adult is scored as
 190 highly paternalistic, intelligible, with neutral discretion and salience.

191 Our core intuition is that these dimensions are continuous and that they can better be used to order
 192 and measure laws rather than categorize them. Accurate models of dimensions allow us to bring into
 193 focus particular aspects of the law. Incorporating all laws onto the same set of axes enables analysis
 194 both within individual bodies of law (i.e., within a single city), but also for comparative analysis
 195 across bodies.

196 5.1 Building LOCUS Scorers

197 For our dimensional analysis, we fine-tune a ModernBERT-base with a linear regression head for each
 198 dimension to score a law. For each dimension, we generate 10,000 scores using 200,000 pairwise
 199 LLM-as-a-judge match-ups between ordinances. During each match, we ask the LLM to compare
 200 the two ordinances along a specific dimension, and return which better exemplifies that dimension.
 201 The model outputs A, B, or Tie. Order can produce bias in pairwise judgement [Liu et al., 2024],

202 so every (A, B) comparison pair is also judged in reverse order (B, A). Pairwise comparison aligns
203 better with human judgement than direct/numeric scoring [Liu et al., 2024], motivating us to use it
204 for dimensional scoring. Each ordinance’s match history is used to compute its latent score along
205 each axis using the Bayesian skill rating system, TrueSkill [Herbrich et al., 2006]. This gives us a
206 total ordering over the sampled ordinances, along with an underlying mean, μ .

207 To train the regression model, we normalize the scores to their z-score by subtracting out the
208 dimension’s mean and dividing its standard deviation. For each dimension, we split the 10,000 scored
209 ordinances into a training set (n=8,000), validation set (n=1,000), and test set (n=1,000). We fine-tune
210 a ModernBERT regression model to predict the normalized TrueSkill score, using mean-squared
211 error as our loss function. To evaluate the model, we compute Pearson correlation on the test set.
212 This technique is inspired by the methodology behind Havelock.ai, an AI-powered orality detector
213 that scores text on how oral or literate it is [Weisenthal, 2026].

214 The dataset for each dimension is constructed using a fixed 10,000 ordinance sample, and 200,000
215 pairwise comparisons using GPT-5.4-nano. We report the Pearson correlation coefficient of the
216 trained BERT models versus the TrueSkill values in Figure 6. Each dimension has a correlation
217 of between 0.82 and 0.94, implying the BERT-based scorers largely capture the dynamics of the
218 TrueSkill model. We provide the prompts plus a sample of high- and low-scoring laws along each
219 dimension in Appendix A. We also provide a website to view the TrueSkill scores for the 10,000
220 laws along each dimension.³ We can use these scores to analyze the laws and correlate them with
221 real-world values of interest, discussed in the next section.

222 5.2 Analysis

223 Figure 5 demonstrates the importance of studying this at a nationwide rather than a single case level.
224 For example, counties are notably more opaque than cities on average and Florida is more than twice
225 as opaque as any other state. Studying multiple dimensions in tandem can unlock new insights into
226 unique laws; opacity and paternalism are only weakly correlated across sections (Pearson $r=0.11$ on
227 $n=2,211,516$).

228 Finding interesting needles in this haystack of laws can be facilitated through this evaluation. For
229 example, curfews are detected with paternalism and a subsequent analysis of the data provides insight
230 into curfew distribution for minors across the United States. Headers containing ‘possession’ and
231 ‘alcoholic’ are associated with paternalistic laws while ‘definitions’ and ‘variances’ are associated
232 with opaque ones.

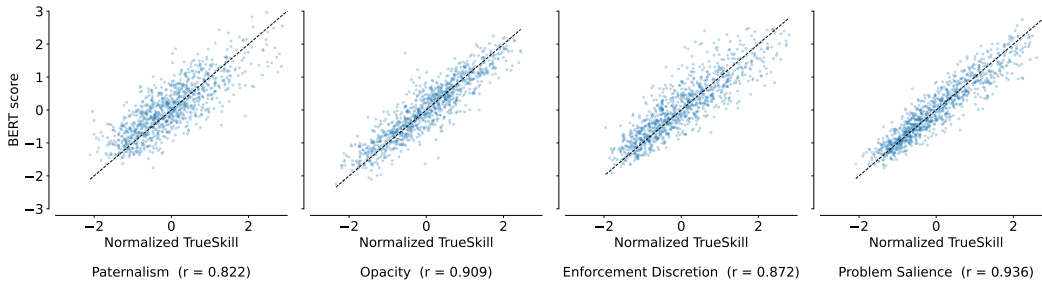


Figure 6: The Pearson correlation between the predicted BERT scores and the normalized TrueSkill scores on 4 distinct test sets (1,000 ordinances per dimension).

233 6 Discussion, Limitations, and Future Work

234 LOCUS-v1 is designed as an access layer, not as a final theory of local legal authority. Its county-
235 harmonized release adopts a transparent simplification: for each county, we select the most substantial
236 available local code among the county code and the code of the county’s largest municipality. This
237 design makes local law searchable and comparable on a national geographic substrate, but it does not
238 determine which rule controls for a particular person, parcel, business, or legal question. In local

³<https://locallaws-locus-leaderboards-web.modal.run>

239 law, authority is layered. State statutes, home-rule provisions, county ordinances, municipal codes,
240 charters, preemption doctrines, and issue-specific delegations may all matter. LOCUS therefore
241 should be understood as infrastructure for retrieval, comparison, and benchmark construction rather
242 than as a substitute for doctrine-sensitive legal analysis.

