

LexGLUE: A Benchmark Dataset for Legal Language Understanding in English

Ilias Chalkidis^{α*} Abhik Jana^β Dirk Hartung^{γ δ} Michael Bommarito^{γ δ}
Ion Androutsopoulos^ε Daniel Martin Katz^{γ δ ζ} Nikolaos Aletras^η

^α University of Copenhagen, Denmark ^β Universität Hamburg, Germany

^γ Bucerius Law School, Hamburg, Germany ^δ CodeX, Stanford Law School, United States

^ε Athens University of Economics and Business, Greece ^η University of Sheffield, UK

^ζ Illinois Tech – Chicago Kent College of Law, United States

Abstract

Laws and their interpretations, legal arguments and agreements are typically expressed in writing, leading to the production of vast corpora of legal text. Their analysis, which is at the center of legal practice, becomes increasingly elaborate as these collections grow in size. Natural language understanding (NLU) technologies can be a valuable tool to support legal practitioners in these endeavors. Their usefulness, however, largely depends on whether current state-of-the-art models can generalize across various tasks in the legal domain. To answer this currently open question, we introduce the Legal General Language Understanding Evaluation (LexGLUE) benchmark, a collection of datasets for evaluating model performance across a diverse set of legal NLU tasks in a standardized way. We also provide an evaluation and analysis of several generic and legal-oriented models demonstrating that the latter consistently offer performance improvements across multiple tasks.

1 Introduction

Law is a field of human endeavor dominated by the use of language. As part of their professional training, law students consume large bodies of text as they seek to tune their understanding of the law and its application to help manage human behavior. Virtually every modern legal system produces massive volumes of textual data (Katz et al., 2020). Lawyers, judges, and regulators continuously author legal documents such as briefs, memos, statutes, regulations, contracts, patents and judicial decisions (Coupette et al., 2021). Beyond the consumption and production of language, law and the art of lawyering is also an exercise centered around the analysis and interpretation of text.

Natural language understanding (NLU) technologies can assist legal practitioners in a variety of legal tasks (Chalkidis and Kampas, 2018; Aletras

THE LEGAL NLP BENCHMARK



Figure 1: LexGLUE: A new benchmark dataset to evaluate the capabilities of NLU models on legal text.

et al., 2019, 2020; Zhong et al., 2020b; Bommarito et al., 2021), from judgment prediction (Aletras et al., 2016; Sim et al., 2016; Katz et al., 2017; Zhong et al., 2018; Chalkidis et al., 2019a; Malik et al., 2021), information extraction from legal documents (Chalkidis et al., 2018, 2019c; Chen et al., 2020; Hendrycks et al., 2021) and case summarization (Bhattacharya et al., 2019) to legal question answering (Ravichander et al., 2019; Kien et al., 2020; Zhong et al., 2020a,c) and text classification (Nallapati and Manning, 2008; Chalkidis et al., 2019b, 2020a). Transformer models (Vaswani et al., 2017) pre-trained on legal, rather than generic, corpora have also been studied (Chalkidis et al., 2020b; Zheng et al., 2021; Xiao et al., 2021).

Pre-trained Transformers, including BERT (Devlin et al., 2019), GPT-3 (Brown et al., 2020), T5 (Raffel et al., 2020), BART (Lewis et al., 2020), DeBERTa (He et al., 2021) and numerous variants, are currently the state of the art in most natural language processing (NLP) tasks. Rapid performance improvements have been witnessed, to the extent that ambitious multi-task benchmarks (Wang et al., 2018, 2019b) are considered almost ‘solved’ a few years after their release and need to be made more challenging (Wang et al., 2019a).

* Corresponding author: ilias.chalkidis@di.ku.dk

Recently, Bommasani et al. (2021) named these pre-trained models (e.g., BERT, DALL-E, GPT-3) *foundation models*. The term may be controversial, but it emphasizes the paradigm shift these models have caused and their interdisciplinary potential. Studying the latter includes the question of how to adapt these models to legal text (Bommarito et al., 2021). As discussed by Zhong et al. (2020b) and Chalkidis et al. (2020b), legal text has distinct characteristics, such as terms that are uncommon in generic corpora (e.g., ‘restrictive covenant’, ‘promissory estoppel’, ‘tort’, ‘novation’), terms that have different meanings than in everyday language (e.g., an ‘executed’ contract is signed and effective, a ‘party’ is a legal entity), older expressions (e.g., pronominal adverbs like ‘herein’, ‘hereto’, ‘wherefore’), uncommon expressions from other languages (e.g., ‘laches’, ‘voir dire’, ‘certiorari’, ‘sub judice’), and long sentences with unusual word order (e.g., “the provisions for termination hereinafter appearing or will at the cost of the borrower forthwith comply with the same”) to the extent that legal language is often classified as a ‘sub-language’ (Tiersma, 1999; Williams, 2007; Haigh, 2018). Furthermore, legal documents are often much longer than the maximum length state-of-the-art deep learning models can handle, including those designed to handle long text (Beltagy et al., 2020; Zaheer et al., 2020; Yang et al., 2020).

Inspired by the recent widespread use of the GLUE multi-task benchmark NLP dataset (Wang et al., 2018, 2019b), the subsequent more difficult SuperGLUE (Wang et al., 2019a), other previous multi-task NLP benchmarks (Conneau and Kiela, 2018; McCann et al., 2018), and similar initiatives in other domains (Peng et al., 2019), we introduce LexGLUE, a benchmark dataset to evaluate the performance of NLP methods in legal tasks. LexGLUE is based on seven English existing legal NLP datasets, selected using criteria largely from SuperGLUE (discussed in Section 3.1).

We anticipate that more datasets, tasks, and languages will be added in later versions of LexGLUE.¹ As more legal NLP datasets become available, we also plan to favor datasets checked thoroughly for validity (scores reflecting real-life performance), annotation quality, statistical power, and social bias (Bowman and Dahl, 2021).

As in GLUE and SuperGLUE (Wang et al.,

2019b,a), one of our goals is to push towards generic (or ‘foundation’) models that can cope with multiple NLP tasks, in our case legal NLP tasks, possibly with limited task-specific fine-tuning. Another goal is to provide a convenient and informative entry point for NLP researchers and practitioners wishing to explore or develop methods for legal NLP. Having these goals in mind, the datasets we include in LexGLUE and the tasks they address have been simplified in several ways, discussed below, to make it easier for newcomers and generic models to address all tasks. We provide Python APIs integrated with Hugging Face (Wolf et al., 2020; Lhoest et al., 2021) to easily import all the datasets we experiment with and evaluate the performance of different models (Section 4.4).

By unifying and facilitating the access to a set of law-related datasets and tasks, we hope to attract not only more NLP experts, but also more interdisciplinary researchers (e.g., law doctoral students willing to take NLP courses). More broadly, we hope LexGLUE will speed up the adoption and transparent evaluation of new legal NLP methods and approaches in the commercial sector, too. Indeed, there have been many commercial press releases in the legal tech industry on high-performing systems, but almost no independent evaluation of the performance of machine learning and NLP-based tools. A standard publicly available benchmark would also allay concerns of undue influence in predictive models, including the use of metadata which the relevant law expressly disregards.

2 Related Work

The rapid growth of the legal text processing field is demonstrated by numerous papers presented in top-tier conferences in NLP and artificial intelligence (Luo et al., 2017; Zhong et al., 2018; Chalkidis et al., 2019a; Valvoda et al., 2021) as well as surveys (Chalkidis and Kampas, 2018; Zhong et al., 2020b; Bommarito et al., 2021). Moreover, specialized workshops on NLP for legal text (Aletras et al., 2019; Di Fatta et al., 2020; Aletras et al., 2020) are regularly organized.

A core task in this area has been legal judgment prediction (forecasting), where the goal is to predict the outcome (verdict) of a court case. In this direction, there have been at least three lines of work. The first one (Aletras et al., 2016; Chalkidis et al., 2019a; Medvedeva et al., 2020, 2021) predicts violations of human rights in cases of the

¹See <https://nllpw.org/resources/> and <https://github.com/thunlp/LegalPapers> for lists of papers, datasets, and other resources related to NLP for legal text.

European Court of Human Rights (ECtHR). The second line of work (Luo et al., 2017; Zhong et al., 2018; Yang et al., 2019) considers Chinese criminal cases where the goal is to predict relevant law articles, criminal charges, and the term of the penalty. The third line of work (Ruger et al., 2004; Katz et al., 2017; Kaufman et al., 2019) includes methods for predicting the outcomes of cases of the Supreme Court of the United States (SCOTUS).

The same or similar tasks have also been studied with court cases in many other jurisdictions including France (Şulea et al., 2017), Philippines (Virtuicio et al., 2018), Turkey (Mumcuoğlu et al., 2021), Thailand (Kowsrihawatt et al., 2018), United Kingdom (Strickson and De La Iglesia, 2020), Germany (Urchs et al., 2021), and Switzerland (Niklaus et al., 2021). Apart from predicting court decisions, there is also work aiming to interpret (explain) the decisions of particular courts (Ye et al., 2018; Chalkidis et al., 2021c; Branting et al., 2021).

Another popular task is legal topic classification. Nallapati and Manning (2008) highlighted the challenges of legal document classification compared to more generic text classification by using a dataset including docket entries of US court cases. Chalkidis et al. (2020a) classify EU laws into EuroVoc concepts, a task earlier introduced by Mencia and Fürnkranz (2007), with a special interest in few- and zero-shot learning. Luz de Araujo et al. (2020) also studied topic classification using a dataset of Brazilian Supreme Court cases. There are similar interesting applications in contract law (Lippi et al., 2019; Tugener et al., 2020).

Several studies (Chalkidis et al., 2018, 2019c; Hendrycks et al., 2021) explored information extraction from contracts, to extract important information such as the contracting parties, agreed payment amount, start and end dates, applicable law, etc. Other studies focus on extracting information from legislation (Cardellino et al., 2017; Angelidis et al., 2018) or court cases (Leitner et al., 2019).

Legal Question Answering (QA) is another task of interest in legal NLP, where the goal is to train models for answering legal questions (Kim et al., 2015; Ravichander et al., 2019; Kien et al., 2020; Zhong et al., 2020a,c; Louis and Spanakis, 2022). Not only is this task interesting for researchers but it could support efforts to help laypeople better understand their legal rights. In the general task setting, this requires identifying relevant legislation, case law, or other legal documents, and extracting

elements of those documents that answer a particular question. A notable venue for legal QA has been the Competition on Legal Information Extraction and Entailment (COLIEE) (Kim et al., 2016; Kano et al., 2017, 2018).

More recently, there have also been efforts to pre-train Transformer-based language models on legal corpora (Chalkidis et al., 2020b; Zheng et al., 2021; Xiao et al., 2021), leading to state-of-the-art results in several legal NLP tasks, compared to models pre-trained on generic corpora.

Overall, the legal NLP literature is overwhelming, and the resources are scattered. Documentation is often not available, and evaluation measures vary across articles studying the same task. Our goal is to create the first unified benchmark to access the performance of NLP models on legal NLU. As a first step, we selected a representative group of tasks, using datasets in English that are also publicly available, adequately documented and have an appropriate size for developing modern NLP methods. We also introduce several simplifications to make the new benchmark more standardized and easily accessible, as already noted.

3 LexGLUE Tasks and Datasets

We present the Legal General Language Understanding² Evaluation (LexGLUE) benchmark, a collection of datasets for evaluating model performance across a diverse set of legal NLU tasks.

3.1 Dataset Desiderata

The datasets of LexGLUE were selected to satisfy the following desiderata:

- **Language:** In this first version of LexGLUE, we only consider English datasets, which also makes experimentation easier for researchers across the globe. We hope to include other languages in future versions of LexGLUE.
- **Substance:**³ The datasets should check the ability of systems to understand and reason about legal text to a certain extent in order to perform tasks that are meaningful for legal practitioners.
- **Difficulty:** The performance of state-of-the-art methods on the datasets should leave large scope for improvements (cf. GLUE and SuperGLUE,

²The term ‘understanding’ is, of course, as debatable as in NLU and GLUE, but is commonly used in NLP to refer to systems that analyze, rather than generate text.

³We reuse this term from the work of Wang et al. (2019a).

