

# ReLeVAnT: Relevance Lexical Vectors for Accurate Legal Text Classification

Anonymous ACL submission

## Abstract

Classifying legal documents by relevance to court filings is a foundational task for downstream applications including motion drafting, docket summarisation, retrieval, litigation surveillance, and training data curation. Existing approaches depend on structured metadata, indicative filenames, or large transformer models. These assumptions break down on noisy, unstructured firm-scale corpora and incur substantial compute cost. The authors propose ReLeVAnT, a lightweight framework for binary legal document classification that models relevance as a discriminative phrase-signal task. ReLeVAnT combines n-gram extraction, contrastive score matching, corpus-specific weighting, and a gradient-boosted tree classifier (XGBoost) into a one-time, interpretable keyword extraction stage followed by a shallow classifier. On the TREC Legal Track, ReLeVAnT achieves 92.7% accuracy and 87.0% Precision; on LexGLUE, 98.9% accuracy and 98.8% Precision. ReLeVAnT runs end-to-end on CPU, classifying documents over 100× faster at inference than competitive embedding baselines while matching or exceeding fine-tuned transformer Precision.

## 1 Introduction

Binary classification of legal documents by relevance underpins a wide range of downstream tasks, including legal workflow automation (Grossman and Cormack, 2010), docket summarisation (Bhattacharya et al., 2021; Saravanan et al., 2008), retrieval systems (Chalkidis et al., 2022; Pipitone and Alami, 2024), litigation surveillance (Ashley, 2017; Katz et al., 2017), and training data curation for larger ML systems (Chalkidis et al., 2019; Zhong et al., 2020). Within the scope of this work, relevance is defined as relatedness to court filings and proceedings. The authors approach the problem from the perspective of robustness, coverage, and computational efficiency, motivated by the practi-

Relevant Document	Irrelevant Document
<p>The motion for <b>summary judgment</b> is hereby granted. The <b>plaintiff</b> has provided <b>substantial evidence of breach of contract</b> and <b>negligence</b>. [CITE] [...] The court finds there was sufficient cause for the <b>trial on damages</b> to proceed.</p> <p>[CITE] [...] Furthermore, the defendant is ordered to produce <b>all discovery documents</b> within 21 days. The settlement conference remains scheduled for March 15th, 2024, and [CITE]; will be held in the District Court. Page 4 / 6</p>	<p>This policy sets forth the company's guidelines on <b>workplace conduct</b>. Employees must follow the established procedures for <b>attendance and punctuality</b>. [CITE] [...] The terms of service establish the acceptable use of the organization's resources. [CITE]</p> <p>Violation of this guide because of failure to comply may result in disciplinary action as in the handbook. Page 3 / 4</p>

Figure 1: An example of the keywords found in excerpts of relevant and irrelevant documents. The '[CITE]' placeholder is left behind during clause filtering. The relevant document contains stronger and more frequent signals of relevance than the irrelevant document.

cal constraints of deploying classification at a firm scale.

Existing methods are constrained by document type, compute cost, or metadata assumptions. [Undavia et al. \(2018\)](#) achieves 72.4% accuracy on SCOTUS classification but requires text in opinion format, limiting use beyond appellate cases. [Li et al. \(2025\)](#) obtain strong results at high speed using filename signals, but break on noisy or out-of-scope filenames common in real corpora. [Limsopatham \(2021\)](#) report up to 20.6% accuracy gains on long legal documents using transformers, at compute costs infeasible at firm scale. Methods such as [de Queiroz Santos Filho et al. \(2025\)](#), [Wang et al. \(2022\)](#), and [Watson et al. \(2023\)](#) report strong results but require curated metadata or structured judgments, which are neither reliably available even at major law firms.

To counter these limitations, the authors model classification as a discriminative task driven by phrase-level signals, building on the observation that relevance in legal corpora is often determined by the presence of highly indicative phrases ([Ash-](#)

066 [ley, 2017](#)) (Fig. 1) rather than broad document-  
067 level similarity. ReLeVAnT explicitly contrasts  
068 phrase frequencies between relevant and irrelevant  
069 documents, captures lexical markers that distin-  
070 guish court filings from non-filings, and feeds these  
071 into a gradient-boosted classifier, thereby avoiding  
072 metadata dependence, structural assumptions, and  
073 expensive representation learning.

074 The contributions of this work are as follows:

- 075 • A formulation of relevance classification as  
076 a contrastive phrase-signal task, combining  
077 document-level and corpus-level frequency  
078 through a contrastive scoring function that ad-  
079 mits interpretable keyword extraction.
- 080 • A lightweight classifier for legal document  
081 relevance that depends solely on document  
082 text, free of metadata, filename, or structural  
083 assumptions, making it deployable on noisy  
084 firm-scale corpora.
- 085 • A comprehensive evaluation across 11 base-  
086 lines spanning classical, embedding-based,  
087 fine-tuned transformer, and recent LLM  
088 (GPT-5.4-Nano) methods, mapping the accu-  
089 racy–cost frontier for binary legal document  
090 classification.
- 091 • State-of-the-art Precision (98.8% on  
092 LexGLUE, 87.0% on TREC Legal Track)  
093 at 1.06ms per-document inference on CPU,  
094 about 100-800× faster than embedding and  
095 LLM baselines, matching or exceeding  
096 fine-tuned transformers.

## 097 2 Related Works

### 098 2.1 Classical Text Classification & Retrieval 099 Methods

100 The heuristics central to ReLeVAnT are grounded  
101 in proven foundations in the Language Modelling  
102 space. Term-frequency signals for document analy-  
103 sis were popularised by ([Sparck Jones, 1972](#)), and  
104 the idea of combining term frequency with term  
105 rarity is central to ReLeVAnT. However, classical  
106 TF-IDF weighs terms globally rather than contrast-  
107 ing classes, and does not explicitly model diminish-  
108 ing returns for repeated occurrences, an important  
109 consideration in long documents.

110 BM25 ([Robertson and Zaragoza, 2009](#)) similarly  
111 leverages term-frequency importance and lexical

112 matching, but its relevance is determined by an ex-  
113 ternal query, whereas ReLeVAnT implicitly mod-  
114 els class-dependent relevance without one. Like  
115 SVMs ([Cortes and Vapnik, 1995](#)), ReLeVAnT as-  
116 sumes linear separability, but instead of learning  
117 weights implicitly, it pre-selects discriminative fea-  
118 tures (keywords, see Section 3), reducing noise and  
119 preserving interpretability. Contrastive learning,  
120 whether across modalities ([Radford et al., 2021](#))  
121 or via noise contrast ([Matsuda et al., 2021](#)), is typ-  
122 ically used to learn semantically meaningful rep-  
123 resentations; ReLeVAnT instead uses contrastive  
124 scoring as a lightweight heuristic to isolate discrim-  
125 inative phrases. To our knowledge, no publicly  
126 available method implements explicit discrimina-  
127 tive filtering to align lexical triggers with legal re-  
128 levance.