243 The corpus itself shows why this distinction matters. City and county codes are not interchangeable
244 legal objects. Across the raw corpus, county codes contain substantially more zoning material,
245 while city codes contain more nuisance and public-order regulation. This pattern is consistent
246 with a functional division of local authority: counties more often regulate land, development, and
247 unincorporated territory, while cities more often regulate density, proximity, and everyday public
248 order. For downstream users, this means that jurisdiction type is not merely provenance metadata.
249 It is part of the substantive representation of local law. Models trained or evaluated on local codes
250 should therefore preserve whether a text comes from a municipal or county source, even when the
251 release is harmonized to a county-level unit of analysis.

252 LOCUS also reveals that local codes share a common representational architecture. When ordinances
253 are ordered by their position in a code, topics tend to appear in a stable sequence: general provisions
254 and governmental structure near the front, followed by business regulation, nuisance and public-order
255 rules, zoning, and building regulation. This finding suggests that local law is not simply a bag of
256 rules. It is organized through a recurring documentary form. That form matters for legal AI. Retrieval
257 systems, chunking strategies, and benchmark designs that ignore position within a code may miss
258 information embedded in the structure of codification itself.

259 At the same time, LOCUS documents the limits of any simple national harmonization. In much
260 of the country, counties and cities follow the functional pattern described above. In the Northeast,
261 however, the relationship changes: counties appear less zoning-heavy and more enforcement-oriented,
262 consistent with a different institutional history in which towns and municipalities retain more primary
263 land-use authority while counties often perform administrative, health, or enforcement functions.
264 The implication is not that harmonization is impossible. Rather, it is that harmonization must be
265 explicit about what it preserves and what it abstracts away. A county-level substrate is useful because
266 counties form a mutually exclusive and exhaustive national geography, but the legal meaning of a
267 county code is not constant across states and regions.

268 These limitations point directly to the next generation of legal AI benchmarks. A system that can
269 answer questions about local law must do more than retrieve a plausible ordinance. It must identify
270 the relevant layer of government, distinguish city from county authority, incorporate state-law context,
271 recognize when multiple sources overlap, and reason about whether a retrieved text is actually
272 controlling for the issue at hand. LOCUS-v1 provides the text, metadata, and geographic substrate
273 needed to build those tasks while preserving a clean separation between corpus construction and
274 future evaluations of legal reasoning.

275 More broadly, LOCUS shows that freeing the law is not only a problem of access. It is a problem of
276 representation. Local ordinances were formally public before LOCUS, but they were not available as
277 a national object of machine reading, systematic comparison, or computational legal analysis. Once
278 made observable at scale, local law appears neither as an undifferentiated mass of rules nor as a
279 set of isolated municipal idiosyncrasies. It has structure: a recurring architecture of codification, a
280 functional division between jurisdictional forms, and regionally specific institutional variation. These
281 are precisely the kinds of structure that legal AI systems must learn to respect if they are to move
282 from text retrieval toward reliable reasoning over public authority.

283 **References**

284 Nikolaos Aletras, Dimitrios Tsarapatsanis, Daniel Preoțiu-Pietro, and Vasileios Lampos. Predict-
285 ing judicial decisions of the European Court of Human Rights: A natural language processing
286 perspective. *PeerJ Computer Science*, 2:e93, 2016. doi: 10.7717/peerj-cs.93.

287 Shuai Bai, Yuxuan Cai, Ruizhe Chen, Keqin Chen, Xionghui Chen, Zesen Cheng, Lianghao
288 Deng, Wei Ding, Chang Gao, Chunjiang Ge, et al. Qwen3-v1 technical report. *arXiv preprint*
289 *arXiv:2511.21631*, 2025.

