

# Faithfulness and Content Selection in Long-Input Multi-Document Summarisation of U.S. Civil Rights Litigation

Anonymous ACL submission

## Abstract

Automatic summarisation of legal cases would reduce the burden on legal professionals and increase the accessibility of the law. However, the abstractive methods which dominate recent research are prone to hallucination. Despite the fact that this is a barrier to practical use, preventing hallucination is currently an understudied area in the legal domain. We conduct the first study at the intersection of legal, multi-document, and faithful summarisation. In particular, by introducing a BERT-based content selection mechanism, we achieve an improvement of 0.2614 in the probability of a generated summary being entailed by its source text compared to a naïve content selection baseline, and observe qualitative improvements. Further, we demonstrate possible improvements of 5.56 ROUGE-1 F1, 5.46 ROUGE-2 F1, 2.7 ROUGE-L F1, and 2.15 BERTScore over the state-of-the-art if a perfectly predictive classifier was used, demonstrating the importance of content selection for summary faithfulness and quality for long-input legal abstractive summarisation.

## 1 Introduction

In common law jurisdictions, judicial decisions are informed by past cases, making identifying relevant precedent cases crucial (Zhong et al., 2020; Shukla et al., 2022; Askari et al., 2021; Althammer et al., 2021). However, the increasing number of precedent cases, each typically hundreds of pages long (Chalkidis et al., 2022), burdens legal professionals (Mumcuoğlu et al., 2021). While popular legal retrieval systems offer case summaries, these are costly and time-consuming to produce manually; automatic summarisation of legal cases using natural language processing tools would significantly benefit legal professionals, and increase the accessibility of the law.

However, abstractive summarisation methods are prone to hallucination - summaries may contain information which is unrelated or unfaithful to the

source text (Feijo and Moreira, 2023). This is a major barrier to practical applicability (Huang et al., 2023; Wang et al., 2020; Fischer et al., 2022; Narayan et al., 2022a; Curran et al., 2023; Norkute et al., 2021), especially in the high-stakes domain of law (Feijo and Moreira, 2019, 2023); a Lexis-Nexis (2024) report found that 57% of respondents were concerned about hallucination. Despite this, the hallucination problem is understudied in relation to legal data. Additionally, the fact that the length of legal texts frequently exceeds transformer-based models' input token limit (Chalkidis et al., 2022) presents a challenging scenario.

Our work confronts these challenges by addressing the following research questions: **RQ1**: Can we improve the quality and faithfulness of abstractive summarisation results by providing a better representation of the source data to the summarisation model - namely, by using a BERT-based content selector trained on OREO labels to identify salient information? **RQ2**: Do transformer based models pretrained in the legal domain further improve results for legal multi document abstractive summarisation?

We contribute to the growing literature on faithfulness in abstractive summarisation, legal summarisation, and multi document summarisation by being the first work at this intersection. Specifically, we: (i) demonstrate that our content selection strategy improves summary faithfulness, through qualitative analysis and an improvement of 0.2614 in the probability of a generated summary being entailed by its source text compared to a naïve content selection baseline; and, (ii) demonstrate possible gains of 5.56 ROUGE-1 F1, 5.46 ROUGE-2 F1, 2.70 ROUGE-L F1, and 2.15 BERTScore if a perfectly predictive classifier was used for content selection in our methodology.

## 2 Related Work

### 2.1 Summarisation

Automatic summarisation methods aim to condense input text into a fluent shorter text retaining the key information (Feijo and Moreira, 2023; Kornilova and Eidelman, 2019; Bajaj et al., 2021). Extractive methods involve selecting and assembling key information from the source text (Jain et al., 2023), while recent abstractive summarisation methods, which are increasingly based on transformer architectures, generate summaries from scratch, conditioned on the source text.

The majority of legal summarisation research focuses on extractive summarisation, which is not our focus. Abstractive summarisation has been shown to significantly outperform extractive methods (Feijo and Moreira, 2019; Bhattacharya et al., 2019; Klaus et al., 2022), especially as transformer-based models pretrained on legal corpora (Shukla et al., 2022; Mullick et al., 2022; Chalkidis et al., 2020; Niklaus and Giofre, 2023; Zheng et al., 2021) have now been publicly released. Legal abstractive summarisation methodologies have investigated chunking (Shukla et al., 2022; Moro and Ragazzi, 2022), extractive summarisation (Bajaj et al., 2021), multitask learning (Elnaggar et al., 2018), argument roles (Xu et al., 2021; Elaraby and Litman, 2022), and prompt engineering (Pont et al., 2023). Although promising experimental results exist, the literature on legal abstractive summarisation is still relatively small, with multi-document summarisation being particularly understudied.

### 2.2 Faithfulness and Hallucination

Abstractive summarisation can lead to more natural summaries, but it may also introduce content unsupported by the source text, known as 'hallucination' (Huang et al., 2023; Nan et al., 2021b; Ji et al., 2023). Faithfulness refers to the consistency of the generated text with the input text (Sridhar and Visser, 2022), so reducing hallucination corresponds to increasing faithfulness (Ji et al., 2023).