Dataset	Source	Sub-domain	Task Type	Training/Dev/Test Instances	Classes
ECtHR (Task A)	Chalkidis et al. (2019a)	ECHR	Multi-label classification	9,000/1,000/1,000	10+1
ECtHR (Task B)	Chalkidis et al. (2021c)	ECHR	Multi-label classification	9,000/1,000/1,000	10+1
SCOTUS	Spaeth et al. (2020)	US Law	Multi-class classification	5,000/1,400/1,400	14
EUR-LEX	Chalkidis et al. (2021a)	EU Law	Multi-label classification	55,000/5,000/5,000	100
LEDGAR	Tuggenier et al. (2020)	Contracts	Multi-class classification	60,000/10,000/10,000	100
UNFAIR-ToS	Lippi et al. (2019)	Contracts	Multi-label classification	5,532/2,275/1,607	8+1
CaseHOLD	Zheng et al. (2021)	US Law	Multiple choice QA	45,000/3,900/3,900	n/a

Table 1: Statistics of the LexGLUE datasets, including simplifications made.

where top-ranked models now achieve average scores higher than 90%). Unlike SuperGLUE (Wang et al., 2019a), we did not rule out, but rather favored, datasets requiring domain (in our case legal) expertise.

- **Availability & Size:** We consider only publicly available datasets, documented by published articles, avoiding proprietary, untested, poorly documented datasets. We also excluded very small datasets, e.g., with fewer than 5K documents. Although large pre-trained models often perform well with relatively few task-specific training instances, newcomers may wish to experiment with simpler models that may perform disappointingly with small training sets. Small test sets may also lead to unstable and unreliable results.

3.2 Tasks and Datasets

LexGLUE comprises seven datasets. Table 1 shows core information for each of the LexGLUE datasets and tasks, described in detail below.⁴

ECtHR Tasks A & B The European Court of Human Rights (ECtHR) hears allegations that a state has breached human rights provisions of the European Convention of Human Rights (ECHR). We use the dataset of Chalkidis et al. (2019a, 2021c), which contains approx. 11K cases from the ECtHR public database. The cases are chronologically split into training (9k, 2001–2016), development (1k, 2016–2017), and test (1k, 2017–2019). For each case, the dataset provides a list of *factual* paragraphs (facts) from the case description. Each case is mapped to *articles* of the ECHR that were violated (if any). In Task A, the input to a model is the list of facts of a case, and the output is the set of violated articles. In the most recent version of the dataset (Chalkidis et al., 2021c), each case is also mapped to articles of ECHR that were *allegedly* violated (considered by the court). In Task B, the input is again the list of facts of a case, but the output is the set of allegedly violated articles.

The total number of ECHR articles is currently 66. Several articles, however, cannot be violated, are rarely (or never) discussed in practice, or do not depend on the facts of a case and concern procedural technicalities. Thus, we use a simplified version of the label set (ECHR articles) in both Task A and B, including only 10 ECHR articles that can be violated and depend on the case’s facts.

SCOTUS The US Supreme Court (SCOTUS)⁵ is the highest federal court in the United States of America and generally hears only the most controversial or otherwise complex cases which have not been sufficiently well solved by lower courts. We release a new dataset combining information from SCOTUS opinions⁶ with the Supreme Court DataBase (SCDB)⁷ (Spaeth et al., 2020). SCDB provides metadata (e.g., decisions, issues, decision directions) for all cases (from 1946 up to 2020). We opted to use SCDB to classify the court *opinions* in the available 14 *issue areas* (e.g., Criminal Procedure, Civil Rights, Economic Activity, etc.). This is a single-label multi-class classification task (Table 1). The 14 issue areas cluster 278 issues whose focus is on the subject matter of the controversy (dispute). The SCOTUS cases are chronologically split into training (5k, 1946–1982), development (1.4k, 1982–1991), test (1.4k, 1991–2016) sets.

EUR-LEX European Union (EU) legislation is published in the EUR-Lex portal.⁸ All EU laws are annotated by EU’s Publications Office with multiple concepts from EuroVoc, a multilingual thesaurus maintained by the Publications Office.⁹ The current version of EuroVoc contains more than 7k concepts referring to various activities of the EU and its Member States (e.g., economics, health-care, trade). We use the English part of the dataset of Chalkidis et al. (2021a), which comprises 65k EU laws (documents) from EUR-Lex. Given a

⁴In Appendix G, we provide examples, i.e., pairs of (inputs, outputs), for all datasets and tasks.

⁵<https://www.supremecourt.gov>

⁶<https://www.courtlistener.com>

⁷<http://scdb.wustl.edu>

⁸<http://eur-lex.europa.eu/>

⁹<http://eurovoc.europa.eu/>

Method	Source	# Params	Vocab. Size	Max Length	Pretrain Specs	Pre-training Corpora
BERT	(Devlin et al., 2019)	110M	32K	512	1M / 256	(16GB) Wiki, BC
RoBERTa	(Liu et al., 2019)	125M	50K	512	100K / 8K	(160GB) Wiki, BC, CC-News, OWT
DeBERTa	(He et al., 2021)	139M	50K	512	1M / 256	(160GB) Wiki, BC, CC-News, OWT
Longformer*	(Beltagy et al., 2020)	149M	50K	4096	65K / 64	(160GB) Wiki, BC, CC-News, OWT
BigBird*	(Zaheer et al., 2020)	127M	50K	4096	1M / 256	(160GB) Wiki, BC, CC-News, OWT
Legal-BERT	(Chalkidis et al., 2020b)	110M	32K	512	1M / 256	(12GB) Legislation, Court Cases, Contracts
CaseLaw-BERT	(Zheng et al., 2021)	110M	32K	512	2M / 256	(37GB) US Court Cases

Table 2: Key specifications of the examined models. We report the number of parameters, the size of vocabulary, the maximum sequence length, the core pre-training specifications (training steps and batch size), and the training corpora (OWT = OpenWebText, BC = BookCorpus). Starred models have been warm-started from RoBERTa.

document, the task is to predict its EuroVoc labels (concepts). The dataset is chronologically split in training (55k, 1958–2010), development (5k, 2010–2012), test (5k, 2012–2016) subsets. It supports four different label granularities, comprising 21, 127, 567, 7390 EuroVoc concepts, respectively. We use the 100 most frequent concepts from level 2, which has a highly skewed label distribution and temporal concept drift (Chalkidis et al., 2021a), making it sufficiently difficult for an entry point.

LEDGAR Tuggener et al. (2020) introduced LEDGAR (Labeled EDGAR), a dataset for contract provision (paragraph) classification. The contract provisions come from contracts obtained from the US Securities and Exchange Commission (SEC) filings, which are publicly available from EDGAR¹⁰ (Electronic Data Gathering, Analysis, and Retrieval system). The original dataset includes approx. 850k contract provisions labeled with 12.5k categories. Each label represents the single main topic (theme) of the corresponding contract provision, i.e., this is a single-label multi-class classification task. In LexGLUE, we use a subset of the original dataset with 80k contract provisions, considering only the 100 most frequent categories as a simplification. We split the new dataset chronologically into training (60k, 2016–2017), development (10k, 2018), and test (10k, 2019) sets.

UNFAIR-ToS The UNFAIR-ToS dataset (Lippi et al., 2019) contains 50 Terms of Service (ToS) from on-line platforms (e.g., YouTube, Ebay, Facebook, etc.). The dataset has been annotated on the sentence-level with 8 types of *unfair contractual terms*, meaning terms (sentences) that potentially violate user rights according to EU consumer law.¹¹ The input to a model is a sentence, the output is the set of unfair types (if any). We split the dataset chronologically into training (5.5k, 2006–2016), development (2.3k, 2017), test (1.6k, 2017) sets.

¹⁰<https://www.sec.gov/edgar/>

¹¹Art. 3 of Direct. 93/13, Unfair Terms in Consumer Contracts (<http://data.europa.eu/eli/dir/1993/13/oj>).

CaseHOLD The CaseHOLD (Case Holdings on Legal Decisions) dataset (Zheng et al., 2021) contains approx. 53k multiple choice questions about holdings of US court cases from the Harvard Law Library case law corpus. *Holdings* are short summaries of legal rulings that accompany referenced decisions relevant for the present case, e.g.:

“... to act pursuant to City policy, re d 503, 506-07 (3d Cir.1985)(**holding that for purposes of a class certification motion the court must accept as true all factual allegations in the complaint and may draw reasonable inferences therefrom**).”

The input consists of an *excerpt* (or prompt) from a court decision, containing a reference to a particular case, where the *holding* statement (in boldface) is masked out. The model must identify the correct (masked) holding statement from a selection of five choices. We split the dataset in training (45k), development (3.9k), test (3.9k) sets, excluding samples that are shorter than 256 tokens. Chronological information is missing from CaseHOLD, thus we cannot perform a chronological re-split.

4 Models Considered

4.1 Linear SVM

Our first baseline model is a linear Support Vector Machine (SVM) (Cortes and Vapnik, 1995) with TF-IDF features for the top- K frequent n -grams of the training set, where $n \in [1, 2, 3]$.

4.2 Pre-trained Transformer Models

We experiment with Transformer-based (Vaswani et al., 2017) pre-trained language models, which achieve state of the art performance in most NLP tasks (Bommasani et al., 2021) and NLU benchmarks (Wang et al., 2019a). These models are pre-trained on very large unlabeled corpora to predict masked tokens (masked language modeling) and typically also to perform other pre-training tasks that still do not require any manual annotation (e.g., predicting if two sentences were adjacent in the corpus or not, dubbed next sentence prediction).

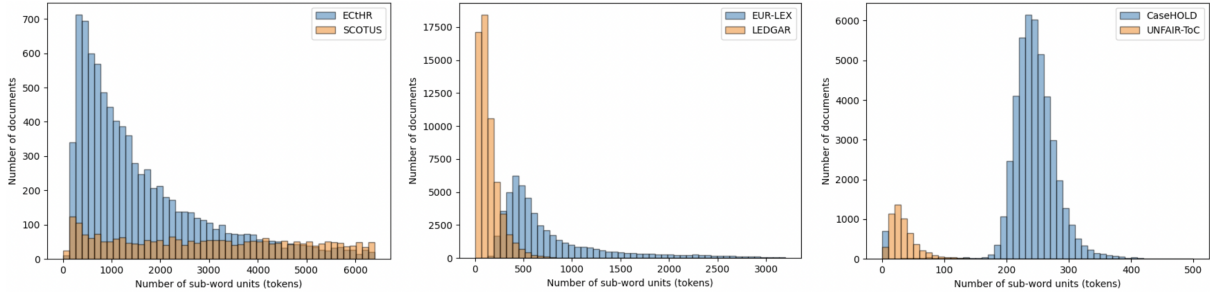


Figure 2: Distribution of text input length, measured in BERT sub-word units, across LexGLUE datasets.

The pre-trained models are then fine-tuned (further trained) on task-specific (typically much smaller) annotated datasets, after adding task-specific layers. We fine-tune and evaluate the performance of the following publicly available models (Table 2).

BERT (Devlin et al., 2019) is the best-known pre-trained Transformer-based language model. It is pre-trained to perform masked language modeling and next sentence prediction.

RoBERTa (Liu et al., 2019) is also a pre-trained Transformer-based language model. Unlike BERT, RoBERTa uses dynamic masking, it eliminates the next sentence prediction pre-training task, uses a larger vocabulary, and has been pre-trained on much larger corpora. Liu et al. (2019) reported improved results on NLU benchmarks using RoBERTa, compared to BERT.

DeBERTa (He et al., 2021) is another improved BERT model that uses disentangled attention, i.e., four separate attention mechanisms considering the content and the relative position of each token, and an enhanced mask decoder, which explicitly considers the absolute position of the tokens. DeBERTa has been reported to outperform BERT and RoBERTa in several NLP tasks (He et al., 2021).

Longformer (Beltagy et al., 2020) extends Transformer-based models to support longer sequences, using sparse-attention. The latter is a combination of local (window-based) attention and global (dilated) attention that reduces the computational complexity of the model and thus can be deployed in longer documents (up to 4096 tokens). Longformer outperforms RoBERTa on long document tasks and QA benchmarks.

BigBird (Zaheer et al., 2020) is another sparse-attention based transformer that uses a combination of a local (window-based) attention, global (dilated), and random attention, i.e., all tokens also attend a number of random tokens on top of those

in the same neighborhood (window) and the global ones. BigBird has been reported to outperform Longformer on QA and summarization tasks.

Legal-BERT (Chalkidis et al., 2020b) is a BERT model pre-trained on English legal corpora, consisting of legislation, contracts, and court cases. It uses the original pre-training BERT configuration. The sub-word vocabulary of Legal-BERT is built from scratch, to better support legal terminology.