- 290 Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. Neural legal judgment prediction
291 in English. In *Proceedings of the 57th Annual Meeting of the Association for Computational
292 Linguistics (ACL)*, pages 4317–4323, 2019.
- 293 Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Martin
294 Katz, and Nikolaos Aletras. LexGLUE: A benchmark dataset for legal language understanding
295 in English. In *Proceedings of the 60th Annual Meeting of the Association for Computational
296 Linguistics (ACL)*, 2022.
- 297 Matthew Dahl, Varun Magesh, Mirac Suzgun, and Daniel E. Ho. Large legal fictions: Profiling legal
298 hallucinations in large language models. *Journal of Legal Analysis*, 16(1):64–93, 2024.
- 299 Georgetown Law Library. State legal research: General and multi-jurisdictional — local govern-
300 ment. <https://guides.ll.georgetown.edu/statelegalresearch/localgovernment>,
301 2026. Last updated February 27, 2026; accessed May 5, 2026.
- 302 Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex
303 Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, et al. Legalbench: A
304 collaboratively built benchmark for measuring legal reasoning in large language models. In
305 *Advances in Neural Information Processing Systems (NeurIPS)*, 2023.
- 306 Harvard Law School Library Innovation Lab. Caselaw Access Project. <https://case.law/>, 2018.
- 307 Peter Henderson, Mark S. Krass, Lucia Zheng, Neel Guha, Christopher D. Manning, Dan Jurafsky,
308 and Daniel E. Ho. Pile of law: Learning responsible data filtering from the law and a 256GB
309 open-source legal dataset. In *Advances in Neural Information Processing Systems (NeurIPS)*,
310 2022.
- 311 Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. CUAD: An expert-annotated NLP
312 dataset for legal contract review. In *Advances in Neural Information Processing Systems (NeurIPS),
313 Datasets and Benchmarks Track*, 2021.
- 314 Ralf Herbrich, Tom Minka, and Thore Graepel. Trueskill™: a bayesian skill rating system. *Advances
315 in neural information processing systems*, 19, 2006.
- 316 Oliver Wendell Holmes, Jr. The path of the law. *Harvard Law Review*, 10(8):457–478, March 1897.
- 317 Nils Holzenberger, Andrew Blair-Stanek, and Benjamin Van Durme. A dataset for statutory reasoning
318 in tax law entailment and question answering. In *Proceedings of the Natural Legal Language
319 Processing Workshop*, 2020.
- 320 Alistair E. W. Johnson, Tom J. Pollard, Lu Shen, Li-wei H. Lehman, Mengling Feng, Mohammad
321 Ghassemi, Benjamin Moody, Peter Szolovits, Leo Anthony Celi, and Roger G. Mark. Mimic-iii, a
322 freely accessible critical care database. *Scientific Data*, 3(1):160035, 2016. doi: 10.1038/sdata.
323 2016.35.
- 324 Alistair E. W. Johnson, Lucas Bulgarelli, Lu Shen, Alvin Gow, Tom Pollard, Steven Horng, Leo An-
325 thony Celi, and Roger Mark. Mimic-iv, a freely accessible electronic health record dataset.
326 *Scientific Data*, 10(1):1, 2023. doi: 10.1038/s41597-022-01899-x.
- 327 Yuta Koreeda and Christopher D. Manning. ContractNLI: A dataset for document-level natural
328 language inference for contracts. In *Findings of the Association for Computational Linguistics:
329 EMNLP*, 2021.
- 330 Yinhong Liu, Han Zhou, Zhijiang Guo, Ehsan Shareghi, Ivan Vulić, Anna Korhonen, and Nigel
331 Collier. Aligning with human judgement: The role of pairwise preference in large language model
332 evaluators. In *Conference on Language Modeling (COLM)*, 2024.
- 333 Michael A. Livermore, Allen B. Riddell, and Daniel N. Rockmore. The supreme court and the
334 judicial genre. *Arizona Law Review*, 59:837–901, 2017.
- 335 OpenAI. Introducing GPT-5.4. <https://openai.com/index/introducing-gpt-5-4/>, March
336 2026. Accessed: 2026-05-07.

337 Jake Poznanski, Luca Soldaini, and Kyle Lo. olmocr 2: Unit test rewards for document ocr. *arXiv*
338 *preprint arXiv:2510.19817*, 2025.

339 Ralf Steinberger, Bruno Pouliquen, Anna Widiger, Camelia Ignat, Tomaz Erjavec, Dan Tufiş, and
340 Dániel Varga. The JRC-Acquis: A multilingual aligned parallel corpus with 20+ languages. In
341 *Proceedings of the 5th International Conference on Language Resources and Evaluation (LREC)*,
342 2006.

343 Supreme Court of the United States. Georgia v. public.resource.org, inc. Slip Opinion No. 18-1150,
344 April 2020. URL https://www.supremecourt.gov/opinions/19pdf/18-1150_7m58.pdf.
345 590 U.S. 255, 140 S. Ct. 1498, 206 L. Ed. 2d 732.

346 Said Taghadouini, Adrien Cavaillès, and Baptiste Aubertin. Lightonocr: A 1b end-to-end multilingual
347 vision-language model for state-of-the-art ocr. *arXiv preprint arXiv:2601.14251*, 2026.

348 Don Tuggener, Pius von Däniken, Thomas Peetz, and Mark Cieliebak. LEDGAR: A large-scale
349 multi-label corpus for text classification of legal provisions in contracts. In *Proceedings of the 12th*
350 *Language Resources and Evaluation Conference (LREC)*, 2020.

351 Benjamin Warner, Antoine Chaffin, Benjamin Clavié, Orion Weller, Oskar Hallström, Said
352 Taghadouini, Alexis Gallagher, Raja Biswas, Faisal Ladhak, Tom Aarsen, et al. Smarter, better,
353 faster, longer: A modern bidirectional encoder for fast, memory efficient, and long context finetun-
354 ing and inference. In *Proceedings of the 63rd Annual Meeting of the Association for Computational*
355 *Linguistics (Volume 1: Long Papers)*, pages 2526–2547, 2025.

356 Joe Weisenthal. Havelock ai. <https://havelock.ai>, 2026. Accessed: 2026-05-06.

357 Lianmin Zheng, Wei-Lin Chiang, Ying Sheng, Siyuan Zhuang, Zhanghao Wu, Yonghao Zhuang,
358 Zi Lin, Zhuohan Li, Dacheng Li, Eric P. Xing, Hao Zhang, Joseph E. Gonzalez, and Ion Stoica.
359 Judging LLM-as-a-Judge with MT-Bench and Chatbot Arena. In *Advances in Neural Information*
360 *Processing Systems*, volume 36, 2023.