### 129 2.2 Legal Relevance and Document Filtering

130 Identifying relevant documents within large le-  
131 gal corpora has been studied extensively in  
132 Technology-Assisted Review (TAR) and the TREC  
133 Legal Track ([Cormack et al., 2010](#); [Grossman et al.,  
134 2011](#)), which formalise relevance as responsiveness  
135 to a given matter and rank documents accordingly.  
136 These frameworks are inherently query-driven, re-  
137 quiring a predefined topic description to guide rele-  
138 vance estimation. The authors repurpose the TREC  
139 dataset for the task by labelling certain sets as rele-  
140 vant and irrelevant (Section 4). Continuous Active  
141 Learning ([Cormack and Grossman, 2015](#)) extends  
142 this paradigm with human-in-the-loop refinement,  
143 achieving high recall but relying on iterative la-  
144 labelling, seed queries, and interactive workflows  
145 that are operationally complex. In contrast, ReLe-  
146 VAnT formulates legal relevance as a fixed classi-  
147 fication problem, independent of external queries  
148 or iterative feedback, making it better suited to  
149 large, noisy, unstructured corpora where metadata  
150 or human supervision may be unavailable.

### 151 2.3 Metadata and Structure-based Methods

152 Several top-performing methods in legal classifi-  
153 cation rely on metadata, structured datasets, or cu-  
154 rated features. Filename signals combined with  
155 TF-IDF and lightweight models ([Li et al., 2025](#)) of-  
156 fer a cheap solution similar in spirit to ReLeVAnT  
157 but limit generalisability across corpora. The ap-  
158 proach of [Sebastiani \(2002\)](#) assumes clean meta-  
159 data on authors, keywords, and tags, which is often  
160 inconsistent or missing in legal data ([Ismaylovna,  
161 2024](#)). Domain-specific work in animal protection

(Watson et al., 2023) pioneered keyword-based relevance signals but depends on structured judgments and header cues unavailable in many court filings, while graph-based methods such as (Wang et al., 2022) require more document structure than real-world data typically provides and are sensitive to extraction errors. ReLeVAnT relies solely on document text, generalises beyond narrow domains, and uses NER filtering only to remove non-legal entities.

## 2.4 Deep Learning-based Methods

Deep learning has found application in legal text classification, but a common drawback is the computational and data costs. These factors are largely independent of the proposed method. Transformer-based approaches (Limsopatham, 2021) model contextual and long-range dependencies, neither of which is strictly necessary for relevance estimation, and incur substantial scalability overhead. SCOTUS classification using Word2Vec with autoregressive architectures Undavia et al. (2018) captures token-level patterns but requires large labelled datasets, optimises a global semantic representation rather than relevance-specific signals, and is restricted to structured court opinions. The shallow n-gram model of Joulin et al. (2017) is closer to ReLeVAnT in spirit but operates on character-level n-grams, which carry less significance than word-level n-grams from a legal-relevance perspective (e.g., *complaint & compliance, dismiss & dismissal*). The use of discriminatory word-level n-grams is a novel feature of the proposed method.

## 3 Methodology

This section details the proposed methodology of ReLeVAnT. Building on the idea that discriminatory words are strong signals of relevance, the authors propose a 2-stage pipeline comprising the Keyword Extraction pipeline and the XGBoost-based classification pipeline (CLS).

As shown in Fig. 2, the Keyword Extraction pipeline (KE) begins with NER-based filtering (Nadeau and Sekine, 2007). This removes all person names, but preserves tags of locations, clauses and organisations. This is particularly important in longer documents where entity names, such as names of plaintiffs, defendants, petitioners, respondents, appellants, appellees, occur very often, and can affect the keywords due to the extraction being frequency-based, as explained further. The

remaining clauses and constitutional references are removed to filter out any residual entities left by NER. Then, n-grams are constructed from these filtered texts. These n-grams are key for representing relevance since the signal "Motion to Dismiss" is stronger than individual words like "Motion" or "Dismiss". With these extracted n-grams, contrastive scores are computed. This can be formalised as:

$$CSM(t) = \frac{f_t^+}{f_t^+ + (f_t^-)^p + \epsilon} \quad (1)$$

where  $t$  is the candidate term (extracted n-gram),  $f_t^+$  is the total frequency of  $t$  across all relevant documents,  $f_t^-$  is the total frequency of  $t$  across all irrelevant documents,  $\epsilon$  is the smoothing constant and  $p$  is the penalty exponent applied to negative frequency to control the aggressiveness of punishing terms also in irrelevant docs. Finally,  $CSM(t)$  is the contrastive score for  $t$  in the range  $[0, 1]$ , where 1 is highly specific to relevant documents, and 0 is common in irrelevant documents. Using these scores, the document-level frequency is computed. This is formalised as:

$$DF(t) = \frac{r_t^+}{r_t^+ + r_t^- + \epsilon} \quad (2)$$

where  $r_t^+ = d_t^+ / N^+$  represents the fraction of relevant documents containing the term  $t$ ,  $r_t^- = d_t^- / N^-$  is the fraction of irrelevant documents containing the term  $t$ ,  $d_t^+$  and  $d_t^-$  are the number of positive and negative documents containing  $t$ , and  $N^+$  and  $N^-$  are the total number of positive and negative documents in the corpus. Moreover, a hard filter is implemented to reject terms that appear in any irrelevant document.  $DF(t)$  counts the fraction of relevant vs irrelevant documents containing the term  $t$ . This captures terms that might appear many times in a few irrelevant documents (inflating the negative frequency), but are actually spread across many relevant documents. Here, it should be noted that  $CSM(t)$  does not indicate absence, but rather a lack of lexical signal related to legal relevance in that case.