Only one existing work (Feijo and Moreira, 2023) attempts to tackle the problem of hallucination for legal domain summarisation. Feijo and Moreira (2023) propose the LegalSumm method where summaries are generated for multiple distinct chunks of the source text, and a textual entailment model scores chunk-summary pairs to select the most faithful summary. However, this approach has limitations, such as the training examples to

assess faithfulness not being reflective of real hallucination patterns, and the fact that the final summary derived from only one chunk may not include all salient information. Further, this method is not suitable for all judicial documents due to its use of specific case structure in the chunking process.

Various techniques to control hallucination have been proposed in the general domain, including filtering training examples (Matsumaru et al., 2020; Chaudhury et al., 2022), maximising faithfulness metrics during training (Nan et al., 2021b), modifying beam search (Zhao et al., 2020; Sridhar and Visser, 2022; Chaudhury et al., 2022; King et al., 2022), post-generation fact correction (Huang et al., 2023; Ji et al., 2023), and including additional information to guide generation (Dong et al., 2022; Cao et al., 2018; Narayan et al., 2022a, 2021).

## 3 Dataset

This study uses the Multi-LexSum dataset, the first dataset for legal multi-document summarisation (Shen et al., 2022). Multi-LexSum contains 9,280 expert-written summaries in accessible language pertaining to 4,539 U.S. civil rights lawsuits between 1950 and 2021, obtained from the Civil Rights Litigation Clearinghouse (CRLC). Case law was chosen for its practical application and volume. The documents to be summarised for each case include complaints, motions, and settlement agreements. Each of a case's documents can be over 100 pages, with a single case potentially involving hundreds of documents. A mean of 99378.2 words (10.3 documents) must be summarised per case, giving a very high compression ratio of 840.7.

Multi-LexSum contains multiple levels of summary granularity (examples in Appendix A); we focus on short summaries (mean 130 words), as long summaries (mean 646.5 words) frequently exceed the maximum decoder token length (1024) for the transformer model we use (PEGASUS). The short summaries cover the background, involved parties, and the case's outcome in a single paragraph. Writing summaries for standard cases takes 1-4 hours, while complex cases can take over 10 hours for an experienced lawyer.

Multi-LexSum is a relatively understudied dataset. Shen et al. (2022) conduct a preliminary study using Multi-LexSum using off-the-shelf models. While their results indicate that longer input lengths improve model performance, this is likely due to the content selection method used to handle

maximum token length resulting in salient information not being included in the model input, which may also have led to hallucination by weakening the coupling between input-summary pairs during training. Human evaluation suggested that an alternative content selection strategy could thus enhance model performance and reduce hallucination.

### 3.1 Preprocessing

We applied preprocessing steps to the noisy Multi-LexSum data initially extracted using OCR. We detail this process and provide an annotated example before and after cleaning in Appendix B. The cleaning process allows the source text to be correctly segmented into sentences and paragraphs, which is vital in our methodology. As data filtering has been shown to minimize hallucinations (Nan et al., 2021a; Ji et al., 2023; Dong et al., 2022; Chaudhury et al., 2022; Narayan et al., 2021), we removed cases where less than 75% of legally salient entities in the summary occur in the source text. Additionally, to augment the dataset, we integrate long summaries (Shen et al., 2022) which are under 671 words, the maximum length for the short summary subset. Table 1 shows the dataset splits after cleaning, filtering, and augmentation.

	Complete Dataset	Short Summaries (Original)	Short Summaries (Preprocessed)
Train	4,539	3,138	3,436
Val.	3,177 (70%)	2,210 (70%)	2,508 (73%)
Test	454 (10%)	312 (10%)	312 (9%)
Total	908 (20%)	616 (20%)	616 (18%)

Table 1: Size of dataset splits after preprocessing.

## 4 Overview

The chosen task of abstractive summarisation involves generating a short summary  $S_i$  of the a set of  $N$  documents denoted as  $D_i = \{D_{i_1}, D_{i_2}, \dots, D_{i_N}\}$  belonging to the same case. We concatenate the documents  $D_i$  for each case in chronological order; the dates of each document were scraped from the CRLC website as these were not generally extractable from the text.

## 5 Models

We use PEGASUS (Zhang et al., 2020a), a state-of-the-art sequence-to-sequence transformer encoder-decoder model as our backbone abstractive summarisation model. Jointly with the Masked Lan-

guage Modeling objective, PEGASUS has a pre-training objective designed specifically for abstractive summarisation - Gap Sentence Generation. Key sentences, selected based on ROUGE-F1, are masked from the input text during training, and the model must reproduce them; these key sentences are similar to a summary (Zhang et al., 2020a).

PEGASUS has a legal-pretrained variant, Legal-PEGASUS<sup>1</sup>, trained on U.S. case law. Pretraining on legal data has been shown to increase performance on legal NLP tasks (Zhong et al., 2020; Shukla et al., 2022; Niklaus and Giofre, 2023). Shen et al. (2022) report results on PEGASUS and LED-16384 (a sparse attention transformer able to handle input lengths of up to 16,348 tokens (Beltagy et al., 2020)) that we use as baselines.