361 Lucia Zheng, Neel Guha, Brandon R. Anderson, Peter Henderson, and Daniel E. Ho. When does
362 pretraining help? assessing self-supervised learning for law and the CaseHOLD dataset of 53,000+
363 legal holdings. In *Proceedings of the 18th International Conference on Artificial Intelligence and*
364 *Law (ICAIL)*, 2021.

365 A Scoring Prompts

366 We elicit pairwise judgments from GPT-5.4-nano using a single shared template, parameterized by a
367 rubric for each axis.

Pairwise Comparison System Prompt

```
You are evaluating local laws along the dimension of "{ axis }"
{% include axis + ".tpl" %}
Read the following two laws and determine which has a GREATER degree of { axis
} according to the rubric above.
Respond with the winner and a one sentence explanation of why it is the winner.
-----
Law A
-----
{{ law_a["header"] }}
{{ text_a }}
-----
Law B
-----
{{ law_b["header"] }}
{{ text_b }}
-----
Respond in the following JSON format only:
““json
```

368

```
{
  "winner": "A" or "B" or "Tie"
  "reasoning": "one sentence explanation"
}
'''
```

369

Axis Rubric: Problem Salience

Problem Salience

Problem salience measures how strongly the law represents the regulated issue as important, urgent, or threatening.

A high-salience law uses charged framing (crisis, epidemic, threat), findings/preambles emphasizing severity, or heightened penalties signaling gravity.

A low-salience law treats the issue as routine, technical, or administrative without rhetorical emphasis on stakes.

370

Axis Rubric: Paternalism vs. Externality Orientation

Paternalism vs. Externality Orientation

Paternalism vs. externality orientation measures whether the law is primarily protecting the regulated actor from themselves or protecting others/the public from the actor's conduct.

A highly paternalistic law targets self-regarding behavior (harms or risks borne mainly by the actor).

A law oriented toward externalities targets conduct whose harms fall on third parties or the public at large.

371

Axis Rubric: Opacity / Intelligibility

Opacity / Intelligibility

Opacity / intelligibility measures how hard it is for an ordinary person to know what the law requires of them. A highly opaque law relies on dense cross-references, technical jargon, undefined terms, or convoluted structure that obscure the obligations.

A low-opacity (highly intelligible) law states its requirements in plain, self-contained language a layperson can readily understand.

372

Axis Rubric: Enforcement Discretion

Enforcement Discretion

Enforcement discretion measures the degree to which a citizen's exposure to enforcement under a law depends on official choice rather than on the citizen's own conduct.

It is high when two factors compound: (1) breadth of exposure -- the pool of citizens potentially subject to the law is large because its triggering conditions are vague, evaluative, or so commonly met that many qualify as eligible targets; and (2) textual latitude -- the statute's language gives officials wide freedom over whether, when, against whom, and how to act.

A law that exposes a vast pool but mandates uniform enforcement leaves officials little real choice; a law that grants officials sweeping latitude but exposes no citizens to enforcement at all -- e.g., provisions concerning only internal government structure, personnel, contracting, or interagency procedure -- creates no opportunity for capricious wielding.

The score floor is reserved for laws that do not act on private parties; the score ceiling for laws under which many citizens stand exposed and officials choose freely among them.

373

374 **B Annotation Prompt**

375 We prompt gpt-5.4-nano for an initial zero-shot classification, and review anything evaluated flagged
376 annotations (5.5%) with a second pass of gpt-5.4.

Annotation Prompt

```
SYSTEM_PROMPT = (  
    "You are a legal text classifier specializing in municipal and county "  
    "codes. Return only a JSON object that matches the provided schema."  
)  
  
REVIEW_SYSTEM_PROMPT = (  
    "You are a senior legal QA reviewer. Review the first-pass classification "  
    "carefully, correct it when needed, and return only a JSON object that "  
    "matches the provided schema. Treat Process as the label for operative "  
    "administrative procedures, delegated authority, internal governance rules,  
    "  
    "meeting/quorum/voting rules, board composition, appointments, elections, "  
    "hearings, notice requirements, and permitting workflows. Use Structural "  
    "only for non-operative artifacts or formatting noise such as headers, "  
    "tables of contents, HTML fragments, history/source notes, page markers, "  
    "cross-reference lists, and similar text that does not itself state an "  
    "operative rule or procedure."  
)  
  
USER_INSTRUCTIONS = ""  
Task: Classify the provided text by its primary legal function.  
  
Allowed primary_function values:  
- Context: defines terms, scope, or introductory intent.  
- Rules: imposes permissions, obligations, or prohibitions.  
- Process: describes administrative procedure, authority, or government  
  structure.  
- Enforcement: specifies penalties, violations, appeals, or exceptions.  
- Structural: non-substantive artifacts such as section headers, tables of  
  contents, HTML, history/source notes, page numbers, or formatting remnants.  
  
Output rules:  
- is_substantive must be 1 only for Rules or Enforcement. Otherwise use 0.  
- primary_function must be exactly one of: Context, Rules, Process,  
  Enforcement, Structural.  
- sub_category must be null when is_substantive is 0.  
- When is_substantive is 1, sub_category must be exactly one of:  
  Land use, Noise/Nuisance, Housing, Business licensing, Public space,  
  Building/Safety, Other.  
- logic must be one short sentence.  
"".strip()  
  
REVIEW_INSTRUCTIONS = ""  
Task: Review a first-pass legal text classification and produce the final  
corrected classification.  
  
Rules:  
- If the first-pass classification is correct, keep it and set review_outcome to  
  "confirm".  
- If the first-pass classification is incorrect, correct it and set  
  review_outcome to "override".  
- If the first-pass classification is missing or invalid, classify from scratch  
  and set review_outcome to "fresh".  
- Keep the same classification rules and category set as the first pass.  
- Treat operative internal governance text as Process, not Structural. This  
  includes board composition, quorum, voting, meeting procedures, delegation of
```