The combined score is computed by averaging the term-level CSM scores and document-level DF scores.

$$S(t) = w \cdot CSM(t) + (1 - w) \cdot DF(t) \quad (3)$$

where  $w$  is the weighing hyperparameter discussed in the next section. This combined how 'often' and

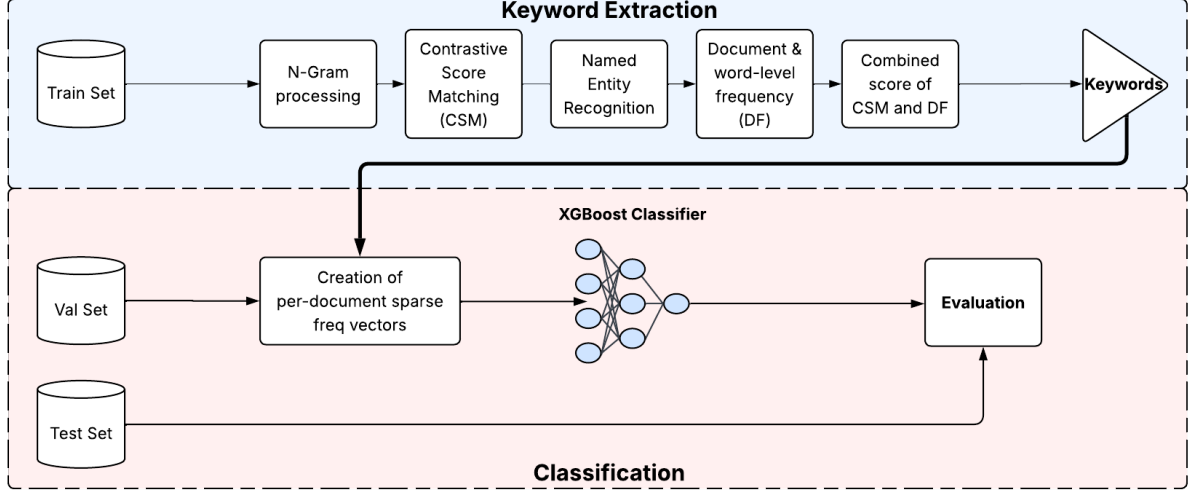


Figure 2: Illustration of the proposed method. The section in Blue highlights the KE stage, and the section in Pink highlights the CLS stage. The given classifier is only for visual purposes, and its exact architecture is detailed in Section 4.

how ‘widely’ a term  $t$  appears in the corpus. This holistic score determines the absolute importance of  $t$  across the entire corpus.

This combined score is weighted by the frequency of the specific keyword to determine the value of the respective index in the resulting vector. This is expressed as:

$$\vec{x}_j = S(t_j) \cdot \text{count}(t_j, \text{doc}) \quad (4)$$

These scores are used to construct a (sparse) vector per document, which is then used by XGBoost to classify documents as relevant or irrelevant. For a given document  $d$  and a set of  $K$  keywords  $t_1, t_2, \dots, t_K$ , the feature vector is formally defined as:

$$\mathbf{v}(d) = \begin{bmatrix} x_1(d) \\ x_2(d) \\ \vdots \\ x_K(d) \end{bmatrix} \quad (5)$$

where each dimension is the score-weighted frequency  $x_j(d) = S(t_j) \cdot \text{count}(t_j, \text{doc})$ .

These feature vectors are constructed for each document, and the documents are classified by an XGBoost gradient-boosted tree ensemble with a logistic objective. The construction of these feature vectors using a well-weighted, relevance-focused discriminative frequency paradigm ensures accurate and reliable classification.

## 4 Experiments

To validate the proposed method, the LexGLUE (Chalkidis et al., 2022) benchmark and TREC Legal Track (Grossman et al., 2011) dataset were used. For TREC, Topic 401 (Online Financial Trading) and Topic 402 (Legality of Derivatives Trading) were included, while Topic 403 (Environmental impact of company activities) was excluded. Topic 403 has substantially lower corpus coverage (1300 judged documents vs 1,444–1,559 for Topics 401 and 402), and with only 44 positive test documents, all methods (including ReLeVAnT, TF-IDF, and fine-tuned LegalBERT) produce 0% relevant-class Precision. It is also thematically distinct from the financial-trading focus of the Enron corpus, making it an outlier. The dataset was split into train, val, and test sets of 2101, 151, and 434 documents (2746, 414, and 1122 pages) respectively, with a relevant-irrelevant ratio of approximately 26.8%-73.2% across all splits. For LexGLUE, ECtHR A and SCOTUS were labelled as relevant (European court cases and US Supreme Court opinions) while EUR-LEX and UNFAIR-ToS were labelled as irrelevant (European legislation and Terms of Service clauses). CaseHOLD and LEDGAR were excluded as they are QA-task-centric and contract-provision-centric, respectively. The dataset was split into 29532, 9675, and 9007 documents (201666, 74538, and 85441 pages) for train, val, and test, with relevant-irrelevant ratios of 47%-53%, 25%-75%, and 27%-73%.

These ratios reflect realistic scenarios where law

Method	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	Pre $\uparrow$
Always positive	26.6	26.6	26.8	26.8
Manual keywords	79.9	64.7	81.7	55.2
Chi <sup>2</sup> +LR	99.8	99.8	90.0	68.4
BM25	41.3	79.5	77.1	45.3
MiniLM+LR	99.4	99.4	85.4	61.9
TF-IDF+LR	99.9	100.0	90.8	92.3
DistilBERT	73.4	0.0	25.8	24.4
LegalBERT	31.5	11.7	37.8	56.6
DistilBERT*	100.0	100.0	90.8	80.0
LegalBERT*	100	100.0	93.2	83.2
GPT-5.4-Nano	98.5	98.5	80.5	51.0
Longformer	–	–	92.8	80.9
ReLeVAnT	<b>98.9</b>	<b>98.8</b>	<b>92.7</b>	<b>87.0</b>

Table 1: Comparison of methods on LexGLUE and TREC. Accuracy and weighted Precision are reported for both datasets. The ‘\*’ denotes results with fine-tuning.

Method	LexGLUE setup (s) $\downarrow$	TREC setup (s) $\downarrow$	Inf. (/doc) $\downarrow$
TF-IDF+LR	34.8	1.2	<1ms
Chi <sup>2</sup> +LR	37.3	1.2	<1ms
BM25	81.8	25	4.6ms
MiniLM+LR	474	44.7	150ms
DistilBERT*	1445	965	2.2ms
LegalBERT*	6420	1440	5.7ms
GPT-5.4-Nano	–	–	870ms
Longformer	–	5220	148.7ms
ReLeVAnT	109.2	369.7	1.06ms

Table 2: Wall-clock setup and inference cost across methods on LexGLUE and TREC. Setup denotes one-time training (or KE+CLS for ReLeVAnT) per corpus; inference is per-document classification cost after setup. The ‘\*’ denotes fine-tuned transformer baselines, which require GPU access. ReLeVAnT is the only competitive method that runs end-to-end on CPU. GPT-5.4-Nano has no setup time since inference is through an API call.

firms hold large volumes of irrelevant data (discovery documents, bills, news reports). ReLeVAnT is invariant to class imbalance and performs well in both regimes, as shown in Table 3. Performance is evaluated using weighted Accuracy and weighted Precision, and all experiments are conducted on an Intel Core Ultra 9 H-class CPU. As described in Section 3, *eyecite* (Cushman et al., 2021) is used to remove clauses and references, the SpaCy core-small model is used for entity filtering, and CSM uses a smoothing constant of 0.01.