CaseLaw-BERT (Zheng et al., 2021) is another law-specific BERT model. It also uses the original pre-training BERT configuration and has been pre-trained from scratch on the Harvard Law case corpus,¹² which comprises 3.4M legal decisions from US federal and state courts. This model is called *Custom Legal-BERT* by Zheng et al. (2021). We call it CaseLaw-BERT to distinguish it from the previously published Legal-BERT of Chalkidis et al. (2020b) and to better signal that it is trained exclusively on case law (court opinions).

Hierarchical Variants Legal documents are usually much longer (i.e., consisting of thousands of words) than other text types (e.g., tweets, customer reviews, news articles) often considered in various NLP tasks. Thus, standard Transformer-based models that can typically process up to 512 sub-word units cannot be directly applied across all LexGLUE datasets, unless documents are severely truncated to the model’s limit. Figure 2 shows the distribution of text input length across all LexGLUE datasets. Even for Transformer-based models specifically designed to handle long text (e.g., Longformer, BigBird), handling longer legal documents remains a challenge.

Given the length of the text input in three of the seven LexGLUE tasks, i.e., ECTHR (A and B) and SCOTUS, we employ a hierarchical variant of each pre-trained Transformer-based model that has not been designed for longer text (BERT, RoBERTa,

¹²<https://case.law/>

DeBERTa, Legal-BERT, CaseLaw-BERT) during fine-tuning and inference. The hierarchical variants are similar to those of Chalkidis et al. (2021c). They use the corresponding pre-trained Transformer-based model to encode each paragraph of the input text independently and obtain the top-level representation $h_{[cls]}$ of each paragraph. A second-level shallow (2-layered) Transformer encoder with always the same (across BERT, RoBERTa, DeBERTa etc.) specifications (e.g., hidden units, number of attention heads) is fed with the paragraph representations to make them context-aware (aware of the surrounding paragraphs). We then max-pool over the context-aware paragraph representations to obtain a document representation, which is fed to a classification layer.¹³

4.3 Task-Specific Fine-Tuning

Text Classification Tasks For EUR-LEX, LEDGAR and UNFAIR-ToS tasks, we feed each document to the pre-trained model (e.g., BERT) and obtain the top-level representation $h_{[cls]}$ of the special [cls] token as the document representation, following Devlin et al. (2019). The latter goes through a dense layer of L output units, one per label, followed by a sigmoid (in EUR-LEX, UNFAIR-ToS) or softmax (in LEDGAR) activation, respectively. For the two ECtHR tasks (A and B) and SCOTUS, where the hierarchical variants are employed, we feed the max-pooled (over paragraphs) document representation to a classification linear layer. The linear layer is again followed by a sigmoid (ECtHR) or softmax (SCOTUS) activation.

Multiple-Choice QA Task For CaseHOLD, we convert each training (or test) instance (the prompt and the five candidate answers) into five input pairs following Zheng et al. (2021). Each pair consists of the prompt and one of the five candidate answers, separated by the special delimiter token [sep]. The top-level representation $h_{[cls]}$ of each pair is fed to a linear layer to obtain a logit, and the five logits are then passed through a softmax yielding a probability distribution over the five candidate answers.

4.4 Data Repository and Code

For reproducibility purposes and to facilitate future experimentation with other models, we pre-process

and release all datasets on Hugging Face Datasets (Lhoest et al., 2021).¹⁴ We also release the code¹⁵ of our experiments, which relies on the Hugging Face Transformers (Wolf et al., 2020) library.¹⁶ Appendix A explains how to load the datasets and run experiments with our code.

5 Experiments

5.1 Experimental Set Up

For TFIDF-based linear SVM models, we use the implementation of Scikit-learn (Pedregosa et al., 2011) and grid-search for hyper parameters (number of features, C , and loss function). For all the pre-trained models, we use publicly available Hugging Face checkpoints.¹⁷ We use the *-base configuration of each pre-trained model, i.e., 12 Transformer blocks, 768 hidden units, and 12 attention heads. We train models with the Adam optimizer (Kingma and Ba, 2015) and an initial learning rate of $3e-5$ up to 20 epochs using early stopping on development data. We use mixed precision (fp16) to decrease the memory footprint in training and gradient accumulation for all hierarchical models. The hierarchical models can read up to 64 paragraphs of 128 tokens each. We use Longformer and BigBird in default settings, i.e., Longformer uses windows of 512 tokens and a single global token ([cls]), while BigBird uses blocks of 64 tokens (windows: $3 \times$ block, random: $3 \times$ block, global: $2 \times$ initial block; each token attends 512 tokens in total). The batch size is 8 in all experiments. We run five repetitions with different random seeds and report the test scores based on the seed with the best scores on development data. We evaluate performance using *micro-F1* (μ -F₁) and *macro-F1* (m-F₁) across all datasets to take into account class imbalance. For completeness, we also report the arithmetic, harmonic, and geometric mean across tasks following Shavrina and Malykh (2021).¹⁸

5.2 Experimental Results

Main Results Table 3 presents the test results for all models across all LexGLUE tasks, while Table 4

¹⁴https://huggingface.co/datasets/lex_glue

¹⁵<https://github.com/coastalcph/lex-glue>

¹⁶<https://huggingface.co/transformers>

¹⁷<http://huggingface.co/models>

¹⁸We acknowledge that the use of scores aggregated over tasks has been criticized in general NLU benchmarks (e.g., GLUE), as models are trained with different numbers of samples, task complexity, and evaluation metrics per task. We believe that the use of a standard common metric (F1) across tasks and averaging with harmonic mean alleviate this issue.

¹³In Appendix D, we present results from preliminary experiments using the standard version of BERT for ECtHR Task A (-12.2%), Task B (-10.6%), and SCOTUS (-3.5%).

Method	ECtHR (A)*		ECtHR (B)*		SCOTUS*		EUR-LEX		LEDGAR		UNFAIR-ToS		CaseHOLD
	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	
TFIDF-SVM	62.6	48.9	73.0	63.8	74.0	64.4	63.4	47.9	87.0	81.4	94.7	75.0	22.4
BERT	71.2	63.6	79.7	73.4	68.3	58.3	71.4	57.2	87.6	81.8	95.6	81.3	70.8
RoBERTa	69.2	59.0	77.3	68.9	71.6	62.0	71.9	57.9	87.9	82.3	95.2	79.2	71.4
DeBERTa	70.0	60.8	78.8	71.0	71.1	62.7	72.1	57.4	88.2	83.1	95.5	80.3	72.6
Longformer	69.9	64.7	79.4	71.7	72.9	64.0	71.6	57.7	88.2	83.0	95.5	80.9	71.9
BigBird	70.0	62.9	78.8	70.9	72.8	62.0	71.5	56.8	87.8	82.6	95.7	81.3	70.8
Legal-BERT	70.0	64.0	80.4	74.7	76.4	66.5	72.1	57.4	88.2	83.0	96.0	83.0	75.3
CaseLaw-BERT	69.8	62.9	78.8	70.3	76.6	65.9	70.7	56.6	88.3	83.0	96.0	82.3	75.4

Table 3: Test results for all examined models across LexGLUE tasks. In starred datasets, we use the hierarchical variant of each model, except for Longformer and BigBird, discussed in Section 4.2.

Method	A-Mean		H-Mean		G-Mean	
	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁
BERT	77.8	69.5	76.7	68.2	77.2	68.8
RoBERTa	77.8	68.7	76.8	67.5	77.3	68.1
DeBERTa	78.3	69.7	77.4	68.5	77.8	69.1
Longformer	78.5	70.5	77.5	69.5	78.0	70.0
BigBird	78.2	69.6	77.2	68.5	77.7	69.0
Legal-BERT	79.8	72.0	78.9	70.8	79.3	71.4
CaseLaw-BERT	79.4	70.9	78.5	69.7	78.9	70.3

Table 4: Test scores aggregated over tasks: arithmetic (A), harmonic (H), and geometric (G) mean.

presents the aggregated (averaged) results. We observe that the two legal-oriented pre-trained models (Legal-BERT, CaseLaw-BERT) perform overall better, especially considering m-F₁ that accounts for class imbalance (considers all classes equally important). Their in-domain (legal) knowledge seems to be more critical in the two datasets relying on US case law data (SCOTUS, CaseHOLD) with an improvement of approx. +2-4% p.p. (m-F₁) over equally sized Transformer-based models, which are pre-trained on generic corpora. These results are explained by the fact that these tasks are more domain-specific in terms of language, compared to the rest. No single model performs best in all tasks, and the results of Table 3 show that there is still large scope for improvement (Section 6).

An exceptional case of the dominance of the pre-trained Transformer models is the SCOTUS dataset, where the TFIDF-based linear SVM performs better than all generic Transformer models. TFIDF-SVM models are domain-specific, since the vocabulary (n-grams) and their IDF scores used to compute TF-IDF scores, are customized per task; which seems to be important for SCOTUS.

Legal-oriented Models Interestingly, the performance of Legal-BERT and CaseLaw-BERT, the two legal-oriented pre-trained models, is almost identical on CaseHOLD, despite the fact that

CaseLaw-BERT is solely trained on US case law. On the other hand, Legal-BERT has been exposed to a wider variety of legal corpora, including EU and UK legislation, ECtHR, ECJ and US court cases, and US contracts. Legal-BERT performs as well as or better than CaseLaw-BERT on all datasets. These results suggest that domain-specific pre-training (and learning a domain-specific sub-word vocabulary) is beneficial, but over-fitting a specific (niche) sub-domain (e.g., US case law), similarly to Zheng et al. (2021), has no benefits.

6 Vision – Future Considerations

Beyond the scope of this work and the examined baseline models, we identify four major factors that could potentially advance the state of the art in LexGLUE and legal NLP more generally:

Long Documents: Several Transformer-based models (Beltagy et al., 2020; Zaheer et al., 2020; Liu et al., 2022) have been proposed to handle long documents by exploring sparse attention mechanisms. These models can handle sequences up to 4096 sub-words, which is largely exceeded in three out of seven LexGLUE tasks (Figure 2). Contrary, the hierarchical model of Section 4.2 can handle sequences up to 8192 sub-words in our experiments, but a part of the model (the additional Transformer blocks that make the paragraph embeddings aware of the other paragraphs) is not pre-trained, which possibly negatively affects performance.

Structured Text: Current models for long documents, like Longformer and BigBird, do not consider the document structure (e.g., sentences, paragraphs, sections). For example, window-based attention may consider a sequence of sentences across paragraph boundaries or even consider truncated sentences. To exploit the document structure, Yang et al. (2020) proposed SMITH, a hierarchi-

cal Transformer model that hierarchically encodes increasingly larger blocks (e.g., words, sentences, documents). SMITH is very similar to the hierarchical model of Section 4.2, but it is pre-trained end-to-end with two objectives: token-level masked and sentence block language modeling.

Large-scale Legal Pre-training: Recent studies (Chalkidis et al., 2020b; Zheng et al., 2021; Bambroo and Awasthi, 2021; Xiao et al., 2021) introduced language models pre-trained on legal corpora, but of relatively small sizes, i.e., 12–36 GB. In the work of Zheng et al. (2021), the pre-training corpus covered only a narrowly defined area of legal documents, US court opinions. The same applies to Lawformer (Xiao et al., 2021), which was pre-trained on Chinese court opinions. Future work could curate and release a legal version of the C4 corpus (Raffel et al., 2020), containing multi-jurisdictional legislation, court decisions, contracts and legal literature at a size of hundreds of GBs. Given such a corpus, a large language model capable of processing long structured text could be pre-trained and it might excel in LexGLUE.

Even Larger Language Models: Scaling up the capacity of pre-trained models has led to increasingly better results in general NLU benchmarks (Kaplan et al., 2020), and models have been scaled up to billions of parameters (Brown et al., 2020; Raffel et al., 2020; He et al., 2021). In Appendix E, we observe that using the large version of RoBERTa leads to substantial performance improvements compared to the base version. The results are comparable or better - in some cases - compared to the legal-oriented language models (Legal-BERT, CaseLaw-BERT). Considering that the two legal-oriented models are much smaller and have been pre-trained with (5–10×) less data (Section 2), we have a strong indication for performance gains by pre-training larger legal-oriented models using larger legal corpora.

7 Limitations and Future Work

Although, our benchmark inevitably cannot cover “*everything in the whole wide (legal) world*” (Raji et al., 2021), we include a representative collection of English datasets that also ground to a certain degree in practically interesting applications.