377

authority, appointment/removal rules, election administration, hearings, notice, application steps, and similar administrative workflows.

- Use Structural only for non-operative artifacts/noise such as section headers, article labels, tables of contents, page markers, history/source notes, HTML, formatting remnants, or cross-reference lists that do not themselves contain operative requirements.
- Derive is_substantive from primary_function: use 1 only for Rules or Enforcement; use 0 for Context, Process, and Structural.
- review_logic must briefly explain why you confirmed, changed, or freshly classified the text.

```
"".strip()
```

```

CLASSIFICATION_SCHEMA = {
  "type": "object",
  "additionalProperties": False,
  "properties": {
    "is_substantive": {"type": "integer", "enum": [0, 1]},
    "primary_function": {
      "type": "string",
      "enum": [
        "Context",
        "Rules",
        "Process",
        "Enforcement",
        "Structural",
      ],
    },
    "sub_category": {
      "anyOf": [
        {
          "type": "string",
          "enum": [
            "Land use",
            "Noise/Nuisance",
            "Housing",
            "Business licensing",
            "Public space",
            "Building/Safety",
            "Other",
          ],
        },
        {"type": "null"},
      ],
    },
    "logic": {"type": "string"},
  },
  "required": [
    "is_substantive",
    "primary_function",
    "sub_category",
    "logic",
  ],
}

```

379 **NeurIPS Paper Checklist**

380 **1. Claims**

381 Question: Do the main claims made in the abstract and introduction accurately reflect the
382 paper’s contributions and scope?

383 Answer: **[Yes]**

384 Justification: We introduce our dataset and the dimensions of evaluation that it enables in
385 our abstract. Section 3 describes our dataset and Section 5 describes our evaluation.

386 Guidelines:

- 387 • The answer **[N/A]** means that the abstract and introduction do not include the claims
388 made in the paper.
- 389 • The abstract and/or introduction should clearly state the claims made, including the
390 contributions made in the paper and important assumptions and limitations. A **[No]** or
391 **[N/A]** answer to this question will not be perceived well by the reviewers.
- 392 • The claims made should match theoretical and experimental results, and reflect how
393 much the results can be expected to generalize to other settings.
- 394 • It is fine to include aspirational goals as motivation as long as it is clear that these goals
395 are not attained by the paper.

396 **2. Limitations**

397 Question: Does the paper discuss the limitations of the work performed by the authors?

398 Answer: **[Yes]**

399 Justification: The Discussion session addresses limitations in two paragraphs (n-1 and n-2
400 in the section).

401 Guidelines:

- 402 • The answer **[N/A]** means that the paper has no limitation while the answer **[No]** means
403 that the paper has limitations, but those are not discussed in the paper.
- 404 • The authors are encouraged to create a separate “Limitations” section in their paper.
- 405 • The paper should point out any strong assumptions and how robust the results are to
406 violations of these assumptions (e.g., independence assumptions, noiseless settings,
407 model well-specification, asymptotic approximations only holding locally). The authors
408 should reflect on how these assumptions might be violated in practice and what the
409 implications would be.
- 410 • The authors should reflect on the scope of the claims made, e.g., if the approach was
411 only tested on a few datasets or with a few runs. In general, empirical results often
412 depend on implicit assumptions, which should be articulated.
- 413 • The authors should reflect on the factors that influence the performance of the approach.
414 For example, a facial recognition algorithm may perform poorly when image resolution
415 is low or images are taken in low lighting. Or a speech-to-text system might not be
416 used reliably to provide closed captions for online lectures because it fails to handle
417 technical jargon.
- 418 • The authors should discuss the computational efficiency of the proposed algorithms
419 and how they scale with dataset size.
- 420 • If applicable, the authors should discuss possible limitations of their approach to
421 address problems of privacy and fairness.
- 422 • While the authors might fear that complete honesty about limitations might be used by
423 reviewers as grounds for rejection, a worse outcome might be that reviewers discover
424 limitations that aren’t acknowledged in the paper. The authors should use their best
425 judgment and recognize that individual actions in favor of transparency play an impor-
426 tant role in developing norms that preserve the integrity of the community. Reviewers
427 will be specifically instructed to not penalize honesty concerning limitations.

428 **3. Theory assumptions and proofs**

429 Question: For each theoretical result, does the paper provide the full set of assumptions and
430 a complete (and correct) proof?