## 6 Content Selection

The self-attention mechanism in transformers limits the input token length to 1024 for PEGASUS (Zhang et al., 2020a). While sparse attention transformers (Beltagy et al., 2020; Guo et al., 2022; Zaheer et al., 2020) can handle longer input sequences and have shown promising performance in general-domain summarisation (Chalkidis et al., 2022; Niklaus and Giofre, 2023), in legal and multi-document cases, the input text often exceeds even these limits; in Multi-LexSum, the average source text length for a case is 83,340 tokens, with a maximum of 4,423,683 tokens. Thus, a content selection strategy to ensure salient information is included within this input token limit is essential; if the input to the summarisation model does not contain the relevant information, summary quality is reduced and hallucination is encouraged, as the input-summary pairs are not tightly coupled.

Previous approaches to handling the input token limit include segmenting the source text in chunks and then concatenating summaries. However, this introduces a number of issues: it is non-trivial to extract the corresponding sentences from the reference summary for each chunk, not all chunks may be equally informative, independent chunk processing may lead to redundancy in the final summary, and for long input texts, summarising every chunk is computationally expensive (Shukla et al., 2022; Moro and Ragazzi, 2022). Similarly, multi-stage frameworks (Zhang et al., 2022), which iteratively use the concatenated summary as the input to another phase of chunking and abstractive summarisa-

<sup>1</sup><https://huggingface.co/nsi319/legal-pegasus>



tion, significantly increase computational complexity and introduce opportunities for hallucination.

To address these issues, we propose a mixed-model approach. We first identify salient information from the source text, and then use this information as input to our backbone model. Importantly in the legal domain, this approach better mirrors the human summarisation process, and hence may contribute to user trust. There is evidence that a human summarising long input text would highlight the important information and then paraphrase this information to form a summary (Bajaj et al., 2021; Norkute et al., 2021; Liu et al., 2018; Jing and McKeown, 1999), and a study on legal text summarisation demonstrates participants’ increased trust in systems for which they understand the summary’s creation process and feel that this process is similar to their own (Norkute et al., 2021; Danilevsky et al., 2020; Adadi and Berrada, 2018).

## 6.1 Oracle Extracts for Gold Labels

We adopt a ranking-based approach to select salient information by training a BERT-based salience classifier to extract relevant sentences from the source text, using the state-of-the-art OREO<sup>2</sup> method to obtain gold standard training labels. This enables us to create a list of all source text sentences ranked by the classifier’s confidence that the sentence contains salient information. During inference, the top-ranked sentences are utilized to construct the input to the PEGASUS model (when finetuning PEGASUS, we instead use the gold standard sentences from OREO as the model input).

To obtain ‘gold standard’ data regarding which sentences of the source text contain salient information for summarisation purposes, we must convert the gold-standard abstractive summaries to their extractive equivalent. As annotations by legal professionals would be prohibitively costly and time-consuming, we use an automatic labeling approach.

Various approaches have been proposed to create oracle extracts, among which greedily maximising the ROUGE overlap with the gold-standard summary is most common (Xu and Lapata, 2022; Bhattacharya et al., 2021; Klaus et al., 2022). However, oracles constructed in this way do not always lead to high-performing summaries (Xu and Lapata, 2022) - indeed, a recent study on legal extractive summarisation (Klaus et al., 2022) suggests that ‘alternative methods to create oracle extractive

summaries’ should be considered. Furthermore, this greedy approach considers only a *single* oracle summary,  $Y^*$ , but there can be *multiple* valid oracle summaries for the same source text; systems trained on greedy oracles are optimised by maximising the probability at  $Y^*$  and assigning zero probability to all other hypotheses, regardless of quality. For this reason, we use the OREO algorithm to create oracles, which incorporates the idea of learning from *multiple* oracle summary hypotheses. Xu and Lapata (2022) showed that OREO led to superior performance compared to the common greedy approach, and that OREO can better guide the learning and inference of an abstractive summarisation system. Further details and hyperparameter details are provided in Appendix C; here we note that OREO is fundamentally ROUGE based.

## 6.2 Sentence Salience Classification

Using the ‘oracle’ sentences output by OREO, we train a classifier to determine the summary-worthiness of sentences (i.e. the binary label assigned by OREO). Conceptually, as obtaining OREO labels requires a case to already have a gold-standard summary, training a classifier to *predict* which sentences of a legal case’s source text are summary-worthy (by training the classifier at a sentence level on OREO labels) allows us to carry out content selection on unseen cases.

We use a legal oriented pre-trained model, CaseLawBERT (Zheng et al., 2021), which provides the best domain match, and achieve an ROC-AUC score of 0.884 (curve in Appendix D). Due to class imbalances (Table 2), we conduct random downsampling, resulting in 68,592 training examples. However, addressing this imbalance in a more sophisticated manner may lead to improved results.

	All Instances	Positive Instances	Negative Instances
Train	6,230,772	34,296 (0.55%)	6,196,476
Val.	1,122,744	4,355 (0.39%)	1,118,389
Test	1,672,233	8,021 (0.48%)	1,664,212

Table 2: Number of instances of each class (binary, assigned by OREO) for sentence salience classification (before downsampling). Positive instances are considered summary-worthy.

## 6.3 Input Construction

To construct the PEGASUS inputs from the ranked list of sentences, we compare several strategies, adding tokens until the limit of 1024 is reached:

<sup>2</sup><https://github.com/yumoxu/oreo>

- Sentences - we add the top scoring sentences with non-zero scores, as in [Xu and Lapata \(2022\)](#).
- Windows - we add the preceding and following sentence for each selected sentence. This provides context, but may lead to irrelevant information being included.
- Paragraphs - we add the whole paragraph the sentence is contained within.