In its current version, LexGLUE can only be used to evaluate English models. As legal documents are typically written in the official language

of the particular country of origin, there is an increasing need for developing models for other languages. The current scarcity of datasets in other languages (with the exception of Chinese) makes a multilingual extension of LexGLUE challenging, but an interesting avenue for future research.

Beyond language barriers, legal restrictions currently inhibit the creation of more datasets. Important document types, such as contracts and scholarly publications are protected by copyright or considered trade secrets. As a result, their owners are concerned with data-leakage when they are used for model training and evaluation. Providing both legal and technical solutions, e.g., using privacy-aware infrastructure and models (Downie, 2004; Feyisetan et al., 2020) is a challenge to be addressed.

Access to court decisions can also be hindered by bureaucratic inertia, outdated technology and data protection concerns, which collectively result in these otherwise public decisions not being publicly available (Pah et al., 2020). While the anonymization of personal data provides a solution to this problem, it is itself an open challenge for legal NLP (Jana and Biemann, 2021). In lack of suitable datasets and benchmarks, we have refrained from including anonymization in this version of LexGLUE, but plan to do so at a later stage.

Another limitation of the current version of LexGLUE is that human evaluation is missing. All datasets rely on *ground truth* labels automatically extracted from data (e.g., court decisions) produced as part of official judicial or archival procedures. These resources should be highly reliable (valid), but we cannot statistically assess their quality. In the future, re-annotating part of the datasets with multiple legal experts would provide an estimation of human level performance and inter-annotator agreement, though the cost would be high, because of the required legal expertise.

While LexGLUE offers a much needed unified testbed for legal NLU, there are several other critical aspects that need to be studied carefully. These include multi-disciplinary research to better understand the limitations and challenges of applying NLP to law (Binns, 2020), while also considering fairness and robustness (Angwin et al., 2016; Dressel and Farid, 2018; Baker Gillis, 2021; Wang et al., 2021; Chalkidis et al., 2022), and broader legal considerations of AI technologies in general (Schwemer et al., 2021; Tsarapatsanis and Aletras, 2021; Delacroix, 2022).

Acknowledgments

This work was partly funded by the Innovation Fund Denmark (IFD)¹⁹ under File No. 0175-00011A and by the German Federal Ministry of Education and Research (BMBF) *kmu-innovativ* program under funding code 01IS18085. We would like to thank Desmond Elliott for providing valuable feedback (baselines for truncated documents presented in Appendix D), Xiang Dai and Joel Niklaus for reviewing and pointing out issues in the new resources (code, datasets).

Ethics Statement

Original Work Attribution

All datasets included in LexGLUE, except SCOTUS, are publicly available and have been previously published. If datasets or the papers that introduced them were not compiled or written by ourselves, we referenced the original work and encourage LexGLUE users to do so as well. In fact, we believe this work should only be referenced, in addition to citing the original work, when experimenting with multiple LexGLUE datasets and using the LexGLUE evaluation infrastructure. Otherwise only the original work should be cited.

Social Impact

We believe that this work does not contain any grounds for ethical concerns. A transparent and rigorous benchmark for NLP in the legal domain might serve as an orientation for scholars and industry researchers. As a result, the capabilities of tools that are trained using natural language data from the legal domain will become clearer, thereby helping their users to better understand them. This increased certainty would also raise the awareness within research and industry communities to potential risks associated with the use of these tools. We regard this contribution to a more realistic, more informed discussion as an important use case of the work presented. Ideally, it could help both beginners and seasoned professionals to understand the limitations of using NLP tools in the legal domain and thereby prevent exaggerated expectations and potential applications that might risk endangering fundamental rights or the rule of law. We currently cannot imagine use cases of this particular work that would lead to ethical concerns or potential harm (Tsaratsanis and Aletras, 2021).

Licensing & Personal Information

LexGLUE comprises seven datasets: ECtHR Task A and B, SCOTUS, EUR-LEX, LEDGAR, UNFAIR-ToS, and CaseHOLD that are available for re-use and re-share with appropriate attribution. The data is in general partially anonymized in accordance with the applicable national law. The data is considered to be in the public sphere from a privacy perspective. This is a very sensitive matter, as the courts try to keep a balance between transparency (the public’s right to know) and privacy (respect for private and family life).

ECtHR contains personal data of the parties and other people involved in the legal proceedings. Its data is processed and made public in accordance with the European data protection laws. This includes either implied consent or legitimate interest to process the data for research purposes. As a result, their processing by us or other future users of the benchmark is not likely to raise ethical concerns.

SCOTUS contains personal data of a similar nature. Again, the data is processed and made available by the US Supreme Court, whose proceedings are public. While this ensures compliance with US law, it is very likely that similarly to the ECtHR any processing could be justified by either implied consent or legitimate interest under European law.

EUR-LEX by contrast is merely a collection of legislation material and therefore not likely to contain personal data, except signatory information (e.g., president of EC). It is openly published by the European Union and processed by the EU’s Publication Office. In addition, since our work qualifies as research, it is privileged pursuant to Art. 6 (1) (f) GDPR.

LEDGAR contains publicly available contract provisions published in the EDGAR database of the US Securities and Exchange Commission (SEC). As far as personal information might be contained, it should equally fall into the public sphere and be covered by research privilege. Our processing does not focus on personal information at all, rather attributing content labels to provisions.

UNFAIR-ToS contains Terms of Services from business entities such as YouTube, Ebay, Facebook, etc., which makes it unlikely for the data to include personal information. These companies keep user data separate from contractual provisions, so to the best of our knowledge not contained in this dataset.

CaseHOLD contains parts of legal decisions

¹⁹<https://innovationsfonden.dk/en>

from US Court decisions, obtained from the Harvard library case law corpus. All of the decisions were previously published in compliance with US law. In addition, most instances (case snippets) are too short to contain identifiable information. Should such data be contained, their processing would equally be covered either by implicit consent or a public interest exception. We use all datasets in accordance with copyright terms and under the licenses set forth by their creators.

Limitations & Potential Harms

We have not employed any crowd-workers or annotators for this work. The paper outlines the main limitations with regard to speaker population (English) and generalizability in a dedicated section (Section 7). As a benchmark paper, our claims naturally match the results of the experiments, which – given the current detail of instructions – should be easily reproduced. We provide several ways of accessing the datasets and running the experiments both with and without Hugging Face infrastructure.

We do not currently foresee any potential harms for vulnerable or marginalized populations and we do not use, to the best of our knowledge, any identifying characteristics for populations of these kinds.

References

- Nikolaos Aletras, Ion Androutsopoulos, Leslie Barrett, Adam Meyers, and Daniel Preotiuc-Pietro, editors. 2020. *Proceedings of the 2nd Natural Legal Language Processing Workshop at KDD 2020*. Online.
- Nikolaos Aletras, Elliott Ash, Leslie Barrett, Daniel Chen, Adam Meyers, Daniel Preotiuc-Pietro, David Rosenberg, and Amanda Stent, editors. 2019. *Proceedings of the 1st Natural Legal Language Processing Workshop at NAACL 2019*. Minneapolis, Minnesota.
- Nikolaos Aletras, Dimitrios Tsarapatsanis, Daniel Preotiuc-Pietro, and Vasileios Lampsos. 2016. [Predicting judicial decisions of the european court of human rights: A natural language processing perspective](#). *PeerJ Computer Science*, 2:e93.
- I. Angelidis, Ilias Chalkidis, and M. Koubarakis. 2018. [Named entity recognition, linking and generation for greek legislation](#). In *JURIX*, Groningen, The Netherlands.
- Julia Angwin, Jeff Larson, Surya Mattu, and Lauren Kirchner. 2016. [Machine bias: There’s software used across the country to predict future criminals. and it’s biased against blacks](#). *ProPublica*.
- Noa Baker Gillis. 2021. [Sexism in the judiciary: The importance of bias definition in NLP and in our courts](#). In *Proceedings of the 3rd Workshop on Gender Bias in Natural Language Processing*, pages 45–54, Online. Association for Computational Linguistics.
- Purbid Bambroo and Aditi Awasthi. 2021. [LegaldB: Long distilbert for legal document classification](#). In *2021 International Conference on Advances in Electrical, Computing, Communication and Sustainable Technologies (ICAECT)*, pages 1–4.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The long-document transformer](#). *CoRR*, abs/2004.05150.
- Paheli Bhattacharya, Kaustubh Hiware, Subham Rajgaria, Nilay Pochhi, Kripabandhu Ghosh, and Saptarshi Ghosh. 2019. A comparative study of summarization algorithms applied to legal case judgments. In *Advances in Information Retrieval*, pages 413–428, Cham. Springer International Publishing.
- Reuben Binns. 2020. [Analogies and disanalogies between machine-driven and human-driven legal judgement](#). *Journal of Cross-disciplinary Research in Computational Law*, 1(1).
- Michael J. Bommarito, Daniel Martin Katz, and Eric M. Detterman. 2021. [Lexnlp: Natural language processing and information extraction for legal and regulatory texts](#). *Research Handbook on Big Data Law*, pages 216–227.
- Rishi Bommasani, Drew A. Hudson, Ehsan Adeli, Russ Altman, Simran Arora, Sydney von Arx, Michael S. Bernstein, Jeannette Bohg, Antoine Bosselut, Emma Brunskill, Erik Brynjolfsson, Shyamal Buch, Dallas Card, Rodrigo Castellon, Niladri Chatterji, Annie Chen, Kathleen Creel, Jared Quincy Davis, Dora Demszky, Chris Donahue, Moussa Doumbouya, Esin Durmus, Stefano Ermon, John Etchemendy, Kawin Ethayarajh, Li Fei-Fei, Chelsea Finn, Trevor Gale, Lauren Gillespie, Karan Goel, Noah Goodman, Shelby Grossman, Neel Guha, Tatsunori Hashimoto, Peter Henderson, John Hewitt, Daniel E. Ho, Jenny Hong, Kyle Hsu, Jing Huang, Thomas Icard, Saahil Jain, Dan Jurafsky, Pratyusha Kalluri, Siddharth Karamcheti, Geoff Keeling, Fereshte Khani, Omar Khattab, Pang Wei Kohd, Mark Krass, Ranjay Krishna, Rohith Kudipudi, Ananya Kumar, Faisal Ladhak, Mina Lee, Tony Lee, Jure Leskovec, Isabelle Levent, Xiang Lisa Li, Xuechen Li, Tengyu Ma, Ali Malik, Christopher D. Manning, Suvir Mirchandani, Eric Mitchell, Zanele Munyikwa, Suraj Nair, Avani Narayan, Deepak Narayanan, Ben Newman, Allen Nie, Juan Carlos Niebles, Hamed Nilforoshan, Julian Nyarko, Giray Ogut, Laurel Orr, Isabel Papadimitriou, Joon Sung Park, Chris Piech, Eva Portelance, Christopher Potts, Aditi Raghunathan, Rob Reich, Hongyu Ren, Frieda Rong, Yusuf Roohani, Camilo Ruiz, Jack Ryan, Christopher Ré, Dorsa Sadigh, Shiori Sagawa, Keshav Santhanam, Andy Shih, Krishnan Srinivasan, Alex Tamkin, Rohan Taori, Armin W. Thomas, Florian Tramèr, Rose E.