431
432
433
434
435
436
437
438
439
440
441
442
443
444
445
446
447
448
449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484

Answer: [N/A]

Justification: No theoretical results.

Guidelines:

- The answer [N/A] means that the paper does not include theoretical results.
- All the theorems, formulas, and proofs in the paper should be numbered and cross-referenced.
- All assumptions should be clearly stated or referenced in the statement of any theorems.
- The proofs can either appear in the main paper or the supplemental material, but if they appear in the supplemental material, the authors are encouraged to provide a short proof sketch to provide intuition.
- Inversely, any informal proof provided in the core of the paper should be complemented by formal proofs provided in appendix or supplemental material.
- Theorems and Lemmas that the proof relies upon should be properly referenced.

4. Experimental result reproducibility

Question: Does the paper fully disclose all the information needed to reproduce the main experimental results of the paper to the extent that it affects the main claims and/or conclusions of the paper (regardless of whether the code and data are provided or not)?

Answer: [Yes]

Justification: We release our code on Github, our data on Hugging Face, and discuss our pipeline in Figure 4. Furthermore, all our prompts and model choices are provided in the appendix.

Guidelines:

- The answer [N/A] means that the paper does not include experiments.
- If the paper includes experiments, a [No] answer to this question will not be perceived well by the reviewers: Making the paper reproducible is important, regardless of whether the code and data are provided or not.
- If the contribution is a dataset and/or model, the authors should describe the steps taken to make their results reproducible or verifiable.
- Depending on the contribution, reproducibility can be accomplished in various ways. For example, if the contribution is a novel architecture, describing the architecture fully might suffice, or if the contribution is a specific model and empirical evaluation, it may be necessary to either make it possible for others to replicate the model with the same dataset, or provide access to the model. In general, releasing code and data is often one good way to accomplish this, but reproducibility can also be provided via detailed instructions for how to replicate the results, access to a hosted model (e.g., in the case of a large language model), releasing of a model checkpoint, or other means that are appropriate to the research performed.
- While NeurIPS does not require releasing code, the conference does require all submissions to provide some reasonable avenue for reproducibility, which may depend on the nature of the contribution. For example
 - (a) If the contribution is primarily a new algorithm, the paper should make it clear how to reproduce that algorithm.
 - (b) If the contribution is primarily a new model architecture, the paper should describe the architecture clearly and fully.
 - (c) If the contribution is a new model (e.g., a large language model), then there should either be a way to access this model for reproducing the results or a way to reproduce the model (e.g., with an open-source dataset or instructions for how to construct the dataset).
 - (d) We recognize that reproducibility may be tricky in some cases, in which case authors are welcome to describe the particular way they provide for reproducibility. In the case of closed-source models, it may be that access to the model is limited in some way (e.g., to registered users), but it should be possible for other researchers to have some path to reproducing or verifying the results.

5. Open access to data and code

485 Question: Does the paper provide open access to the data and code, with sufficient instruc-
486 tions to faithfully reproduce the main experimental results, as described in supplemental
487 material?

488 Answer: [Yes]

489 Justification: We release our dataset on Hugging Face and we release our code on Github
490 for maximal reproducibility and transparency.

491 Guidelines:

- 492 • The answer [N/A] means that paper does not include experiments requiring code.
- 493 • Please see the NeurIPS code and data submission guidelines ([https://neurips.cc/
494 public/guides/CodeSubmissionPolicy](https://neurips.cc/public/guides/CodeSubmissionPolicy)) for more details.
- 495 • While we encourage the release of code and data, we understand that this might not
496 be possible, so [No] is an acceptable answer. Papers cannot be rejected simply for not
497 including code, unless this is central to the contribution (e.g., for a new open-source
498 benchmark).
- 499 • The instructions should contain the exact command and environment needed to run to
500 reproduce the results. See the NeurIPS code and data submission guidelines (<https://neurips.cc/public/guides/CodeSubmissionPolicy>) for more details.
- 501 • The authors should provide instructions on data access and preparation, including how
502 to access the raw data, preprocessed data, intermediate data, and generated data, etc.
- 503 • The authors should provide scripts to reproduce all experimental results for the new
504 proposed method and baselines. If only a subset of experiments are reproducible, they
505 should state which ones are omitted from the script and why.
- 506 • At submission time, to preserve anonymity, the authors should release anonymized
507 versions (if applicable).
- 508 • Providing as much information as possible in supplemental material (appended to the
509 paper) is recommended, but including URLs to data and code is permitted.

511 6. Experimental setting/details

512 Question: Does the paper specify all the training and test details (e.g., data splits, hyperpa-
513 rameters, how they were chosen, type of optimizer) necessary to understand the results?

514 Answer: [Yes]

515 Justification: Section 3 discusses our dataset. Our evaluation code is provided on Github
516 and prompting decisions are meticulously documented in the appendix.

517 Guidelines:

- 518 • The answer [N/A] means that the paper does not include experiments.
- 519 • The experimental setting should be presented in the core of the paper to a level of detail
520 that is necessary to appreciate the results and make sense of them.
- 521 • The full details can be provided either with the code, in appendix, or as supplemental
522 material.

523 7. Experiment statistical significance

524 Question: Does the paper report error bars suitably and correctly defined or other appropriate
525 information about the statistical significance of the experiments?