All baselines, including ReLeVAnT, are compared in Table 2, which demonstrates ReLeVAnT’s feasibility. Each design choice is ablated thoroughly. Excluding lemmatisation and stemming improves performance on both datasets as seen in Table 4.

The n-gram range sweep (Table 5) selects 4-

Class Split	Accuracy $\uparrow$	Pre $\uparrow$
Realistic	99.3	98.7
Original	98.9	98.1

Table 3: ReLeVAnT’s invariance to class imbalance seen in consistent results between less and more aggressive class imbalance scenarios on LexGLUE. The ‘Realistic’ scenario has a relevant-irrelevant class ratio of 47%-53%, whereas the ‘Original’ ratio in the dataset is 19%-81%.

Method	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	Pre $\uparrow$
ReLeVAnT+LS	95.8	92.1	85.1	70.3
ReLeVAnT	<b>98.4</b>	<b>97</b>	<b>87.0</b>	<b>73.1</b>

Table 4: Comparison of ReLeVAnT with and without lemmatisation and stemming on LexGLUE and TREC. The +LS suffix indicates ReLeVAnT with lemmatisation and stemming at KE and CLS.

grams for LexGLUE and 5-grams for TREC as optimal.

The Minimum Term Frequency (MTF) sweep (Table 6) yields optimal thresholds of 30 for LexGLUE and 10 for TREC, reflecting the differing scales of the two datasets.

The penalty exponent ((1)) is central to scoring and is swept in Table 7, while the document frequency weight (Table 8) balances document-level and cross-document presence.

XGBoost tree depth is swept in Table 9, with depth 4 selected.

The final XGBoost classifier uses 400 estimators, max tree depth of 4, learning rate 0.05, subsample ratio 0.8, and column subsample ratio 0.8 per tree, trained with logloss as the evaluation metric and a per-topic decision threshold tuned on a held-out validation set.

To validate that performance is not driven by a few dominant features or underlying bias, three additional configurations were tested (Table 10): using only the top 3 keywords by frequency (LexGLUE: ‘remanded’, ‘testimony’, ‘congres-

N-gram	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	Pre $\uparrow$
2-gram	97.9	96.0	86.5	76.5
3-gram	98.0	96.3	85.7	78.5
4-gram	<b>98.1</b>	<b>96.4</b>	87.1	86.5
5-gram	98.1	96.4	<b>87.7</b>	<b>88.9</b>
6-gram	98.1	96.4	86.4	87.5
7-gram	98.1	96.3	68.1	68.2
15-gram	98.1	96.4	86.6	88.0

Table 5: Comparison of ReLeVAnT across n-gram ranges on LexGLUE and TREC.

MTF	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	Pre $\uparrow$
10	98.4	97.1	<b>87.3</b>	<b>84.0</b>
20	98.4	97.0	85.0	79.8
30	<b>98.4</b>	<b>97.0</b>	84.8	70.1
50	98.1	96.4	84.0	78.3
100	97.4	95.0	82.6	81.5
250	97.0	94.2	80.5	71.1
400	96.3	92.7	74.2	36.8
1000	93.4	86.5	73.2	0.0

Table 6: Comparison of ReLeVAnT across MTF thresholds on LexGLUE and TREC.

PenEx	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	Pre $\uparrow$
1	98.4	97.0	<b>87.7</b>	<b>83.3</b>
2	98.4	97.0	87.4	83.1
10	<b>98.6</b>	<b>97.3</b>	84.2	82.7
50	98.5	97.1	83.9	83.3
$\infty$	98.1	96.4	83.4	73.9

Table 7: Comparison of ReLeVAnT across penalty exponent (PenEx) thresholds on LexGLUE and TREC.

sional’; TREC: ‘CME’, ‘CBOE’, ‘Mini’), the full keyword set excluding those top 3 (ReLeVAnT-K), and random feature vectors. The full ReLeVAnT pipeline outperforms all three, confirming that the classifier draws signal from the broader discriminative keyword set rather than memorising a handful of cues.

## 5 Results

Experimental results in Table 1 reveal the strong performance of ReLeVAnT relative to baselines such as majority-class prediction, always-positive selection, manual keyword selection, TF-IDF with Logistic Regression, Chi-sq with Logistic Regression, BM25 (Robertson and Zaragoza, 2009), MiniLM embeddings (Wang et al., 2020) with Logistic Regression, DistilBERT (Sanh et al., 2019), LegalBERT (Chalkidis et al., 2020) and GPT-5.4-Nano (OpenAI, 2026). ‘Always positive’ refers to a dummy classifier that classifies everything as relevant; this is the anchor of the results, where any value less than this would be considered erroneous. Manually-chosen keywords with an MLP demonstrate subpar classification at <85% accuracy on both datasets (Appendix B). On TREC, ReLeVAnT achieves higher accuracy than TF-IDF (92.7% vs 90.8%), but TF-IDF achieves marginally higher Precision (92.3% vs 87.0%); this is attributed to TF-IDF’s conservative classification behaviour, which yields high Precision at the cost of substantially lower coverage of the relevant class, a tradeoff

DocFreq	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	Pre $\uparrow$
0.25	98.1	96.4	87.5	78.6
0.33	98.9	97.9	87.7	85.4
0.50	99.2	98.5	87.7	83.3
0.66	99.3	98.6	<b>88.4</b>	<b>86.7</b>
0.75	<b>99.3</b>	<b>98.7</b>	87.9	85.3
0.90	98.9	98.0	87.7	79.6

Table 8: Comparison of ReLeVAnT across document frequency thresholds on LexGLUE and TREC.