- Wang, William Wang, Bohan Wu, Jiajun Wu, Yuhuai Wu, Sang Michael Xie, Michihiro Yasunaga, Jiaxuan You, Matei Zaharia, Michael Zhang, Tianyi Zhang, Xikun Zhang, Yuhui Zhang, Lucia Zheng, Kaitlyn Zhou, and Percy Liang. 2021. [On the opportunities and risks of foundation models](#).
- Samuel R. Bowman and George Dahl. 2021. [What will it take to fix benchmarking in natural language understanding?](#) In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 4843–4855, Online.
- L Karl Branting, Craig Pfeifer, Bradford Brown, Lisa Ferro, John Aberdeen, Brandy Weiss, Mark Pfaff, and Bill Liao. 2021. [Scalable and explainable legal prediction](#). *Artificial Intelligence and Law*, 29(2):213–238.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Cristian Cardellino, Milagro Teruel, Laura Alonso Alemany, and Serena Villata. 2017. [Legal NERC with ontologies, Wikipedia and curriculum learning](#). In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 2, Short Papers*, pages 254–259, Valencia, Spain. Association for Computational Linguistics.
- Ilias Chalkidis and Ion Androutsopoulos. 2017. [A deep learning approach to contract element extraction](#). In *Proceedings of the 30th International Conference on Legal Knowledge and Information Systems (JURIX 2017)*, Luxembourg City, Luxembourg.
- Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. 2019a. [Neural legal judgment prediction in English](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 4317–4323, Florence, Italy. Association for Computational Linguistics.
- Ilias Chalkidis, Ion Androutsopoulos, and Achilleas Michos. 2018. [Obligation and prohibition extraction using hierarchical RNNs](#). In *Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)*, pages 254–259, Melbourne, Australia. Association for Computational Linguistics.
- Ilias Chalkidis, Emmanouil Fergadiotis, Prodromos Malakasiotis, and Ion Androutsopoulos. 2019b. [Large-scale multi-label text classification on EU legislation](#). In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pages 6314–6322, Florence, Italy.
- Ilias Chalkidis, Manos Fergadiotis, and Ion Androutsopoulos. 2021a. [MultiEURLEX - a multi-lingual and multi-label legal document classification dataset for zero-shot cross-lingual transfer](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, Online.
- Ilias Chalkidis, Manos Fergadiotis, Sotiris Kotitsas, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020a. [An empirical study on large-scale multi-label text classification including few and zero-shot labels](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 7503–7515, Online.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020b. [LEGAL-BERT: The muppets straight out of law school](#). In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2898–2904, Online.
- Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, and Ion Androutsopoulos. 2019c. [Neural contract element extraction revisited](#). In *Proceedings of the Document Intelligence Workshop at NeurIPS 2019*, Vancouver, Canada.
- Ilias Chalkidis, Manos Fergadiotis, Nikolaos Manginas, Eva Katakalous, and Prodromos Malakasiotis. 2021b. [Regulatory compliance through Doc2Doc information retrieval: A case study in EU/UK legislation where text similarity has limitations](#). In *Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume*, pages 3498–3511, Online. Association for Computational Linguistics.
- Ilias Chalkidis, Manos Fergadiotis, Dimitrios Tsarapatanis, Nikolaos Aletras, Ion Androutsopoulos, and Prodromos Malakasiotis. 2021c. [Paragraph-level rationale extraction through regularization: A case study on european court of human rights cases](#). In *Proceedings of the Annual Conference of the North American Chapter of the Association for Computational Linguistics*, online.
- Ilias Chalkidis and Dimitrios Kampas. 2018. [Deep learning in law: early adaptation and legal word embeddings trained on large corpora](#). *Artificial Intelligence and Law*, 27(2):171–198.
- Ilias Chalkidis, Tommaso Passini, Sheng Zhang, Letizia Tomada, Sebastian Felix Schwemer, and Anders Søgaard. 2022. [Fairlex: A multilingual benchmark for evaluating fairness in legal text processing](#). In *Proceedings of the 60th Annual Meeting of the*

- Association for Computational Linguistics, Dublin, Ireland.
- Yanguang Chen, Yuanyuan Sun, Zhihao Yang, and Hongfei Lin. 2020. [Joint entity and relation extraction for legal documents with legal feature enhancement](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 1561–1571, online.
- Alexis Conneau and Douwe Kiela. 2018. [SentEval: An evaluation toolkit for universal sentence representations](#). In *Proceedings of the Eleventh International Conference on Language Resources and Evaluation (LREC-2018)*, Miyazaki, Japan. European Languages Resources Association (ELRA).
- Corinna Cortes and Vladimir Vapnik. 1995. Support-vector networks. *Machine Learning*, pages 273–297.
- Corinna Coupette, Janis Beckedorf, Dirk Hartung, Michael Bommarito, and Daniel Martin Katz. 2021. [Measuring law over time: A network analytical framework with an application to statutes and regulations in the United States and Germany](#). *Frontiers in Physics*, 9:269.
- Sylvie Delacroix. 2022. [Diachronic interpretability and machine learning systems](#). *Journal of Cross-disciplinary Research in Computational Law*, 1(1).
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. [BERT: Pre-training of deep bidirectional transformers for language understanding](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 4171–4186, Minneapolis, Minnesota. Association for Computational Linguistics.
- Giuseppe Di Fatta, Victor Sheng, and Alfredo Cuzzocrea. 2020. The IEEE ICDM 2020 workshops. In *2020 International Conference on Data Mining Workshops (ICDMW)*, pages 26–29. IEEE.
- John S. Downie. 2004. [IMIRSEL: a secure music retrieval testing environment](#). In *Internet Multimedia Management Systems V*, volume 5601, pages 91 – 99. International Society for Optics and Photonics, SPIE.
- Julia Dressel and Hany Farid. 2018. [The accuracy, fairness, and limits of predicting recidivism](#). *Science Advances*, 4(10).
- Oluwaseyi Feyisetan, Borja Balle, Thomas Drake, and Tom Diethe. 2020. [Privacy and utility-preserving textual analysis via calibrated multivariate perturbations](#). In *Proceedings of the 13th International Conference on Web Search and Data Mining*, pages 178–186.
- Rupert Haigh. 2018. *Legal English*. Routledge.
- Pengcheng He, Xiaodong Liu, Jianfeng Gao, and Weizhu Chen. 2021. [Deberta: Decoding-enhanced bert with disentangled attention](#). In *International Conference on Learning Representations*.
- Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. 2021. [CUAD: An expert-annotated NLP dataset for legal contract review](#). In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 1)*.
- Abhik Jana and Chris Biemann. 2021. [An investigation towards differentially private sequence tagging in a federated framework](#). In *Proceedings of the Third Workshop on Privacy in Natural Language Processing*, pages 30–35.
- Yoshinobu Kano, Mi-Young Kim, Randy Goebel, and Ken Satoh. 2017. Overview of coliee 2017. In *COL-IEE@ ICAIL*, pages 1–8.
- Yoshinobu Kano, Mi-Young Kim, Masaharu Yoshioka, Yao Lu, Juliano Rabelo, Naoki Kiyota, Randy Goebel, and Ken Satoh. 2018. Coliee-2018: Evaluation of the competition on legal information extraction and entailment. In *JSAI International Symposium on Artificial Intelligence*, pages 177–192. Springer.
- Jared Kaplan, Sam McCandlish, Tom Henighan, Tom B. Brown, Benjamin Chess, Rewon Child, Scott Gray, Alec Radford, Jeffrey Wu, and Dario Amodei. 2020. [Scaling laws for neural language models](#). *CoRR*, abs/2001.08361.
- Daniel Martin Katz, Michael J Bommarito, and Josh Blackman. 2017. [A general approach for predicting the behavior of the supreme court of the united states](#). *PloS one*, 12(4):e0174698.
- Daniel Martin Katz, Corinna Coupette, Janis Beckedorf, and Dirk Hartung. 2020. Complex societies and the growth of the law. *Scientific Reports*, 10:18737.
- Aaron Russell Kaufman, Peter Kraft, and Maya Sen. 2019. [Improving supreme court forecasting using boosted decision trees](#). *Political Analysis*, 27(3):381–387.
- Phi Manh Kien, Ha-Thanh Nguyen, Ngo Xuan Bach, Vu Tran, Minh Le Nguyen, and Tu Minh Phuong. 2020. [Answering legal questions by learning neural attentive text representation](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 988–998, online.
- Mi-Young Kim, Randy Goebel, Yoshinobu Kano, and Ken Satoh. 2016. Coliee-2016: evaluation of the competition on legal information extraction and entailment. In *International Workshop on Jurisinformatics (JURISIN 2016)*.
- Mi-young Kim, Ying Xu, and Randy Goebel. 2015. [A Convolutional Neural Network in Legal Question Answering](#). *Ninth International Workshop on Jurisinformatics (JURISIN)*.

- D. P. Kingma and J. Ba. 2015. [Adam: A method for stochastic optimization](#). In *Proceedings of the 5th International Conference on Learning Representations (ICLR)*, San Diego, CA, USA.
- Kankawin Kowsrihawat, Peerapon Vateekul, and Prachya Boonkwan. 2018. [Predicting judicial decisions of criminal cases from thai supreme court using bi-directional gru with attention mechanism](#). In *2018 5th Asian Conference on Defense Technology (ACDT)*, pages 50–55. IEEE.
- Elena Leitner, Georg Rehm, and Julian Moreno-Schneider. 2019. [Fine-grained named entity recognition in legal documents](#). In *Semantic Systems. The Power of AI and Knowledge Graphs*, pages 272–287, Cham. Springer International Publishing.
- Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2020. [BART: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 7871–7880, Online. Association for Computational Linguistics.
- Quentin Lhoest, Albert Villanova del Moral, Yacine Jernite, Abhishek Thakur, Patrick von Platen, Suraj Patil, Julien Chaumond, Mariama Drame, Julien Plu, Lewis Tunstall, Joe Davison, Mario Šaško, Gunjan Chhablani, Bhavitvya Malik, Simon Brandeis, Teven Le Scao, Victor Sanh, Canwen Xu, Nicolas Patry, Angelina McMillan-Major, Philipp Schmid, Sylvain Gugger, Clément Delangue, Théo Matuysière, Lysandre Debut, Stas Bekman, Pierrick Cistac, Thibault Goehringer, Victor Mustar, François Lagunas, Alexander M. Rush, and Thomas Wolf. 2021. [Datasets: A community library for natural language processing](#).
- Marco Lippi, Przemysław Pałka, Giuseppe Contissa, Francesca Lagioia, Hans-Wolfgang Micklitz, Giovanni Sartor, and Paolo Torroni. 2019. [CLAUDETTE: an automated detector of potentially unfair clauses in online terms of service](#). *Artificial Intelligence and Law*, pages 117–139.
- Yang Liu, Jiaxiang Liu, Yuxiang Lu, shikun feng, Yu Sun, Zhida Feng, Li Chen, Hao Tian, hua wu, and Haifeng Wang. 2022. [ERNIE-SPARSE: Robust efficient transformer through hierarchically unifying isolated information](#).
- Yinhan Liu, Myle Ott, Naman Goyal, Jingfei Du, Mandar Joshi, Danqi Chen, Omer Levy, Mike Lewis, Luke Zettlemoyer, and Veselin Stoyanov. 2019. [Roberta: A robustly optimized bert pretraining approach](#). *arXiv preprint arXiv:1907.11692*.
- Daniel Locke and Guido Zuccon. 2018. [A test collection for evaluating legal case law search](#). In *The 41st International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR '18*, page 1261–1264, New York, NY, USA. Association for Computing Machinery.
- Antoine Louis and Gerasimos Spanakis. 2022. [A statutory article retrieval dataset in french](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics*, page To appear. Association for Computational Linguistics.
- Bingfeng Luo, Yansong Feng, Zheng Wang, Zhanxing Zhu, Songfang Huang, Rui Yan, and Dongyan Zhao. 2017. [Learning with noise: Enhance distantly supervised relation extraction with dynamic transition matrix](#). In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 430–439, Vancouver, Canada. Association for Computational Linguistics.
- Pedro Henrique Luz de Araujo, Teófilo Emídio de Campos, Fabricio Ataides Braz, and Nilton Correia da Silva. 2020. [VICTOR: a dataset for Brazilian legal documents classification](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 1449–1458, Marseille, France. European Language Resources Association.
- Vijit Malik, Rishabh Sanjay, Shubham Kumar Nigam, Kripa Ghosh, Shouvik Kumar Guha, Arnab Bhat-tacharya, and Ashutosh Modi. 2021. [ILDC for CJPE: indian legal documents corpus for court judgment prediction and explanation](#). In *Joint Conference of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP 2021)*, online.
- Bryan McCann, Nitish Shirish Keskar, Caiming Xiong, and Richard Socher. 2018. [The natural language decathlon: Multitask learning as question answering](#). *CoRR*, abs/1806.08730.
- Masha Medvedeva, Ahmet Üstun, Xiao Xu, Michel Vols, and Martijn Wieling. 2021. [Automatic judgement forecasting for pending applications of the European Court of Human Rights](#). In *Proceedings of the Fifth Workshop on Automated Semantic Analysis of Information in Legal Text (ASAIL 2021)*.
- Masha Medvedeva, Michel Vols, and Martijn Wieling. 2020. [Using machine learning to predict decisions of the European Court of Human Rights](#). *Artificial Intelligence and Law*, 28(2):237–266.
- Eneldo Loza Mencia and Johannes Fürnkranzand. 2007. [An Evaluation of Efficient Multilabel Classification Algorithms for Large-Scale Problems in the Legal Domain](#). In *Proceedings of the 1st Linguistic Annotation Workshop*, pages 126–132, Halle, Germany.
- Emre Mumcuoğlu, Ceyhun E Öztürk, Haldun M Ozaktas, and Aykut Koç. 2021. [Natural language processing in law: Prediction of outcomes in the higher courts of turkey](#). *Information Processing & Management*, 58(5):102684.