In all cases, we concatenate the extracted information in order of appearance in the temporally ordered source documents. We consider two variants: BERT (ranked list of sentences is obtained from training the BERT classifier on OREO labels), and OREO (ranked list of sentences is obtained directly from OREO, to investigate the *potential* gains our content selection strategy could produce if the classifier was perfectly predictive).

We also consider three baseline methods:

- First-1024 - we take the first 1024 tokens of the temporal concatenation of all the case’s source documents.
- First-K - like in the original MultiLex-sum paper ([Shen et al., 2022](#)), for a case with  $D$  documents, we take the first  $1024/D$  tokens of each. Unlike [Shen et al. \(2022\)](#), the dataset has been cleaned and temporally ordered.
- TextRank - a general-domain unsupervised extractive summarisation method, frequently used as a content selection baseline ([Liu et al., 2018](#); [Bajaj et al., 2021](#); [Klaus et al., 2022](#)).

#### 6.4 Content Selection Preliminary Results

As a preliminary experiment, we investigate the ROUGE recall between the *extracts* produced (which will be used as input to PEGASUS) and the corresponding gold standard summary for the test set (as we have already performed the expensive inference process for the BERT classifier on test set data). We use recall as we wish to consider if the salient information has been selected, not the specificity of salient information. Results are presented in Table 3.

BERT-Sentences and BERT-Windows outperform the naive First-1024 and First-K baselines, with TextRank also performing well. However, the BERT-Paragraphs method performs poorly, likely due to including too much context for each selected sentence and thus being able to include fewer

highly-ranked sentences. We also compare the three BERT-based strategies to their OREO counterparts. OREO-Windows performed best overall in terms of ROUGE-1 and ROUGE-2 recall. All OREO strategies outperformed their BERT-based counterparts in terms of ROUGE-2, although BERT-Sentences outperformed OREO-Sentences in terms of ROUGE-1. We also noted that OREO extracted significantly fewer tokens than BERT in the sentence case, suggesting that an input token length of 1024 tokens is sufficient.

	ROUGE-1	ROUGE-2	ROUGE-L
First-1024	67.51	24.35	41.50
First-K	57.36	19.25	35.94
TextRank	70.28	23.47	43.93
BERT-Sentences	<b>76.61</b>	<b>32.61</b>	<b>46.15</b>
BERT-Windows	73.88	28.00	41.95
BERT-Paragraphs	58.30	19.99	32.66
OREO-Sentences	68.07	32.73	35.43
OREO-Windows	<b>79.43</b>	<b>37.13</b>	<b>45.25</b>
OREO-Paragraphs	73.83	33.24	41.59

Table 3: Mean ROUGE recall scores against corresponding gold standard summary for each strategy tested.

## 7 Experimental Setup

Overall, we vary two dimensions in our experiments, corresponding to our research questions: input representation (RQ1, Section 6), and domain match (RQ2, Section 5).

For comparison to the PEGASUS results reported for Multi-LexSum in [Shen et al. \(2022\)](#), we use the same hyperparameters values where provided: we train for 6 epochs with a learning rate of  $5e-5$ , and for inference we use beam search with 5 beams and n-gram repetition blocks for  $n>3$ . For additional hyperparameters, we trained the models with a batch size of 4, 64 gradient accumulation steps, gradient checkpointing enabled, and a weight decay of 0.01. For our models at inference, we used a minimum of 24 tokens and maximum of 960 tokens for experimental settings with no entity chain, and a minimum of 34 tokens and maximum of 1154 for experimental settings including some form of entity chain, as these were the boundaries observed for our gold-standard data. We also added a length penalty of 2.0 to encourage the generation of long sequences, as [Shen et al. \(2022\)](#) observed that PEGASUS undergenerated the number of words when producing short summaries for Multi-LexSum.

## 8 Evaluation

We evaluate the quality of the produced summaries using standard ROUGE-1, ROUGE-2, and ROUGE-L scores. We also report BERTScore (Zhang et al., 2020b) to capture semantic similarity without relying solely on lexical overlap, as ROUGE fails to capture deeper semantic similarity (Shukla et al., 2022; Bhattacharya et al., 2019; Kanapala et al., 2019; Jain et al., 2023; Zhong et al., 2019; Cohan and Goharian, 2016; Kikuchi et al., 2014; Feijo and Moreira, 2023). We used the DeBERTA model for comparison with previous work.

We evaluate faithfulness using textual entailment following Narayan et al. (2022b,a) and previous studies demonstrating a correlation between entailment scores and human judgements of faithfulness (Narayan et al., 2022b; Fischer et al., 2022; Sridhar and Visser, 2022; Kryscinski et al., 2020; Maynez et al., 2020; Honovich et al., 2022). We report the probability of a generated summary (PEGASUS output) being entailed by its source text (PEGASUS input) returned by a BART-large classifier finetuned on Multi-NLI (Fischer et al., 2022).

## 9 Results and Analysis

### 9.1 Input Representation

To first investigate which content selection approaches are promising, for the standard (not legal pretrained) variant of PEGASUS, Table 4 shows ROUGE and BERTScore F1 scores for our baselines, the three BERT-based strategies, the three additional OREO-based strategies, the PEGASUS and state-of-the-art results reported in Shen et al. (2022), and our reproduction of the PEGASUS results (needed due to incomplete knowledge of hyperparameters used in Shen et al. (2022)).