Max Depth	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	wPre $\uparrow$
2	98.6	98.6	92.2	84.7
3	98.7	98.7	92.6	86.8
4	<b>98.9</b>	<b>98.8</b>	<b>92.7</b>	<b>87.0</b>
6	98.8	98.8	92.6	86.1
9	98.5	98.5	92.2	86.0
12	98.4	98.4	92.1	85.0

Table 9: Comparison of ReLeVAnT (XGBoost) across tree depths on LexGLUE and TREC. Max depth of 4 achieves optimal precision on both datasets, with deeper trees yielding diminishing returns on the small TREC corpus.

poorly suited to the legal setting where missed filings carry real cost. Similarly, Chi<sup>2</sup>+LR matches TF-IDF on accuracy but trails ReLeVAnT by 18.6% in Precision on TREC, reinforcing that statistical feature selection alone is insufficient when class-discriminative phrases are sparse. BM25, designed for retrieval rather than classification, underperforms across both datasets (41.3% accuracy on LexGLUE, 77.1% on TREC), reinforcing the limitation of query-driven scoring in the absence of a query. MiniLM with Logistic Regression performs strongly on LexGLUE (99.4% accuracy) but drops to 85.4% on TREC with only 61.9% Precision, showing that semantic embeddings struggle when the discriminative signal is lexical and topic-specific rather than stylistic. Moreover, Zero-shot transformer inference fails because pretrained objectives, such as masked language modelling, confer no knowledge of task-specific discriminative

Vector Config	LexGLUE		TREC	
	Acc $\uparrow$	Pre $\uparrow$	Acc $\uparrow$	wPre $\uparrow$
Top 3 keywords	87.9	70.8	78.7	70.4
Random vectors	73.4	0	28.9	21.5
ReLeVAnT-K	98.4	97.0	87.8	75.3
ReLeVAnT	<b>98.9</b>	<b>98.8</b>	<b>92.7</b>	<b>87.0</b>

Table 10: Comparison of ReLeVAnT across specific vector configurations on LexGLUE and TREC. ReLeVAnT-K refers to the entire keyword list except the top 3 keywords.

features. Fine-tuning addresses this by providing a supervised signal that adjusts representations toward class-separating directions in the embedding space; on TREC, fine-tuned DistilBERT trails ReLeVAnT on both metrics (90.8% vs 92.7% accuracy and 80.0% vs 87.0% Precision), while fine-tuned LegalBERT marginally surpasses ReLeVAnT on accuracy (93.2% vs 92.7%) but still trails on Precision (83.2% vs 87.0%). ReLeVAnT achieves a similar goal by explicitly contrasting frequencies upstream, allowing a lightweight gradient-boosted classifier to operate on already-discriminative features rather than learning them end-to-end from raw text. To evaluate whether a modern lightweight LLM can match contrastive lexical signals out-of-the-box, the authors include OpenAI’s GPT-5.4-Nano (OpenAI, 2026) as a zero-shot baseline. GPT-5.4-Nano is the smallest, lowest-latency variant of the GPT-5.4 family, marketed by OpenAI for short-turn classification, data extraction, and ranking tasks, making it a natural fit for binary relevance prediction. It performs well on LexGLUE because the binary framing reduces to genre identification (court opinions vs legislation/ToS), a distinction salient from pretraining alone. It fails on TREC because relevance there requires distinguishing documents within a single corpus that share genre, register, and vocabulary; a discrimination that depends on corpus-specific lexical contrasts the model has no access to in a zero-shot setting.

Beyond accuracy and Precision, Table 2 highlights the practical cost advantages of ReLeVAnT. At inference, ReLeVAnT classifies a document in 1.06ms on CPU, roughly  $141\times$  faster than MiniLM+LR (150ms),  $5.4\times$  faster than fine-tuned LegalBERT (5.7ms on GPU), and  $820\times$  faster than GPT-5.4-Nano (870ms), while operating on commodity hardware rather than requiring sustained GPU access. Setup costs reveal an even sharper contrast: on LexGLUE, ReLeVAnT completes setup in 109.2s on CPU,  $4.3\times$  faster than MiniLM+LR (474s),  $13.2\times$  faster than DistilBERT fine-tuning (1445s on GPU), and  $58.8\times$  faster than LegalBERT fine-tuning (6420s on GPU). On TREC, ReLeVAnT (369.7s on CPU) completes  $2.6\times$  faster than DistilBERT (965s on GPU) and  $3.9\times$  faster than LegalBERT (1440s on GPU). Across both datasets, the cost asymmetry is compounded by the hardware difference: ReLeVAnT runs end-to-end on commodity CPU while the transformer baselines require sustained GPU access. ReLeVAnT does incur a higher setup

cost than the lightest baselines, TF-IDF+LR and Chi-sq+LR, which finish setup in under 40s on LexGLUE and  $\sim 1$ s on TREC, but those baselines pay this back in classification quality, with TF-IDF lagging ReLeVAnT by 1.9% on TREC accuracy and Chi-sq+LR by 18.6% on TREC Precision. ReLeVAnT thus occupies a favourable position on the accuracy–cost frontier: it matches or exceeds fine-tuned transformer Precision on both datasets while running at one to two orders of magnitude lower compute cost, and it materially outperforms the cheapest classical baselines without straying far from their cost regime.

ReLeVAnT is robust to class imbalance on LexGLUE (Table 3), with a 0.4% accuracy and 0.6% Precision difference between regimes, owing to the near independence of negative frequency at CSM (Eq. 1). Tables 5–9 show single-hyperparameter sweeps; the final configuration values appear in Table ??.

To explore the impact of vocabulary scope, the first set of experiments focused on lemmatisation and stemming. As evident in Table 4, ReLeVAnT without this processing step performs considerably better on both datasets, by upto 2.4% in accuracy and 4.3% in Precision on LexGLUE and 1.1% in accuracy and 2.8% in Precision on TREC. Legal language is precise, and morphological variants carry distinct procedural or substantive meaning that lemmatisation and stemming erode, reducing the contrastive signal CSM relies on.

Table 5 justifies the optimal choice of n-grams for each dataset. For LexGLUE, performance improves up to 4-grams, then stagnates, while the keyword count remains stable. TREC demonstrates optimal results at 5-grams, after which performance steadily declines. The disagreement over the optimal n-gram choice between the datasets stems from corpus size: TREC has fewer documents than LexGLUE, leading to longer n-grams in TREC that are less sparse, allowing vectors to capture just enough semantic information for reliable classification.

Table 6 demonstrates optimal performance at MTF 30 for LexGLUE and MTF 10 for TREC. The MTF threshold resolves a tension inherent to the CSM: low-frequency terms can produce extreme contrast scores, but those scores may reflect either genuine but rare topic indicators or statistical coincidences from a single repeated mention. Setting MTF too low admits the latter as features, while setting it too high discards the former. The

disagreement between LexGLUE and TREC is explained by the discriminative signal in LexGLUE being stylistic; formal legal documents distinguish themselves through standardised, repeated phraseology that occurs many times across the corpus, so terms that fail to clear MTF 30 are almost certainly noise. In contrast, TREC’s discriminative signal is topic-specific terminology within a uniform conversational genre, where genuinely informative terms may appear in only a small fraction of the relevant set, demanding a lower MTF to be retained.