526 Answer: [N/A]

527 Justification: Our paper primarily details our dataset and its uses. Our Hugging Face datacard
528 provides summary statistics of our columns.

529 Guidelines:

- 530 • The answer [N/A] means that the paper does not include experiments.
- 531 • The authors should answer [Yes] if the results are accompanied by error bars, confidence
532 intervals, or statistical significance tests, at least for the experiments that support the
533 main claims of the paper.
- 534 • The factors of variability that the error bars are capturing should be clearly stated (for
535 example, train/test split, initialization, random drawing of some parameter, or overall
536 run with given experimental conditions).

- 537 • The method for calculating the error bars should be explained (closed form formula,
538 call to a library function, bootstrap, etc.)
- 539 • The assumptions made should be given (e.g., Normally distributed errors).
- 540 • It should be clear whether the error bar is the standard deviation or the standard error
541 of the mean.
- 542 • It is OK to report 1-sigma error bars, but one should state it. The authors should
543 preferably report a 2-sigma error bar than state that they have a 96% CI, if the hypothesis
544 of Normality of errors is not verified.
- 545 • For asymmetric distributions, the authors should be careful not to show in tables or
546 figures symmetric error bars that would yield results that are out of range (e.g., negative
547 error rates).
- 548 • If error bars are reported in tables or plots, the authors should explain in the text how
549 they were calculated and reference the corresponding figures or tables in the text.

550 8. Experiments compute resources

551 Question: For each experiment, does the paper provide sufficient information on the com-
552 puter resources (type of compute workers, memory, time of execution) needed to reproduce
553 the experiments?

554 Answer: [Yes]

555 Justification: Section 4.2 describes our annotation set-up. Section 4.3 describes our OCR
556 pipeline.

557 Guidelines:

- 558 • The answer [N/A] means that the paper does not include experiments.
- 559 • The paper should indicate the type of compute workers CPU or GPU, internal cluster,
560 or cloud provider, including relevant memory and storage.
- 561 • The paper should provide the amount of compute required for each of the individual
562 experimental runs as well as estimate the total compute.
- 563 • The paper should disclose whether the full research project required more compute
564 than the experiments reported in the paper (e.g., preliminary or failed experiments that
565 didn't make it into the paper).

566 9. Code of ethics

567 Question: Does the research conducted in the paper conform, in every respect, with the
568 NeurIPS Code of Ethics <https://neurips.cc/public/EthicsGuidelines?>

569 Answer: [Yes]

570 Justification: We mitigate the harms of our research in so far as data is concerned. We do
571 not have human test subjects. We maximize the transparency of our data in the interest of
572 positive societal impact.

573 Guidelines:

- 574 • The answer [N/A] means that the authors have not reviewed the NeurIPS Code of
575 Ethics.
- 576 • If the authors answer [No], they should explain the special circumstances that require a
577 deviation from the Code of Ethics.
- 578 • The authors should make sure to preserve anonymity (e.g., if there is a special consid-
579 eration due to laws or regulations in their jurisdiction).

580 10. Broader impacts

581 Question: Does the paper discuss both potential positive societal impacts and negative
582 societal impacts of the work performed?

583 Answer: [Yes]

584 Justification: We discuss both our desired positive impacts and the challenges of working
585 with data and creating an accurate harmonization in Section 6.

586 Guidelines:

- 587 • The answer [N/A] means that there is no societal impact of the work performed.

- 588 • If the authors answer [N/A] or [No], they should explain why their work has no societal
589 impact or why the paper does not address societal impact.
- 590 • Examples of negative societal impacts include potential malicious or unintended uses
591 (e.g., disinformation, generating fake profiles, surveillance), fairness considerations
592 (e.g., deployment of technologies that could make decisions that unfairly impact specific
593 groups), privacy considerations, and security considerations.
- 594 • The conference expects that many papers will be foundational research and not tied
595 to particular applications, let alone deployments. However, if there is a direct path to
596 any negative applications, the authors should point it out. For example, it is legitimate
597 to point out that an improvement in the quality of generative models could be used to
598 generate Deepfakes for disinformation. On the other hand, it is not needed to point out
599 that a generic algorithm for optimizing neural networks could enable people to train
600 models that generate Deepfakes faster.
- 601 • The authors should consider possible harms that could arise when the technology is
602 being used as intended and functioning correctly, harms that could arise when the
603 technology is being used as intended but gives incorrect results, and harms following
604 from (intentional or unintentional) misuse of the technology.
- 605 • If there are negative societal impacts, the authors could also discuss possible mitigation
606 strategies (e.g., gated release of models, providing defenses in addition to attacks,
607 mechanisms for monitoring misuse, mechanisms to monitor how a system learns from
608 feedback over time, improving the efficiency and accessibility of ML).

609 11. Safeguards

610 Question: Does the paper describe safeguards that have been put in place for responsible
611 release of data or models that have a high risk for misuse (e.g., pre-trained language models,
612 image generators, or scraped datasets)?

613 Answer: [Yes]

614 Justification: We license the data under CC by NC 4.0 to minimize risk. Furthermore, we
615 withhold a large set of laws for researchers to mimic the best practices of MIMIC [Johnson
616 et al., 2023].