- Ramesh Nallapati and Christopher D. Manning. 2008. [Legal docket classification: Where machine learning stumbles](#). In *Proceedings of the 2008 Conference on Empirical Methods in Natural Language Processing*, pages 438–446, Honolulu, Hawaii. Association for Computational Linguistics.
- Joel Niklaus, Ilias Chalkidis, and Matthias Stürmer. 2021. [Swiss-Court-Predict: A Multilingual Legal Judgment Prediction Benchmark](#). In *Proceedings of the 3rd Natural Legal Language Processing Workshop Workshop*, Online.
- Adam R Pah, David L Schwartz, Sarath Sanga, Zachary D Clopton, Peter DiCola, Rachel Davis Mersey, Charlotte S Alexander, Kristian J Hammond, and Luís A Nunes Amaral. 2020. [How to build a more open justice system](#). *Science*, 369(6500):134–136.
- F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel, B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer, R. Weiss, V. Dubourg, J. Vanderplas, A. Passos, D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in Python. *Journal of Machine Learning Research*, 12:2825–2830.
- Yifan Peng, Shankai Yan, and Zhiyong Lu. 2019. [Transfer Learning in Biomedical Natural Language Processing: An Evaluation of BERT and ELMo on Ten Benchmarking Datasets](#). In *Proceedings of the 2019 Workshop on Biomedical Natural Language Processing (BioNLP 2019)*.
- Colin Raffel, Noam Shazeer, Adam Roberts, Katherine Lee, Sharan Narang, Michael Matena, Yanqi Zhou, Wei Li, and Peter J. Liu. 2020. [Exploring the limits of transfer learning with a unified text-to-text transformer](#). *Journal of Machine Learning Research*, 21(140):1–67.
- Inioluwa Deborah Raji, Emily Denton, Emily M. Bender, Alex Hanna, and Amandalynne Paullada. 2021. [AI and the everything in the whole wide world benchmark](#). In *Thirty-fifth Conference on Neural Information Processing Systems Datasets and Benchmarks Track (Round 2)*.
- Pranav Rajpurkar, Jian Zhang, Konstantin Lopyrev, and Percy Liang. 2016. [SQuAD: 100,000+ questions for machine comprehension of text](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 2383–2392, Austin, Texas. Association for Computational Linguistics.
- Abhilasha Ravichander, Alan W Black, Shomir Wilson, Thomas Norton, and Norman Sadeh. 2019. [Question answering for privacy policies: Combining computational and legal perspectives](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 4947–4958, Hong Kong, China.
- Theodore W Ruger, Pauline T Kim, Andrew D Martin, and Kevin M Quinn. 2004. [The supreme court forecasting project: Legal and political science approaches to predicting supreme court decisionmaking](#). *Columbia Law Review*, pages 1150–1210.
- Sebastian Felix Schwemer, Letizia Tomada, and Tommaso Pasini. 2021. [Legal ai systems in the eu’s proposed artificial intelligence act](#). In *In Joint Proceedings of the Workshops on Automated Semantic Analysis of Information in Legal Text (ASAIL 2021) and AI and Intelligent Assistance for Legal Professionals in the Digital Workplace (LegalAIIA 2021)*.
- Tatiana Shavrina and Valentin Malykh. 2021. [How not to lie with a benchmark: Rearranging NLP leaderboards](#).
- Yanchuan Sim, Bryan Routledge, and Noah A. Smith. 2016. [Friends with motives: Using text to infer influence on SCOTUS](#). In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pages 1724–1733, Austin, Texas. Association for Computational Linguistics.
- Harold J. Spaeth, Lee Epstein, Jeffrey A. Segal Andrew D. Martin, Theodore J. Ruger, and Sara C. Benesh. 2020. [Supreme Court Database, Version 2020 Release 01](#). Washington University Law.
- Benjamin Strickson and Beatriz De La Iglesia. 2020. [Legal judgement prediction for uk courts](#). In *Proceedings of the 2020 The 3rd International Conference on Information Science and System*, pages 204–209.
- Peter M Tiersma. 1999. *Legal language*. University of Chicago Press.
- Dimitrios Tsarapatsanis and Nikolaos Aletras. 2021. [On the ethical limits of natural language processing on legal text](#). In *Findings of the Association for Computational Linguistics: ACL-IJCNLP 2021*, pages 3590–3599, Online. Association for Computational Linguistics.
- Don Tuggener, Pius von Däniken, Thomas Peetz, and Mark Cieliebak. 2020. [LEDGAR: A large-scale multi-label corpus for text classification of legal provisions in contracts](#). In *Proceedings of the 12th Language Resources and Evaluation Conference*, pages 1235–1241, Marseille, France. European Language Resources Association.
- Stefanie Urchs, Jelena Mitrović, and Michael Granitzer. 2021. [Design and Implementation of German Legal Decision Corpora](#). In *Proceedings of the 13th International Conference on Agents and Artificial Intelligence*, pages 515–521, Online. SCITEPRESS - Science and Technology Publications.
- Josef Valvoda, Tiago Pimentel, Niklas Stoehr, Ryan Cotterell, and Simone Teufel. 2021. [What about the precedent: An information-theoretic analysis of](#)

- common law. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 2275–2288, Online. Association for Computational Linguistics.
- Ashish Vaswani, Noam Shazeer, Niki Parmar, Jakob Uszkoreit, Llion Jones, Aidan N. Gomez, Lukasz Kaiser, and Illia Polosukhin. 2017. [Attention is all you need](#). In *Proceedings of the 31st International Conference on Neural Information Processing Systems*, pages 6000–6010, Long Beach, California, USA.
- Michael Benedict L. Virtucio, Jeffrey A. Aborot, John Kevin C. Abonita, Roxanne S. Avinante, Rother Jay B. Copino, Michelle P. Neverida, Vanesa O. Osiana, Elmer C. Peramo, Joanna G. Syjuco, and Glenn Brian A. Tan. 2018. [Predicting decisions of the philippine supreme court using natural language processing and machine learning](#). In *2018 IEEE 42nd Annual Computer Software and Applications Conference (COMPSAC)*, volume 2, pages 130–135. IEEE.
- Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. [Superglue: A stickier benchmark for general-purpose language understanding systems](#). In *Advances in Neural Information Processing Systems*, volume 32. Curran Associates, Inc.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2018. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *Proceedings of the 2018 EMNLP Workshop BlackboxNLP: Analyzing and Interpreting Neural Networks for NLP*, pages 353–355, Brussels, Belgium. Association for Computational Linguistics.
- Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R. Bowman. 2019b. [GLUE: A multi-task benchmark and analysis platform for natural language understanding](#). In *International Conference on Learning Representations*, New Orleans, Louisiana, USA.
- Yuzhong Wang, Chaojun Xiao, Shirong Ma, Haoxi Zhong, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2021. [Equality before the law: Legal judgment consistency analysis for fairness](#). *Science China - Information Sciences*.
- Christopher Williams. 2007. *Tradition and change in legal English: Verbal constructions in prescriptive texts*, volume 20. Peter Lang.
- Thomas Wolf, Lysandre Debut, Victor Sanh, Julien Chaumond, Clement Delangue, Anthony Moi, Pierric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz, Joe Davison, Sam Shleifer, Patrick von Platen, Clara Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven Le Scao, Sylvain Gugger, Mariama Drame, Quentin Lhoest, and Alexander Rush. 2020. [Transformers: State-of-the-art natural language processing](#). In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing: System Demonstrations*, pages 38–45, Online. Association for Computational Linguistics.
- Chaojun Xiao, Xueyu Hu, Zhiyuan Liu, Cunchao Tu, and Maosong Sun. 2021. [Lawformer: A pre-trained language model for chinese legal long documents](#). *CoRR*, abs/2105.03887.
- Liu Yang, Mingyang Zhang, Cheng Li, Michael Bendersky, and Marc Najork. 2020. [Beyond 512 Tokens: Siamese Multi-Depth Transformer-Based Hierarchical Encoder for Long-Form Document Matching](#), page 1725–1734. Association for Computing Machinery, New York, NY, USA.
- Wenmian Yang, Weijia Jia, Xiaojie Zhou, and Yutao Luo. 2019. [Legal judgment prediction via multi-perspective bi-feedback network](#). In *Proceedings of the Twenty-Eighth International Joint Conference on Artificial Intelligence, IJCAI-19*, pages 4085–4091. International Joint Conferences on Artificial Intelligence Organization.
- Hai Ye, Xin Jiang, Zhunchen Luo, and Wenhan Chao. 2018. [Interpretable Charge Predictions for Criminal Cases: Learning to Generate Court Views from Fact Descriptions](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1854–1864, New Orleans, Louisiana. Association for Computational Linguistics.
- Manzil Zaheer, Guru Guruganesh, Kumar Avinava Dubey, Joshua Ainslie, Chris Alberti, Santiago Ontanon, Philip Pham, Anirudh Ravula, Qifan Wang, Li Yang, and Amr Ahmed. 2020. [Big Bird: Transformers for longer sequences](#). In *Advances in Neural Information Processing Systems*, pages 17283–17297, online.
- Lucia Zheng, Neel Guha, Brandon R. Anderson, Peter Henderson, and Daniel E. Ho. 2021. [When does pretraining help? assessing self-supervised learning for law and the casehold dataset](#). In *Proceedings of the 18th International Conference on Artificial Intelligence and Law*. Association for Computing Machinery.
- Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Chaojun Xiao, Zhiyuan Liu, and Maosong Sun. 2018. [Legal judgment prediction via topological learning](#). In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 3540–3549, Brussels, Belgium.
- Haoxi Zhong, Yuzhong Wang, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020a. [Iteratively questioning and answering for interpretable legal judgment prediction](#). In *Proceedings of the AAAI Conference on Artificial Intelligence*, pages 1250–1257.

Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020b. [How does nlp benefit legal system: A summary of legal artificial intelligence](#). In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics*, pages 5218–5230.

Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang Zhang, Zhiyuan Liu, and Maosong Sun. 2020c. [JEC-QA: A legal-domain question answering dataset](#). In *The Thirty-Fourth AAAI Conference on Artificial Intelligence*, pages 9701–9708, New York, NY, USA.

Octavia-Maria Şulea, Marcos Zampieri, Mihaela Vela, and Josef van Genabith. 2017. [Predicting the Law Area and Decisions of French Supreme Court Cases](#). In *Proceedings of the International Conference Recent Advances in Natural Language Processing, RANLP 2017*, pages 716–722, Varna, Bulgaria. INCOMA Ltd.

A Datasets, Code, and Participation

Where are the datasets? We provide access to LexGLUE on Hugging Face Datasets (Lhoest et al., 2021) at https://huggingface.co/datasets/lex_glue. For example, to load the SCOTUS dataset, you first simply install the datasets Python library and then make the following call:

```
from datasets import load_dataset
dataset = load_dataset('lex_glue', task='scotus')
```

How do I run experiments? To make reproducing the results of the already examined models or future models even easier, we release our code on GitHub (<https://github.com/coastalcph/lex-glue>). In that repository (in the folder /EXPERIMENTS), there are Python scripts, relying on the Hugging Face Transformers library (Wolf et al., 2020), to run and evaluate any Transformer-based model (e.g., BERT, RoBERTa, LegalBERT, and their hierarchical variants, as well as, Longformer, and BigBird). We also provide bash scripts to replicate the experiments for each dataset with 5 random seeds, as we did for the reported results.

B No labeling as an additional class

In ECtHR Tasks A & B and UNFAIR-ToS, there are unlabeled samples. Concretely, in ECtHR Task A, a possible event is *no violation*, i.e., the court ruled that the defendant did not violate any ECHR article. Contrary, *no violation* is not a possible event in the original ECtHR Task B dataset, i.e., at least a single ECHR article is allegedly violated

(considered by the court) in every case; however, there is such a rare scenario after the simplifications we introduced, i.e., some cases were originally labeled only with rare labels that were excluded from our benchmark (Section 3.2). In UNFAIR-ToS, the vast majority of sentences are not labeled with any type of *unfairness* (unfair term against users), i.e., most sentences do not raise any questions of possible violations of the European consumer law.