Table 7 points to the choice of the penalty exponent in the CSM modelling, as explained in Eq. 1. A higher penalty exponent imposes stricter punishment on keywords that also occur in negative documents. Tuning this to 10 produces the best results for LexGLUE, while tuning to 1 gives maximal returns for TREC. The disagreement is because LexGLUE’s irrelevant documents share a large overlapping vocabulary in legal language, sentence structure, and writing style with the relevant set, hence such terms need to be penalised more. In TREC, which has a smaller corpus, an excess penalty discards relevant but rarer terms, as seen in the decline in performance.

Table 8 showcases the best performance of ReLeVAnT at 0.75 weight for LexGLUE and 0.66 for TREC. This indicates that preference for in-document frequency results in better performance over frequency across documents for both datasets, owing to the size of the dataset and the length of documents. This reveals an interesting finding: local frequency is more influential than global frequency across corpora with longer documents, and vice versa.

Table 9 details the tree depths explored for the XGBoost classifier. A maximum depth of 4 achieves optimal precision on both datasets. Shallower trees (depth 2–3) underfit the feature space, while deeper trees (depth 6–12) progressively overfit on TREC’s smaller corpus. LexGLUE is largely insensitive to depth given the corpus’s size and stylistic separability.

To investigate dependency on top keywords, the authors isolated the top 3 most influential keywords from each dataset. The exact same setup with just these 3 keywords demonstrates passable performance of 87.9% accuracy and 70.8% Precision on LexGLUE, and 78.7% accuracy and 70.4% Precision on TREC. This supports the initial assumption of discriminatory words being a strong, separable signal of legal relevance: ‘remanded’, ‘testimony’,

‘congressional’ for LexGLUE and ‘CME’, ‘CBOE’, ‘Mini’ for TREC (see Appendix A.6 for glosses), all intuitively strong markers of legal relevance. Removing these 3 to assess generalizability reveals that XGBoost can compensate for the lack of very strong signals and performs very similarly to the full keyword list, with only a 0.9% drop in accuracy and 1.7% drop in Precision on LexGLUE, and 0.4% drop in accuracy and 1.4% drop in Precision on TREC. To assess the classifier’s randomness, random feature vectors were generated for each document, yielding an exact class-split accuracy of 73.9% and a Precision of 0. These results further reinforce the hypothesis about legal relevance being measurable through the frequency of discriminatory signals.

These ablations isolate three findings. First, the contrastive keyword selection carries most of the signal: three keywords alone recover 87.9% accuracy on LexGLUE, while random vectors collapse to 0% Precision. Second, the classifier’s role is to integrate redundant signals across the keyword set, which is why removing the top three keywords results in only a 0.9% drop in accuracy. Third, the optimal hyperparameter configuration tracks the discriminative signal in each corpus: LexGLUE’s stylistic signal rewards higher MTF thresholds, stronger penalty exponents, and shorter n-grams, while TREC’s topic-specific signal rewards the opposite. Across both corpora, in-document frequency dominates cross-document frequency, suggesting local repetition outweighs corpus-wide spread when documents are long. Combined with the cost-frontier results, these findings support the broader claim that explicit upstream feature engineering remains competitive with end-to-end transformer learning when the underlying signal is lexical rather than semantic.

## 6 Conclusion and Future Works

This work introduced ReLeVAnT, a lightweight, contrastive, phrase-driven framework that leverages a shallow classifier to achieve state-of-the-art results on binary classification tasks for court filings. It is completely independent of metadata, filenames, document and directory structure, allowing for massive coverage and robustness at an exponentially lower cost than other methods in the legal text field.



## 7 Limitations

The authors outline several limitations of ReLeVAnT and its evaluation, and encourage future work to address them.

### Performance on small, imbalanced corpora.

ReLeVAnT, along with the baselines, performs poorly on TREC Topic 403 (Environmental impact of company activities), where the relevant class comprises only 44 documents (10.6%) of the test set. As discussed in Section 4, this topic produced near-zero relevant-class Precision across nearly all hyperparameter configurations, motivating its exclusion from the main evaluation. This points to a real failure mode of the method: when the relevant pool is too small to populate the contrastive frequency table reliably, CSM scores become dominated by noise rather than signal. Methods designed for low-resource settings (few-shot fine-tuning, active learning) may be more appropriate in this regime. This problem is generally not encountered with top-level law firms.

**Binary classification only.** ReLeVAnT is formulated as binary relevance classification (relevant vs irrelevant). Multi-class and multi-label legal classification can be applied similarly. For example, identifying the specific filing type, jurisdiction, or legal issue is not addressed. The contrastive scoring formulation in Eq. 1 extends naturally to one-vs-rest schemes, but the trade-offs of doing so (keyword overlap across classes, threshold calibration, computational cost scaling with class count) have not been studied here.

**Multi-seed results.** All reported numbers reflect training runs with results averaged across 10 seeds. The authors do not report variance across seeds since it is  $\leq 0.1$  for each result.

**Transformer baselines under input-length constraints.** DistilBERT and LegalBERT were fine-tuned with documents truncated to 512 tokens, following standard practice. Legal documents in both datasets frequently exceed this length, which limits the context available to the transformer baselines. Methods that handle long documents (e.g., Longformer, hierarchical BERT) exist but introduce additional architectural and computational overhead beyond the scope of this comparison; ReLeVAnT processes documents at full length without truncation. This means the gap between ReLeVAnT and the transformer baselines partly reflects trun-

cation rather than purely modelling capacity, and the comparison would need to be revisited against long-context models.

**Dataset scope.** Both datasets are English-language. TREC corresponds to the Enron corpus with two financial-trading topics, and LexGLUE in the binary framing reduces to coarse domain classification (US/European court opinions vs EU legislation and Terms of Service). Neither tests ReLeVAnT against a corpus where relevant and irrelevant documents share a single domain, language, and style. For example, distinguishing motions from briefs within a single case. Generalisation to such fine-grained, intra-domain relevance distinctions is not established.

**Interpretability vs explanation.** ReLeVAnT produces an inspectable keyword list, which the authors treat as a form of interpretability. The authors do not provide per-document explanations (e.g., SHAP-style feature attributions for individual classifications), nor have the authors evaluated whether the extracted keywords align with how legal practitioners actually judge relevance. User studies with domain experts would be needed to validate the method’s practical interpretability.