617 Guidelines:

- 618 • The answer [N/A] means that the paper poses no such risks.
- 619 • Released models that have a high risk for misuse or dual-use should be released with
620 necessary safeguards to allow for controlled use of the model, for example by requiring
621 that users adhere to usage guidelines or restrictions to access the model or implementing
622 safety filters.
- 623 • Datasets that have been scraped from the Internet could pose safety risks. The authors
624 should describe how they avoided releasing unsafe images.
- 625 • We recognize that providing effective safeguards is challenging, and many papers do
626 not require this, but we encourage authors to take this into account and make a best
627 faith effort.

628 12. Licenses for existing assets

629 Question: Are the creators or original owners of assets (e.g., code, data, models), used in
630 the paper, properly credited and are the license and terms of use explicitly mentioned and
631 properly respected?

632 Answer: [Yes]

633 Justification: All vendors allow data downloading under terms of use, and we follow any
634 robots.txt guidance. Furthermore, fair use supports the collection of this data. We license the
635 data under CC by NC 4.0 and release a fraction of the data as an access layer. Additionally,
636 we cite the models that we use for annotation.

637 Guidelines:

- 638 • The answer [N/A] means that the paper does not use existing assets.
- 639 • The authors should cite the original paper that produced the code package or dataset.
- 640 • The authors should state which version of the asset is used and, if possible, include a
641 URL.

- 642 • The name of the license (e.g., CC-BY 4.0) should be included for each asset.
- 643 • For scraped data from a particular source (e.g., website), the copyright and terms of
- 644 service of that source should be provided.
- 645 • If assets are released, the license, copyright information, and terms of use in the
- 646 package should be provided. For popular datasets, paperswithcode.com/datasets
- 647 has curated licenses for some datasets. Their licensing guide can help determine the
- 648 license of a dataset.
- 649 • For existing datasets that are re-packaged, both the original license and the license of
- 650 the derived asset (if it has changed) should be provided.
- 651 • If this information is not available online, the authors are encouraged to reach out to
- 652 the asset’s creators.

653 13. New assets

654 Question: Are new assets introduced in the paper well documented and is the documentation

655 provided alongside the assets?

656 Answer: [Yes]

657 Justification: Section 3 describes our dataset and our Hugging Face release provides full

658 transparency of the LOCUS harmonization layer.

659 Guidelines:

- 660 • The answer [N/A] means that the paper does not release new assets.
- 661 • Researchers should communicate the details of the dataset/code/model as part of their
- 662 submissions via structured templates. This includes details about training, license,
- 663 limitations, etc.
- 664 • The paper should discuss whether and how consent was obtained from people whose
- 665 asset is used.
- 666 • At submission time, remember to anonymize your assets (if applicable). You can either
- 667 create an anonymized URL or include an anonymized zip file.

668 14. Crowdsourcing and research with human subjects

669 Question: For crowdsourcing experiments and research with human subjects, does the paper

670 include the full text of instructions given to participants and screenshots, if applicable, as

671 well as details about compensation (if any)?

672 Answer: [N/A]

673 Justification: We do not do research with human subjects.

674 Guidelines:

- 675 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
- 676 with human subjects.
- 677 • Including this information in the supplemental material is fine, but if the main contribu-
- 678 tion of the paper involves human subjects, then as much detail as possible should be
- 679 included in the main paper.
- 680 • According to the NeurIPS Code of Ethics, workers involved in data collection, curation,
- 681 or other labor should be paid at least the minimum wage in the country of the data
- 682 collector.

683 15. Institutional review board (IRB) approvals or equivalent for research with human

684 subjects

685 Question: Does the paper describe potential risks incurred by study participants, whether

686 such risks were disclosed to the subjects, and whether Institutional Review Board (IRB)

687 approvals (or an equivalent approval/review based on the requirements of your country or

688 institution) were obtained?

689 Answer: [N/A]

690 Justification: We do not work with human subjects.

691 Guidelines:

- 692 • The answer [N/A] means that the paper does not involve crowdsourcing nor research
- 693 with human subjects.

- 694
- 695
- 696
- 697
- 698
- 699
- 700
- 701
- Depending on the country in which research is conducted, IRB approval (or equivalent) may be required for any human subjects research. If you obtained IRB approval, you should clearly state this in the paper.
 - We recognize that the procedures for this may vary significantly between institutions and locations, and we expect authors to adhere to the NeurIPS Code of Ethics and the guidelines for their institution.
 - For initial submissions, do not include any information that would break anonymity (if applicable), such as the institution conducting the review.

702 **16. Declaration of LLM usage**

703 Question: Does the paper describe the usage of LLMs if it is an important, original, or
704 non-standard component of the core methods in this research? Note that if the LLM is used
705 only for writing, editing, or formatting purposes and does *not* impact the core methodology,
706 scientific rigor, or originality of the research, declaration is not required.

707 Answer: [Yes]

708 Justification: We use LLMs as core methodology for our annotation (Section 4.2) and
709 evaluation (Section 5.1).

710 Guidelines:

- 711
- 712
- 713
- 714
- The answer [N/A] means that the core method development in this research does not involve LLMs as any important, original, or non-standard components.
 - Please refer to our LLM policy in the NeurIPS handbook for what should or should not be described.