Among BERT-based strategies, BERT-Windows, our most effective method and improved **ROUGE-1 by 0.82** compared to the reported PEGASUS performance. However, we do not observe improvements with respect to other metrics.

On all metrics apart from ROUGE-2, BERT-Windows was the most effective of the six tested input strategies. This is likely due to the balance of the number of relevant sentences included and providing context for each sentence. First-1024, First-K, BERT-Sentences, and BERT-Windows all outperform our reproduction of Shen et al. (2022)’s with respect to ROUGE-1, as expected, and BERT-Sentences and BERT-Windows outperform the reproduction baseline with respect to ROUGE-2. TextRank fails to outperform this baseline, which is consistent with its poor performance as a content selector for abstractive summarisation in Bajaj et al. (2021). BERT-Paragraph also fails to outperform the baseline, likely due to including longer context for each sentence, which limits the relevant information that can be included.

None of our 6 proposed strategies outperform the reproduction baseline on ROUGE-L or BERTScore metrics. We expected a greater improvement from the First-K baseline over the reproduction baseline, which intuitively should improve results as its only difference to the content selection strategy in Shen et al. (2022) is the introduction of dataset cleaning and temporal ordering. We hypothesise that the dataset filtering process may have resulted in decreased ROUGE scores<sup>3</sup>, consistent with Nan et al. (2021a) - although this is likely to contribute to increased faithfulness.

To investigate the potential of content selection strategies independently of the BERT classifier (i.e. if the BERT salience classifier was perfectly predictive of OREO labels), we also analyzed the model’s performance using the oracles from OREO as inputs. The OREO strategies outperformed BERT counterparts, with OREO-Sentences surpassing the SOTA by up to **4.45 ROUGE-1, 4.39 ROUGE-2, 1.40 ROUGE-L, and 0.27 BERTScore**.

As OREO-Sentences extracts typically consist of far fewer tokens (mean 264.15) than BERT-sentences (mean 1000.78), yet BERT-Sentences extracts have a greater ROUGE recall with the reference summary (see Section 6.4), this suggests that the specificity and saliency of the inputs provided to PEGASUS is key. Indeed, when measuring the mean of ROUGE-1 and ROUGE-2 precision between the OREO and BERT extracts used as input to PEGASUS with the gold summary, the OREO extracts display a greater precision (28.57 vs 7.67 for sentences, 10.68 vs 7.73 for windows, and 12.82 vs 12.15 for paragraphs). The increasing similarity in precision scores between OREO and BERT variants as the number of sentences for which information is included in the extracts decreases also suggests that the BERT classifier performs best for its high confidence outputs.

Overall, we establish that content selection does have the potential to improve summarisation out-

<sup>3</sup>This was *not* due to the augmentation process - we performed an ablation study without dataset augmentation for the Lead-K baseline and achieved poorer results: ROUGE-1 F1 42.56, ROUGE-2 F1 18.41, ROUGE-L F1 27.93.

	ROUGE-1	ROUGE-2	ROUGE-L	BERTScore	Entailment
<b>Baselines</b>					
PEGASUS (reproduced)	43.23	19.26	29.35	36.15	0.2937
PEGASUS	43.35	19.91	29.99	37.88	-
LED-16384 (SOTA)	46.54	22.08	31.91	40.00	-
<b>PEGASUS</b>					
First-1024	43.39	18.96	28.42	34.47	-
First-K	43.24	18.96	28.40	34.94	-
TextRank	42.36	17.23	27.31	33.45	-
BERT-Sentences	43.61	19.33	27.58	34.52	0.5134
BERT-Windows	44.17	19.28	28.53	35.62	0.5551
BERT-Paragraphs	40.14	16.28	25.95	31.39	-
OREO-Sentences	50.99	26.47	33.31	40.27	0.4915
OREO-Windows	47.97	23.28	31.55	38.92	0.5457
OREO-Paragraphs	47.15	22.42	30.83	37.84	-
<b>Legal-PEGASUS</b>					
BERT-Sentences	42.77	19.08	27.25	34.81	0.4954
BERT-Windows	44.34	19.55	28.91	36.35	0.5551
OREO-Sentences	52.10	27.54	34.61	42.15	0.4680
OREO-Windows	48.41	23.72	31.91	39.44	0.5469

Table 4: The upper part shows results for PEGASUS summaries, the lower part for Legal-PEGASUS summaries. We report Mean ROUGE and BERTScore F1 scores with respect to the corresponding reference summary. The last column shows entailment scores (not calculated for all experimental setups due to limited compute resources). We highlight the best scores for OREO and BERT in red and blue respectively, for PEGASUS and Legal-PEGASUS.

puts, but that the salience classifier performance limits these improvements in practice.