## References

- Kevin D Ashley. 2017. *Artificial intelligence and legal analytics: new tools for law practice in the digital age*. Cambridge University Press.
- Paheli Bhattacharya, Soham Poddar, Koustav Rudra, Kripabandhu Ghosh, and Saptarshi Ghosh. 2021. Incorporating domain knowledge for extractive summarization of legal case documents. In *Proceedings of the eighteenth international conference on artificial intelligence and law*, pages 22–31.
- Ilias Chalkidis, Emmanouil Fergadiotis, Prodromos Malakasiotis, and Ion Androutsopoulos. 2019. 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.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. **LEGAL-BERT: The muppets straight out of law school**. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2898–2904, Online. Association for Computational Linguistics.
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. 2022. Lexglue: A benchmark

707	dataset for legal language understanding in english.	Nut Limsopatham. 2021. Effectively leveraging bert for	759
708	In <i>Proceedings of the 60th Annual Meeting of the</i>	legal document classification. In <i>Proceedings of the</i>	760
709	<i>Association for Computational Linguistics (Volume</i>	<i>natural legal language processing workshop 2021,</i>	761
710	<i>1: Long Papers)</i> , pages 4310–4330.	pages 210–216.	762
711	Gordon V. Cormack and Maura R. Grossman. 2015.	Takeru Matsuda, Masatoshi Uehara, and Aapo Hyvari-	763
712	Autonomy and reliability of continuous active learn-	nen. 2021. Information criteria for non-normalized	764
713	ing for technology-assisted review. <i>arXiv preprint</i>	models. <i>Journal of Machine Learning Research,</i>	765
714	<i>arXiv:1504.06868</i> .	22(158):1–33.	766
715	Gordon V. Cormack, Maura R. Grossman, Bruce Hedin,	David Nadeau and Satoshi Sekine. 2007. A survey of	767
716	and Douglas W. Oard. 2010. Overview of the trec	named entity recognition and classification. <i>Lingvis-</i>	768
717	2010 legal track. In <i>Proceedings of the Text REtrieval</i>	<i>ticae Investigationes,</i> 30(1):3–26.	769
718	<i>Conference (TREC)</i> . National Institute of Standards	OpenAI. 2026. Introducing gpt-5.4 mini	770
719	and Technology (NIST).	and nano. <a href="https://openai.com/index/introducing-gpt-5-4-mini-and-nano/">https://openai.com/index/introducing-gpt-5-4-mini-and-nano/</a> . Ac-	771
720	Corinna Cortes and Vladimir Vapnik. 1995. Support-	cessed: 2026-05-07.	772
721	vector networks. <i>Machine learning,</i> 20(3):273–297.	Nicholas Pipitone and Ghita Houir Alami. 2024.	773
722	Jack Cushman, Matthew Dahl, and Michael Lissner.	Legalbench-rag: A benchmark for retrieval-	774
723	2021. <i>eyecite: A tool for parsing legal citations.</i>	augmented generation in the legal domain. <i>arXiv</i>	775
724	<i>Journal of Open Source Software,</i> 6(66):3617.	<i>preprint arXiv:2408.10343</i> .	776
725	José Jorge de Queiroz Santos Filho, Filipe Araújo Dan-	Alec Radford, Jong Wook Kim, Chris Hallacy, Aditya	777
726	tas, Melquezedeqe da Silva Lima, Shirley Barbosa	Ramesh, Gabriel Goh, Sandhini Agarwal, Girish Sas-	778
727	dos Santos, Galileu Genesis, Maria Gabriely Lima	try, Amanda Askell, Pamela Mishkin, Jack Clark, and	779
728	da Salva, Álvaro Farias Pinheiro, and Erayl-	1 others. 2021. Learning transferable visual models	780
729	son Galdino da Silva. 2025. Comparing machine	from natural language supervision. In <i>International</i>	781
730	learning and an expert system for legal document	<i>conference on machine learning,</i> pages 8748–8763.	782
731	classification. In <i>Conference on Digital Government</i>	PmLR.	783
732	<i>Research,</i> volume 26.	Stephen Robertson and Hugo Zaragoza. 2009. <i>The prob-</i>	784
733	Maura R Grossman and Gordon V Cormack. 2010.	<i>abilistic relevance framework: BM25 and beyond,</i>	785
734	Technology-assisted review in e-discovery can be	volume 4. Now Publishers Inc.	786
735	more effective and more efficient than exhaustive	Victor Sanh, Lysandre Debut, Julien Chaumond, and	787
736	manual review. <i>Rich. JL &amp; Tech.,</i> 17:1.	Thomas Wolf. 2019. Distilbert, a distilled version	788
737	Maura R. Grossman, Gordon V. Cormack, Bruce Hedin,	of bert: smaller, faster, cheaper and lighter. <i>ArXiv,</i>	789
738	and Douglas W. Oard. 2011. Overview of the trec	abs/1910.01108.	790
739	2011 legal track. In <i>Proceedings of the Text REtrieval</i>	Murali Saravanan, Balaraman Ravindran, and S Ra-	791
740	<i>Conference (TREC)</i> . National Institute of Standards	man. 2008. Automatic identification of rhetorical	792
741	and Technology (NIST).	roles using conditional random fields for legal doc-	793
742	BJ Ismaylova. 2024. Problems of admissibility and	ument summarization. In <i>Proceedings of the Third</i>	794
743	reliability of metadata as evidence. international jour-	<i>International Joint Conference on Natural Language</i>	795
744	nal of. <i>Law and Policy,</i> 2(8):1.	<i>Processing: Volume-I.</i>	796
745	Armand Joulin, Edouard Grave, Piotr Bojanowski, and	Fabrizio Sebastiani. 2002. Machine learning in auto-	797
746	Tomáš Mikolov. 2017. Bag of tricks for efficient text	mated text categorization. In <i>ACM Computing Sur-</i>	798
747	classification. In <i>Proceedings of the 15th conference</i>	<i>veys,</i> volume 34, pages 1–47. ACM.	799
748	<i>of the European chapter of the association for com-</i>	Karen Sparck Jones. 1972. A statistical interpretation	800
749	<i>putational linguistics: volume 2, short papers,</i> pages	of term specificity and its application in retrieval.	801
750	427–431.	<i>Journal of documentation,</i> 28(1):11–21.	802
751	Daniel Martin Katz, Michael J Bommarito II, and Josh	Samir Undavia, Adam Meyers, and John E Ortega. 2018.	803
752	Blackman. 2017. A general approach for predicting	A comparative study of classifying legal documents	804
753	the behavior of the supreme court of the united states.	with neural networks. In <i>2018 Federated conference</i>	805
754	<i>PloS one,</i> 12(4):e0174698.	<i>on computer science and information systems (FedC-</i>	806
755	Zhijian Li, Stefan Larson, and Kevin Leach. 2025. Doc-	<i>SIS),</i> pages 515–522. IEEE.	807
756	ument classification using file names. In <i>Proceedings</i>	Qiqi Wang, Kaiqi Zhao, Robert Amor, Benjamin Liu,	808
757	<i>of the 2025 ACM Symposium on Document Engineer-</i>	and Ruofan Wang. 2022. D2gclf: Document-to-	809
758	<i>ing,</i> pages 1–10.	graph classifier for legal document classification. In	810
		<i>Findings of the Association for Computational Lin-</i>	811
		<i>guistics: NAACL 2022,</i> pages 2208–2221.	812
			813