In multi-label classification, the set of labels per instance is represented as a one-hot vector $Y = [y_1, y_2, \dots, y_L]$, where $y_i = 1$ if the instance is labeled with the i -th class, and $y_i = 0$ otherwise. If an instance is not labeled with any class, its Y includes only zeros. During training, binary cross-entropy correctly penalizes such instances, if the predictions ($\hat{Y} = [\hat{y}_1, \hat{y}_2, \dots, \hat{y}_L]$) diverge from zeros. During evaluation, however, the F1-score ($F1 = \frac{TP}{TP + \frac{1}{2}(FP + FN)}$) ignores instances with $Y = \hat{Y} = [0, 0, \dots, 0]$, because it considers only the true positives (TP), false positives (FP), and false negatives (FN), and instances where $Y = \hat{Y} = [0, 0, \dots, 0]$ contribute no TPs, FPs, FNs. In order to make F1 sensitive to the correct labeling of such examples, during evaluation (not training) we include an additional label (y_0 or \hat{y}_0) in both targets (Y) and predictions (\hat{Y}), whose value is 1 (positive) if the original (without y_0, \hat{y}_0) Y and \hat{Y} are $Y = [0, 0, \dots, 0]$ or $\hat{Y} = [0, 0, \dots, 0]$, respectively, and 0 (negative) otherwise. This is particularly important for proper evaluation, as across three datasets a considerable portion of the examples are unlabeled (11.5% in ECtHR Task A, 1.6% in ECtHR Task B, and 95.5% in UNFAIR-ToS).

C Additional Results

Tables 5 and 6 show *development* results for all examined models across datasets. We report the mean and standard deviations (\pm) for the three seeds (among the five used) with the best development scores per model to exclude catastrophic failures, i.e., runs with severely low performance. The standard deviations are relatively low across models and datasets (up to 0.5% for μ -F₁ and up to 1% for m-F₁). The development results are generally higher compared to the test ones (cf. Table 3) in many cases, as one would expect.

Table 7 reports training times per dataset and model; both the time per epoch (T/e), and the total training time (T) across all epochs. All full-attention BERT models, except Longformer and

Method	ECtHR (A)*	ECtHR (B)*	SCOTUS*	EUR-LEX	LEDGAR	UNFAIR-ToS	CaseHOLD
TFIDF-SVM	65.0	75.3	78.6	73.7	86.8	94.1	22.4
BERT	71.0 \pm 0.7	79.6 \pm 0.5	72.7 \pm 0.2	77.3 \pm 0.2	87.9 \pm 0.1	95.5 \pm 0.0	72.8 \pm 0.1
RoBERTa	70.4 \pm 0.5	78.4 \pm 0.7	76.9 \pm 0.6	77.6 \pm 0.0	88.1 \pm 0.1	94.8 \pm 0.2	74.1 \pm 0.2
DeBERTa	69.3 \pm 0.7	79.0 \pm 0.3	76.1 \pm 0.5	77.8 \pm 0.1	88.3 \pm 0.2	95.5 \pm 0.1	73.8 \pm 0.1
Longformer	71.0 \pm 0.3	80.4 \pm 0.9	76.9 \pm 0.0	77.5 \pm 0.0	88.1 \pm 0.2	95.1 \pm 0.2	73.9 \pm 0.2
BigBird	71.0 \pm 0.2	80.1 \pm 0.5	75.9 \pm 0.2	77.3 \pm 0.1	88.0 \pm 0.1	95.2 \pm 0.4	73.7 \pm 0.2
Legal-BERT	71.9 \pm 0.4	79.8 \pm 0.2	80.4 \pm 0.3	77.6 \pm 0.1	88.5 \pm 0.0	95.1 \pm 0.2	76.4 \pm 0.3
CaseLaw-BERT	72.1 \pm 0.3	79.6 \pm 0.0	81.3 \pm 0.6	77.2 \pm 0.1	88.4 \pm 0.2	95.3 \pm 0.4	77.4 \pm 0.2

Table 5: Development μ -F₁ results for all examined models across all LexGLUE tasks. We report the mean and standard deviations (\pm) for the three seeds with the best development scores per model. In starred datasets, we use the hierarchical variant of each model, except for Longformer and BigBird, as discussed in Section 4.2.

Method	ECtHR (A)*	ECtHR (B)*	SCOTUS*	EUR-LEX	LEDGAR	UNFAIR-ToS	CaseHOLD
TFIDF-SVM	55.6	64.1	71.2	56.9	79.6	69.4	22.0
BERT	65.4 \pm 1.2	74.8 \pm 0.6	65.9 \pm 0.8	62.6 \pm 0.8	81.8 \pm 0.1	75.8 \pm 1.3	72.8 \pm 0.1
RoBERTa	65.4 \pm 0.2	74.2 \pm 1.1	69.5 \pm 0.8	63.5 \pm 0.4	81.9 \pm 0.2	74.4 \pm 0.7	74.1 \pm 0.2
DeBERTa	63.5 \pm 0.9	74.0 \pm 0.4	68.4 \pm 0.8	63.6 \pm 0.3	82.0 \pm 0.5	77.1 \pm 1.2	73.8 \pm 0.1
Longformer	65.5 \pm 1.6	77.7 \pm 1.0	70.4 \pm 0.5	63.8 \pm 0.5	82.0 \pm 0.3	75.2 \pm 1.2	73.9 \pm 0.2
BigBird	65.8 \pm 1.1	74.1 \pm 0.5	69.1 \pm 0.2	63.0 \pm 0.3	81.7 \pm 0.2	76.5 \pm 1.8	73.7 \pm 0.2
Legal-BERT	68.0 \pm 0.2	76.1 \pm 0.5	72.7 \pm 0.2	62.0 \pm 0.9	82.2 \pm 0.3	76.9 \pm 1.3	76.4 \pm 0.3
CaseLaw-BERT	67.1 \pm 0.7	74.6 \pm 0.5	74.0 \pm 1.2	62.9 \pm 0.3	82.3 \pm 0.3	76.5 \pm 0.3	77.4 \pm 0.2

Table 6: Development m-F₁ results for all examined models across all LexGLUE tasks. We report the mean and standard deviation (\pm) for the three seeds with the best development scores per model. In starred datasets, we use the hierarchical variant of each model, except for Longformer and BigBird, as discussed in Section 4.2.

Big-Bird, have comparable times with the exception of DeBERTa that has four separate attention mechanisms. We observe that when the hierarchical variant of these models is deployed, i.e., in ECtHR tasks and SCOTUS, it is approximately twice (2 \times) as fast compared to Longformer and BigBird.

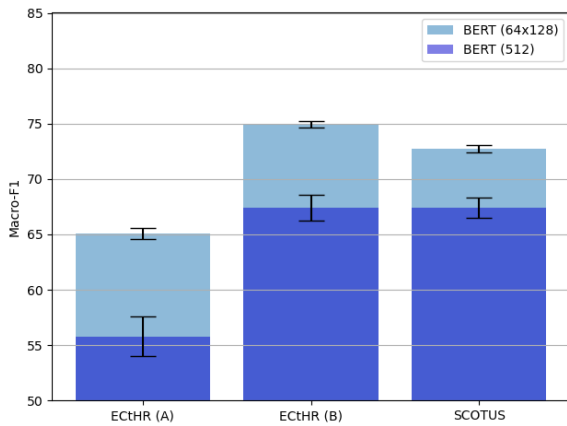


Figure 3: Development m-F₁ scores of standard BERT (up to 512 tokens) and its hierarchical variant (Section 4.2, 64 \times 128 tokens) in ECtHR (Task A, B) and SCOTUS, i.e., the datasets with long documents. Light blue denotes the average score across 5 runs for the hierarchical variant (used in Table 3 for these datasets), while dark blue corresponds to standard BERT (not used in Table 3 for these datasets). The error bars show the standard error.

D Use of 512-token BERT models

In Figure 3, we show results for the standard BERT model of Devlin et al. (2019), which can process up to 512 tokens, compared to its hierarchical variant (Section 4.2), which can process up to 64 \times 128 tokens. We observe that across all datasets that contain long documents (ECtHR A & B, SCOTUS, cf. Fig. 2(a)), the hierarchical variant clearly outperforms the standard model fed with truncated documents (ECtHR A: +10.2% p.p., ECtHR B: 7.5% p.p., SCOTUS: 4.9% p.p.). Compared to the ECtHR tasks, the gains are lower in SCOTUS, a topic classification task where long-range reasoning is not needed; by contrast, for ECtHR multiple distant facts need to be combined. Based on these results, we conclude that using severely truncated documents is not a plausible option for LexGLUE, and other directions for processing long documents should be considered in the future, ideally fully pre-trained hierarchical models, contrary to our semi-pre-trained hierarchical models (Section 6).

E Use of Roberta Large

We additionally evaluate RoBERTa-large, i.e., 24 Transformer blocks, 1024 hidden units, and 18 at-

Method	ECtHR (A)*		ECtHR (B)*		SCOTUS*		EUR-LEX		LEDGAR		CaseHOLD	
	<i>T</i>	<i>T/e</i>	<i>T</i>	<i>T/e</i>	<i>T</i>	<i>T/e</i>	<i>T</i>	<i>T/e</i>	<i>T</i>	<i>T/e</i>	<i>T</i>	<i>T/e</i>
BERT	3h 42m	28m	3h 9m	28m	1h 24m	11m	3h 36m	19m	6h 9m	21m	4h 24m	24m
RoBERTa	4h 11m	27m	3h 43m	27m	2h 46m	17m	3h 36m	19m	6h 22m	21m	4h 21m	24m
DeBERTa	7h 43m	46m	6h 48m	46m	3h 42m	29m	5h 34m	36m	9h 29m	40m	6h 42m	45m
Longformer	6h 47m	56m	7h 31m	56m	6h 27m	34m	11h 10m	45m	15h 47m	50m	4h 45m	30m
BigBird	8h 41m	1h 2m	8h 17m	1h 2m	5h 51m	37m	3h 57m	24m	8h 13m	27m	6h 4m	49m
Legal-BERT	3h 52m	28m	3h 2m	28m	2h 2m	17m	3h 22m	19m	5h 23m	21m	4h 13m	23m
CaseLaw-BERT	3h 2m	28m	2h 57m	28m	2h 34m	34m	3h 40m	19m	6h 8m	21m	4h 21m	24m

Table 7: Training time in total (*T*) and per epoch (*T/e*) across LexGLUE tasks. In starred datasets, we use the hierarchical variant of each model, except for Longformer and BigBird, as described in Section 4.2.

Method	ECtHR (A)*		ECtHR (B)*		SCOTUS*		EUR-LEX		LEDGAR		UNFAIR-ToS		CaseHOLD
	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁	m-F ₁	μ -F ₁ /m-F ₁
Results on Development Set													
RoBERTa (B)	70.6	65.7	79.3	75.8	77.5	64.1	77.6	70.4	88.0	82.1	94.6	75.2	74.3
RoBERTa (L)	72.7	69.3	81.1	77.0	74.6	56.9	78.0	74.5	88.5	82.8	95.8	80.3	76.8
Legal-BERT	72.5	68.2	79.7	76.8	77.6	63.3	80.8	72.9	88.5	82.6	95.3	78.2	76.6
CaseLaw-BERT	71.8	67.7	79.5	74.9	77.3	63.1	82.1	75.6	88.7	82.7	95.7	76.9	77.7
Results on Test Set													
RoBERTa (B)	69.2	59.0	77.3	68.9	71.6	62.0	71.9	57.9	87.9	82.3	95.2	79.2	71.4
RoBERTa (L)	73.8	67.6	79.8	71.6	75.5	66.3	72.5	58.1	88.6	83.6	95.8	81.6	74.4
Legal-BERT	70.0	64.0	80.4	74.7	76.4	66.5	72.1	57.4	88.2	83.0	96.0	83.0	75.3
CaseLaw-BERT	69.8	62.9	78.8	70.3	76.6	65.9	70.7	56.6	88.3	83.0	96.0	82.3	75.4

Table 8: Development and test results across LexGLUE tasks. In starred datasets, we use the hierarchical variant of each model, discussed in Section 4.2. (B) and (L) denote the base and large version of RoBERTa, respectively.

tention heads, to better understand the dynamics between domain specificity and model size. In this case, we use the AdamW optimizer with a 1e-5 maximum learning rate, warm-up ratio of 0.1, and a weight decay rate of 0.06, and we use a similar mini-batch size of 8 examples.²⁰

Table 8 reports the development and test results using the seed (run) with the best development scores. We observe that using the large version of RoBERTa, dubbed RoBERTa (L), with more than 2× parameters (355M), leads to substantial performance improvements compared to the base version of RoBERTa, dubbed RoBERTa (B), across all tasks. The results are comparable, or better in some cases, compared to the legal-oriented language models (Legal-BERT, CaseLaw-BERT).