## 9.2 Domain-Specific Pretraining

As the sentence and window-based strategies offer the most promising results, we only report Legal-PEGASUS results for these strategies. The lower part of Table 4 shows results for Legal-PEGASUS. With legal pretraining, we observe improvements in BERTScore and ROUGE-1 for all input settings. Our best results for the complete pipeline are given by BERT-Windows. However, these results still only outperform the PEGASUS results reported in Shen et al. (2022) with respect to ROUGE-1, by **0.99 F1**. In contrast, OREO-Sentences further outperforms the state of the art, achieving an improvement of **5.56 ROUGE-1 F1**, **5.46 ROUGE-2 F1**, **2.7 ROUGE-L F1**, and **2.15 BERTScore**. Overall, we observed greater improvements for better content selection strategies. Our results again indicate the importance of content selection, and the importance of the domain match at pretraining.

## 9.3 Faithfulness

Entailment scores are reported in Table 4. While we do not have entailment scores for the exact PEGASUS setup in Shen et al. (2022) as we do not

have access to the original model outputs, and we acknowledge that our reproduction leads to slightly different results, it is evident that all our experimental setups vastly improve the probability of the source text entailing the summary text in comparison to this reproduction baseline (mean entailment probability 0.2937); our BERT-Windows content selection strategy improves entailment probability by 0.2614 for both PEGASUS and Legal-PEGASUS. This suggests that content selection is effective in improving summary faithfulness.

Interestingly, Legal-PEGASUS led to reduced entailment scores compared to vanilla PEGASUS. Overall, BERT-based methods consistently exhibited higher faithfulness than OREO-based counterparts, and window-based methods showed higher faithfulness than sentence-based methods. Our findings align with the literature in that ROUGE does not correlate with faithfulness; although OREO-Sentences receives the worst entailment scores, this method performs best on ROUGE and BERTScore.

## 9.4 Qualitative Analysis

Although human expert evaluation of the summaries is infeasible, to better understand our models' behaviour and failure modes, we manually analysed generated summaries for a sample of 10 cases



589 across experimental settings (example in Appendix  
590 G). In general, the outputs of the reproduction of  
591 the PEGASUS method in Shen et al. (2022) were  
592 comparatively good at reproducing the correct date  
593 when the case began, as this is frequently men-  
594 tioned at the start of the document. Background  
595 information for the case (often at the start of the  
596 initial document) are also reflected fairly reliably.  
597 However, the summaries often hallucinate the law  
598 which is alleged to be violated, which is extremely  
599 vital, and struggle to accurately represent the case’s  
600 procedure. This is likely as this information is not  
601 included in the input text captured using the naïve  
602 content selection strategy.

603 In contrast, our models produce longer sum-  
604 maries which better match the reference summaries  
605 in content (not reflected by the ROUGE results).  
606 We observed limited variation across input strate-  
607 gies, including OREO-based strategies. In general,  
608 our models perform well for the background and  
609 laws involved in the case, but performance often  
610 declines for a case’s procedural actions, with key  
611 information being missed, and the reasoning for de-  
612 cisions failing to be provided. This may be because  
613 these aspects follow a less standard format, so are  
614 less easily identified by the BERT classifier (Zhong  
615 et al., 2019). While our models contain less hal-  
616 lucinatory content than our reproduction of Shen  
617 et al. (2022)’s approach, two common hallucina-  
618 tion scenarios remain. Firstly, dates and monetary  
619 amounts which occur in the source text are often  
620 contained in the summary in the incorrect context;  
621 such intrinsic hallucination is non-trivial to combat.  
622 The second scenario stems from issues relating to  
623 case understanding. A notable subtype of this is the  
624 inclusion of information from *cited cases* as if it  
625 pertains to the main case under discussion; this may  
626 be because discussions of cited cases often include  
627 a high density of common legal keywords. As both  
628 BERT and OREO methods make this mistake, this  
629 suggests that selecting relevant sentences may be  
630 more suited to human annotation than automatic  
631 overlap-based methods; while full-scale human an-  
632 notation would be infeasible, semi-supervised ap-  
633 proaches, which have been applied with success  
634 in other areas of legal AI (Branting et al., 2019),  
635 may be promising. At the BERT classifier stage,  
636 including context for the sentence under considera-  
637 tion may help to distinguish between information  
638 relating to main or cited cases.

639 Another issue is our models’ tendency to include  
640 large extractive fragments from the source text, ev-

641 idenced by artefacts (such as numerals) remain-  
642 ing despite the cleaning process being reproduced.  
643 This limits the readability of the summaries in some  
644 cases by replicating complex legal terminology and  
645 syntax from the source text. This high degree of  
646 extractivity may be due to limited text occurring  
647 in the model input for each point, and that these  
648 fragments are not well-flowing text, which models  
649 such as PEGASUS are trained on.

## 650 10 Discussion and Conclusion

651 We conduct the first study at the intersection of  
652 legal, multi-document, and faithful summarisation,  
653 investigating the impact of content selection and  
654 legal pretraining on the abstractive summarisation  
655 of U.S. civil rights litigation, with PEGASUS as  
656 our backbone model. Our full test-time pipeline  
657 outperforms the PEGASUS results in Shen et al.  
658 (2022) by 0.99 ROUGE-1 F1. We show that using  
659 oracle extracts vastly outperforms the state-of-the-  
660 art, with legal pretraining further boosting results:  
661 we achieve an improvement of 5.56 ROUGE-1 F1,  
662 5.46 ROUGE-2 F1, 2.7 ROUGE-L F1, and 2.15  
663 BERTScore. Our content selection strategy also  
664 leads to an improvement of 0.2614 in the probabil-  
665 ity of a generated summary being entailed by its  
666 source text, compared to a naïve content selection  
667 baseline. Overall, we provide evidence that con-  
668 tent selection has the ability to improve summary  
669 faithfulness and quality. However, the generated  
670 summaries can still contain hallucinations and omit  
671 key information. Several issues, such as the qual-  
672 ity of the content selection method and addressing  
673 specific hallucination scenarios, remain to be ad-  
674 dressed for such automatic summarisation to see  
675 real-world adoption.