814	Wenhui Wang, Furu Wei, Li Dong, Hangbo Bao, Nan	<b>A.4 Penalty Exponent</b>	865
815	Yang, and Ming Zhou. 2020. <a href="#">MiniLM: Deep self-</a>	The vocabulary overlap that motivates LexGLUE’s	866
816	<a href="#">attention distillation for task-agnostic compression</a>	higher penalty exponent comes from multi-	867
817	<a href="#">of pre-trained transformers</a> . In <i>Proceedings of the</i>	ple sources: shared legal terminology across	868
818	<i>17th International Conference on Machine Learning</i>	court opinions and legislation, common sentence-	869
819	<i>(ICML)</i> .	structural patterns (‘the court holds’, ‘pursuant to’),	870
820	Joe Watson, Guy Aglionby, and Samuel March. 2023.	and overlapping stylistic register across formal le-	871
821	Using machine learning to create a repository of judg-	gal writing. Without aggressive penalisation, these	872
822	ments concerning a new practice area: a case study	shared terms produce high CSM scores despite be-	873
823	in animal protection law. <i>Artificial Intelligence and</i>	ing non-discriminative. TREC’s smaller, narrower	874
824	<i>Law</i> , 31(2):293–324.	corpus has less such overlap, and an excess penalty	875
825	Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang	instead discards genuinely relevant but rarer terms	876
826	Zhang, Zhiyuan Liu, and Maosong Sun. 2020. How	<b>A.5 XGBoost Tree Depth</b>	877
827	does nlp benefit legal system: A summary of legal	The underfit-overfit behaviour across depths is	878
828	artificial intelligence. In <i>Proceedings of the 58th</i>	visible in Table 9: at depth 2 to 3, the model	879
829	<i>annual meeting of the association for computational</i>	fires too broadly across keyword dimensions, pro-	880
830	<i>linguistics</i> , pages 5218–5230.	ducing lower precision; at depth 6 to 12, the	881
831	<b>A Detailed Ablation Analysis</b>	trees memorise the training noise on TREC’s	882
832	<b>A.1 Lemmatisation and Stemming</b>	2101-document training set rather than general-	883
833	The loss of semantic meaning from lemmatisation	ising, causing precision to decrease by 1 to 2%.	884
834	and stemming is concrete in legal text. For ex-	LexGLUE plateaus between depths 3 and 6, con-	885
835	ample, ‘futures’ and ‘future’ are both lemmatised	sistent with its larger corpus and stronger stylistic	886
836	to ‘futur’; however, ‘futures’ is a domain-specific	separation between classes. Extended per-ablation	887
837	term in securities trading cases, whereas ‘future’	analysis with worked examples appears in Ap-	888
838	could refer to plans or an event in a news report.	pendix A.	889
839	Similarly, stemming collapses ‘proceeding’ (a legal	<b>A.6 Top-3 Keyword Glosses</b>	890
840	noun indicating regulatory action) and ‘proceed’ (a	The top three keywords by frequency in each cor-	891
841	general verb) into the same form. These collisions	pus relate to legally salient entities and concepts	892
842	remove precisely the morphological distinctions	as follows. For LexGLUE: ‘remanded’ refers to	893
843	CSM uses to separate relevant from irrelevant doc-	a higher court sending a case back to a lower	894
844	uments.	court for further proceedings, and is characteris-	895
845	<b>A.2 N-gram Range</b>	tic of appellate opinions; ‘testimony’ refers to ev-	896
846	An additional experiment with 15-grams was con-	idence given by a witness under oath, central to	897
847	ducted to test for edge cases of very long repeated	court records; ‘congressional’ relates to the United	898
848	sentences (e.g., boilerplate clauses or signature	States Congress, frequently cited in SCOTUS op-	899
849	blocks). As shown in the final row of Table 5,	inions. For TREC: ‘CME’ refers to the Chicago	900
850	performance is statistically indistinguishable from	Mercantile Exchange, a major US derivatives mar-	901
851	the 4-gram and 5-gram configurations, indicating	ketplace; ‘CBOE’ refers to the Chicago Board Op-	902
852	that very long phrases offer no additional signal	tions Exchange, a major US options market; ‘Mini’	903
853	beyond what mid-length n-grams already capture.	refers to the E-mini futures contract, a class of	904
854	<b>A.3 Minimum Term Frequency</b>	electronically-traded derivatives. All six terms are	905
855	The topic-specific terms that drive TREC classifi-	highly diagnostic of the relevant document classes	906
856	cation at MTF 10 include regulator acronyms (e.g.,	in their respective corpora.	907
857	‘CFTC’, ‘SEC’), trading-instrument names (e.g.,	<b>B Manual Keyword Baseline</b>	908
858	‘CME’, ‘CBOE’, ‘E-mini’), and Enron-internal	The following keywords were used with an OR-	909
859	project codenames. These appear in only a small	match rule: a document is classified as relevant if	910
860	fraction of the relevant set but are highly discrim-	any keyword appears in the text.	911
861	inative when present. Raising MTF above 10		
862	progressively discards these terms, explaining the		
863	steep accuracy decline on TREC visible in Table 6		
864	from MTF 100 onward.		

912 *applicant, court, judgment, petition, cer-*  
913 *tiorari, affirmed, reversed, held, constitu-*  
914 *tional, amendment, justice, appeal, ver-*  
915 *dict, plaintiff, defendant, proceedings,*  
916 *statute, ruling, tribunal, conviction*

917 **TREC Legal Track**

918 **Topic 401** (EnronOnline / Online Financial Trading  
919 Systems):

920 *enron online, enrononline, eon, online*  
921 *trading, online system, trading system,*  
922 *electronic trading, online platform, trad-*  
923 *ing platform*

924 **Topic 402** (Legality of OTC Derivatives and Finan-  
925 cial Instrument Trading):

926 *derivative, derivatives, otc, over the*  
927 *counter, over-the-counter, financial in-*  
928 *strument, swap, swaps, futures contract,*  
929 *options contract, hedge, hedging, no-*  
930 *tional, credit default*

931 All keywords were selected by a subject matter  
932 expert without consulting the training data. Clas-  
933 sification uses an OR rule: a document is labelled  
934 relevant if any keyword appears verbatim in the  
935 document text.