Considering that the two legal-oriented models are much smaller and have been pre-trained with (5–10×) less data (Section 2), we have a strong indication for expected performance gains by pre-training larger legal-oriented models using larger legal corpora (Section 6).

²⁰Large models tend to be very sensitive to parameter updates, especially in the initial training steps; hence a smaller learning rate and warm up steps are very crucial.

F Other Tasks and Datasets Considered

We considered including the Contract Understanding Atticus Dataset (CUAD) (Hendrycks et al., 2021), an expertly curated dataset that comprises 510 contracts annotated with 41 valuable contractual insights (e.g., agreement date, parties, governing law). The task is formulated as a SQUAD-like question answering task, where given a *question* (the name of an insight) and a *paragraph* from the contract, the model has to identify the answer span in the paragraph.²¹ The original dataset follows the SQUAD v2.0 setting, including unanswerable questions. Following SQUAD v1.1 (Rajpurkar et al., 2016), we simplified the task by removing all unanswerable pairs (question, paragraph), which are the majority in the original dataset. We also excluded pairs whose answers exceeded 128 full words to alleviate the imbalance between short and long answers. We then re-split the dataset chronologically into training (5.2k, 1994–2019), development (572, 2019–2020), and test (604, 2020) sets.

Following Devlin et al. (2019), and similarly to Hendrycks et al. (2021), for each training (or test)

²¹The question mostly resembles a *prompt*, rather than a natural question, as there is a closed set of 41 alternatives.

instance, we consider pairs that consist of a question and a paragraph, separated by the special delimiter token [sep]. The top-level representations $[h_1, \dots, h_N]$ of the tokens of the paragraph are fed into a linear layer to obtain two logits per token (for the token being the start or end of the answer span), which are then passed through a softmax activation (separately for start and end) to obtain probability distributions. The tokens with the highest start and end probabilities are selected as boundaries of the answer span. We evaluated performance with token-level F1 score, similarly to SQUAD.

We trained all the models of Table 2, which scored approx. 10-20% in token-level F1, with Legal-BERT performing slightly better than the rest (+5% F1).²² In the paper that introduced CUAD (Hendrycks et al., 2021), several other measures (Precision@ N% Recall, AUPR, Jaccard similarity) are used to more leniently estimate a model’s ability to approximately locate answers in context paragraphs. Through careful manual inspection of the dataset, we noticed the following points that seem to require more careful consideration.

- Contractual insights (categories, shown in italics below) include both entity-level (short) answers (e.g., “SERVICE AGREEMENT” for *Document Name*, and “Imprimis Pharmaceuticals, Inc.” for *Parties*) and paragraph-level (long) answers (e.g., “If any of the conditions specified in Section 8 shall not have been fulfilled when and as required by this Agreement, or by the Closing Date, or waived in writing by Capital Resources, this Agreement and all of Capital Resources obligations hereunder may be canceled [...] except as otherwise provided in Sections 2, 7, 9 and 10 hereof.” for *Termination for Convenience*). These two different types of answers (short and paragraph-long) seem to require different models and different evaluation measures, unlike how they are treated in the original CUAD paper.
- Some contractual insights (categories), e.g., *Parties*, have been annotated with both short (e.g., “Imprimis Pharmaceuticals, Inc.”) and long (e.g., “together, Blackwell and Munksgaard shall be referred to as ‘the Publishers.’”) answers. Annotations of this kind introduce noise during both training and evaluation. For example, it becomes unclear when a short (finer/strict) or a long (loose) annotation should be taken to be the correct one.
- Annotations may include indirect mentions, e.g., ‘Franchisee’, ‘Service Provider’ for *Parties*, instead of the actual entities (the company name).
- Annotations may include semi-redacted text (e.g., “_____, 1996” for *Agreement Date*), or even fully redacted text (e.g., “_____” for *Parties*). This practice may be necessary to hide sensitive information, but for the purposes of a benchmark dataset such cases could have been excluded.

The points above, which seem to require revisiting the annotations of CUAD, and the very low F1 scores of all models led us to exclude CUAD from LexGLUE. We also note that there is related work covering similar topics, such as Contract Element Extraction (Chalkidis and Androutsopoulos, 2017), Contractual Obligation Extraction (Chalkidis et al., 2018), and Contractual Provision Classification (Tuggenier et al., 2020), where models perform much better (in terms of accuracy), relying on simpler (separate) more carefully designed tasks and much bigger datasets. Thus we believe that the points mentioned above, which blur the task definition of CUAD and introduce noise, and the limited (compared to larger datasets) number of annotations strongly affect the performance of the models on CUAD, underestimating their true potential.

We also initially considered some very interesting legal Information Retrieval (IR) datasets (Locke and Zuccon, 2018; Chalkidis et al., 2021b) that aim to examine crucial real-life tasks (relevant case law retrieval, regulatory compliance). However, we decided to exclude them from the first version of LexGLUE, because they rely on processing multiple long documents and require more task-specific neural network architectures (e.g., siamese networks), and different evaluation measures. Hence, they would make LexGLUE more complex and a less attractive entry point for newcomers to legal NLP. We plan, however, to include more demanding tasks in future LexGLUE versions, as the legal NLP community will be growing.

G Dataset Examples

In Table 9, we present training examples, i.e., pairs of input(s), output(s), for LexGLUE datasets and tasks. More examples can be inspected using the dataset preview functionality provided in the online dataset card of Hugging Face.²³

²²F1 is one of the two official SQUAD measures. In the second one, Exact Answer Accuracy, all models scored 0%.

²³https://huggingface.co/datasets/lex_glue

Dataset	Input(s)	Output(s) / Label(s)
ECtHR	<p>Text: 12. In 1987 the applicant association published a book entitled Euskadi at war. There were four versions – Basque, English, Spanish and French – and the book was distributed in numerous countries, including France and Spain. According to the applicant association, this was a collective work containing contributions from a number of academics with specialist knowledge of the Basque Country and giving an account of the historical, cultural, linguistic and socio-political aspects of the Basque cause. It ended with a political article entitled “Euskadi at war, a promise of peace” by the Basque national liberation movement.</p> <p>13. The book was published in the second quarter of 1987. On 29 April 1988 a ministerial order was issued by the French Ministry of the Interior under section 14 of the Law of 29 July 1881, as amended by the decree of 6 May 1939, banning the circulation, distribution and sale of the book in France in any of its four versions on the ground that “the circulation in France of this book, which promotes separatism and vindicates recourse to violence, is likely to constitute a threat to public order”. On 6 May 1988, pursuant to the aforementioned order, the département director of the airport and border police refused to allow over two thousand copies of the book to be brought into France. [...]</p>	<p>3 (<i>Right to a fair trial</i>)</p> <p>6 (<i>Freedom of expression</i>)</p>
SCOTUS	<p>Text: 329 U.S. 29 67 S.Ct. 1 91 L.Ed. 22</p> <p>CHAMPLIN REFINING CO v. UNITED STATES et al. No. 21. Argued Oct. 18, 21, 1946. Decided Nov. 18, 1946.</p> <p>Appeal from the District Court of the United States for the Western District of Oklahoma. Messrs. Dan Moody, of Austin, Tex., and Harry O. Glasser, of Enid, Okla., for appellant. Mr. Edward Dumbauld, of Washington, D.C., for appellees.</p> <p>Mr. Justice JACKSON delivered the opinion of the Court.</p> <p>1 The Interstate Commerce Commission, acting under § 19a of the Interstate Commerce Act,1 ordered the appellant to furnish certain inventories, schedules, maps and charts of its pipe line property.</p> <p>2 Champlin's objections that the Act does not authorize the order, or if it be construed to do so is unconstitutional, were overruled by the Commission and again by the District Court which dismissed the company's suit for an injunction.3 These questions of law are brought here by appeal. [...]</p>	<p>7 (<i>Economic Activity</i>)</p>
EUR-LEX	<p>Text: Commission Regulation (EC) No 1156/2001 of 13 June 2001 fixing the export refunds on white sugar and raw sugar exported in its unaltered state</p> <p>THE COMMISSION OF THE EUROPEAN COMMUNITIES</p> <p>Having regard to the Treaty establishing the European Community, Having regard to Council Regulation (EC) No 2038/1999 of 13 September 1999 on the common organisation of the markets in the sugar sector(1), as amended by Commission Regulation (EC) No 1527/2000(2), and in particular point (a) of the second subparagraph of Article 18(5) thereof,</p> <p>Whereas: (1) Article 18 of Regulation (EC) No 2038/1999 provides that the difference between quotations or prices on the world market for the products listed in Article 1(1)(a) of that Regulation and prices for those products within the Community may be covered by an export refund. (2) Regulation (EC) No 2038/1999 provides that when refunds on white and raw sugar, undenatured and exported in its unaltered state, are being fixed account must be taken of the situation on the Community and world markets in sugar and in particular of the price and cost factors [...]</p>	<p>28 (<i>Trade Policy</i>),</p> <p>93 (<i>Beverages and Sugar</i>),</p> <p>94 (<i>Foodstuff</i>)</p>
LEDGAR	<p>Text: The validity or unenforceability of any provision or provisions of this Agreement shall not affect the validity or enforceability of any other provision hereof, which will remain in full force and effect. Should a court or other body of competent jurisdiction determine that any provision of this Agreement is excessive in scope or otherwise illegal, invalid, void or unenforceable, such provision shall be adjusted rather than voided, if possible, so that it is enforceable to the maximum extent possible.</p>	<p>79 (<i>Severability</i>)</p>
UNFAIR-ToS	<p>Text: By creating a tinder account or by using the tinder imessage app (“tinder stacks”), whether through a mobile device , mobile application or computer (collectively, the “service”) you agree to be bound by (i) these terms of use, (ii) our privacy policy and safety tips, each of which is incorporated by reference into this agreement, and (ii) any terms disclosed and agreed to by you if you purchase additional features, products or services we offer on the service (collectively, this “agreement”).</p>	<p>4 (<i>Contract by Using</i>)</p>
CaseHOLD	<p>Context: Drapeau's cohorts, the cohort would be a “victim” of making the bomb. Further, firebombs are inherently dangerous. There is no peaceful purpose for making a bomb. Felony offenses that involve explosives qualify as “violent crimes” for purposes of enhancing the sentences of career offenders. See 18 U.S.C. § 924(e)(2)(B)(ii) (defining a “violent felony” as: “any crime punishable by imprisonment for a term exceeding one year ... that ... involves use of explosives”). Courts have found possession of a bomb to be a crime of violence based on the lack of a nonviolent purpose for a bomb and the fact that, by its very nature, there is a substantial risk that the bomb would be used against the person or property of another. See United States v. Newman, 125 F.3d 863 (10th Cir.1997) (unpublished) ([HOLDING]); United States v. Dodge, 846 F.Supp. 181</p> <p>Choices (Holdings):</p> <p>(A) “holding that possession of a pipe bomb is a crime of violence for purposes of 18 usc 3142f1”,</p> <p>(B) “holding that bank robbery by force and violence or intimidation under 18 usc 2113a is a crime of violence”,</p> <p>(C) “holding that sexual assault of a child qualified as crime of violence under 18 usc 16”,</p> <p>(D) “holding for the purposes of 18 usc 924e that being a felon in possession of a firearm is not a violent felony as defined in 18 usc 924e2b”,</p> <p>(E) “holding that a court must only look to the statutory definition not the underlying circumstances of the crime to determine whether a given offense is by its nature a crime of violence for purposes of 18 usc 16”</p>	<p>0 (<i>Choice A</i>)</p>

Table 9: Training examples (pairs of inputs, outputs) for LeXGLUE datasets and tasks.

H Updated version V4 (09/11/2022)

We updated the results for the TFIDF-SVM method in Tables 3, 5 and 6. There was a bug in our code base,²⁴ where TFIDF-SVM grid search function was re-fitting (re-training) the model given both the training and development sets. This was inconsistent with our general practices leading to overestimated development and test scores. The new corrected results make clear that TFIDF-SVM models are significantly outperformed in most cases, even on the SCOTUS dataset. Thanks to Yu-Chen and Daniel Gonzalez for pointing out this issue.

²⁴<https://github.com/coastalcph/lex-glue>