676 Our study’s limitations and error cases suggest  
677 several future research areas. Further research into  
678 content selection is promising - investigating meth-  
679 ods of content selection not fundamentally based on  
680 ROUGE, such as using human salience annotations  
681 in a semi-supervised framework, could be fruitful.  
682 More generally, the application of our generated  
683 summaries as the input to other legal NLP tasks  
684 could be studied. Also, future work could conduct  
685 similar investigations on different legal domains  
686 and jurisdictions, or using different backbone mod-  
687 els; for example, it would be interesting to observe  
688 the effect of content selection on models able to  
689 handle longer inputs, such as LED, or GPT-based  
690 models.



## 11 Limitations

Our project is limited by the resources available. Firstly, while presenting a realistic summarisation scenario, the OCR process used to construct the Multi-LexSum dataset from the original court documents introduces noise, which despite dataset cleaning, can adversely affect both summarisation outputs and the metrics used to evaluate these outputs. The availability of computational resources also limits the range of experiments that can be conducted and the hyperparameter settings used. Perhaps most importantly, the lack of a thorough human evaluation of our models' outputs by domain experts limits our interpretation of our findings, as metrics such as ROUGE and entailment are only *proxies* for summary quality and faithfulness. This factor is especially important in the legal domain, where a lack of correlation between automatic metrics and human expert judgements has previously been demonstrated (Shukla et al., 2022; Bhattacharya et al., 2019), and as the utility of automatic metrics to judge faithfulness remains a topic of research debate.

Our work also has limitations pertaining to the intended use case of legal summarisation - namely, by legal professionals or ordinary civilians without resources or expertise in machine learning. The powerful GPUs required for finetuning and performing inference for the transformer models used throughout our pipeline are unlikely to be available in non-academic environments. Furthermore, our methodology's performance on other datasets (for example - for other legal areas, jurisdictions, or languages) has not yet been tested.

## References

- Amina Adadi and Mohammed Berrada. 2018. [Peeking Inside the Black-Box: A Survey on Explainable Artificial Intelligence \(XAI\)](#). *IEEE Access*, 6:52138–52160. Conference Name: IEEE Access.
- Sophia Althammer, Arian Askari, Suzan Verberne, and Allan Hanbury. 2021. [DoSSIER@COLIEE 2021: Leveraging dense retrieval and summarization-based re-ranking for case law retrieval](#). ArXiv:2108.03937 [cs].
- A. A. Askari, S. V. Verberne, O. Alonso, S. Marchesin, M. Najork, and G. Silvello. 2021. [Combining lexical and neural retrieval with longformer-based summarization for effective case law retrieval](#). In *Proceedings of the second international conference on design of experimental search & information REtrieval systems*, pages 162–170. CEUR.
- Ahsaas Bajaj, Pavitra Dangati, Kalpesh Krishna, Pradhiksha Ashok Kumar, Rheeeya Uppaal, Bradford Windsor, Eliot Brenner, Dominic Dotterer, Rajarshi Das, and Andrew McCallum. 2021. [Long Document Summarization in a Low Resource Setting using Pre-trained Language Models](#). In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing: Student Research Workshop*, pages 71–80, Online. Association for Computational Linguistics.
- Claire Barale, Michael Rovatsos, and Nehal Bhuta. 2023. [Automated Refugee Case Analysis: A NLP Pipeline for Supporting Legal Practitioners](#). In *Findings of the Association for Computational Linguistics: ACL 2023*, pages 2992–3005, Toronto, Canada. Association for Computational Linguistics.
- Iz Beltagy, Matthew E. Peters, and Arman Cohan. 2020. [Longformer: The Long-Document Transformer](#). ArXiv:2004.05150 [cs].
- 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*, Lecture Notes in Computer Science, pages 413–428, Cham. Springer International Publishing.
- Paheli Bhattacharya, Soham Poddar, Koustav Rudra, Kripabandhu Ghosh, and Saptarshi Ghosh. 2021. [Incorporating domain knowledge for extractive summarization of legal case documents](#). In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, ICAIL '21, pages 22–31, New York, NY, USA. Association for Computing Machinery.
- Steven Bird and Edward Loper. 2004. [NLTK: The Natural Language Toolkit](#). In *Proceedings of the ACL Interactive Poster and Demonstration Sessions*, pages 214–217, Barcelona, Spain. Association for Computational Linguistics.
- Michael James Bommarito, Daniel Martin Katz, and Eric Detterman. 2018. [LexNLP: Natural Language Processing and Information Extraction For Legal and Regulatory Texts](#).
- K. Branting, B. Weiss, B. Brown, C. Pfeifer, A. Chakraborty, L. Ferro, M. Pfaff, and A. Yeh. 2019. [Semi-Supervised Methods for Explainable Legal Prediction](#). In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law*, ICAIL '19, pages 22–31, New York, NY, USA. Association for Computing Machinery.
- Ziqiang Cao, Furu Wei, Wenjie Li, and Sujian Li. 2018. [Faithful to the original: fact-aware neural abstractive summarization](#). In *Proceedings of the Thirty-Second AAAI Conference on Artificial Intelligence and Thirtieth Innovative Applications of Artificial Intelligence Conference and Eighth AAAI Symposium*

799				
800		on Educational Advances in Artificial Intelligence,	ings of the 29th International Conference on Com-	856
801		AAAI'18/IAAI'18/EAAI'18, pages 4784–4791, New	putational Linguistics, pages 6187–6194, Gyeongju,	857
		Orleans, Louisiana, USA. AAAI Press.	Republic of Korea. International Committee on Com-	858
			putational Linguistics.	859
802	Ilias Chalkidis, Manos Fergadiotis, Prodromos Malaka-			
803	siotis, Nikolaos Aletras, and Ion Androutsopoulos.		Ahmed Elnaggar, Christoph Gebendorfer, Ingo Glaser,	860
804	2020. <a href="#">LEGAL-BERT: The Muppets straight out of</a>		and Florian Matthes. 2018. <a href="#">Multi-Task Deep Learn-</a>	861
805	<a href="#">Law School</a> . In <i>Findings of the Association for Com-</i>		ing for Legal Document Translation, Summarization	862
806	<i>putational Linguistics: EMNLP 2020</i> , pages 2898–		and Multi-Label Classification. In <i>Proceedings of</i>	863
807	2904, Online. Association for Computational Lin-		<i>the 2018 Artificial Intelligence and Cloud Computing</i>	864
808	guistics.		<i>Conference, AICCC '18</i> , pages 9–15, New York, NY,	865
			USA. Association for Computing Machinery.	866
809	Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael			
810	Bommarito, Ion Androutsopoulos, Daniel Katz, and		Diego Feijo and Viviane Moreira. 2019. <a href="#">Summarizing</a>	867
811	Nikolaos Aletras. 2022. <a href="#">LexGLUE: A Benchmark</a>		<a href="#">Legal Rulings: Comparative Experiments</a> . In <i>Pro-</i>	868
812	<a href="#">Dataset for Legal Language Understanding in En-</a>		<i>ceedings of the International Conference on Recent</i>	869
813	<a href="#">glish</a> . In <i>Proceedings of the 60th Annual Meeting of</i>		<i>Advances in Natural Language Processing (RANLP</i>	870
814	<i>the Association for Computational Linguistics (Vol-</i>		2019), pages 313–322, Varna, Bulgaria. INCOMA	871
815	<i>ume 1: Long Papers</i> ), pages 4310–4330, Dublin,		Ltd.	872
816	Ireland. Association for Computational Linguistics.			
			Diego de Vargas Feijo and Viviane P. Moreira. 2023.	873
817	Subhajit Chaudhury, Sarathkrishna Swaminathan, Chu-		<a href="#">Improving abstractive summarization of legal rulings</a>	874
818	laka Gunasekara, Maxwell Crouse, Srinivas Rav-		<a href="#">through textual entailment</a> . <i>Artificial Intelligence</i>	875
819	ishankar, Daiki Kimura, Keerthiram Murugesan,		<i>and Law</i> , 31(1):91–113.	876
820	Ramón Fernandez Astudillo, Tahira Naseem, Pa-			
821	van Kapanipathi, and Alexander Gray. 2022. <a href="#">X-</a>		Tim Fischer, Steffen Remus, and Chris Biemann. 2022.	877
822	<a href="#">FACTOR: A Cross-metric Evaluation of Factual Cor-</a>		<a href="#">Measuring Faithfulness of Abstractive Summaries</a> .	878
823	<a href="#">rectness in Abstractive Summarization</a> . In <i>Proce-</i>		In <i>Proceedings of the 18th Conference on Natural</i>	879
824	<i>edings of the 2022 Conference on Empirical Methods</i>		<i>Language Processing (KONVENS 2022)</i> , pages 63–	880
825	<i>in Natural Language Processing</i> , pages 7100–7110,		73, Potsdam, Germany. KONVENS 2022 Organizers.	881
826	Abu Dhabi, United Arab Emirates. Association for			
827	Computational Linguistics.		Mandy Guo, Joshua Ainslie, David Uthus, Santiago On-	882
			tanon, Jianmo Ni, Yun-Hsuan Sung, and Yinfei Yang.	883
828	Arman Cohan and Nazli Goharian. 2016. <a href="#">Revisiting</a>		2022. <a href="#">LongT5: Efficient Text-To-Text Transformer</a>	884
829	<a href="#">Summarization Evaluation for Scientific Articles</a> . In		<a href="#">for Long Sequences</a> . In <i>Findings of the Associa-</i>	885
830	<i>Proceedings of the Tenth International Conference</i>		<i>tion for Computational Linguistics: NAACL 2022</i> ,	886
831	<i>on Language Resources and Evaluation (LREC'16)</i> ,		pages 724–736, Seattle, United States. Association	887
832	pages 806–813. European Language Resources As-		for Computational Linguistics.	888
833	sociation (ELRA).			
			Or Honovich, Roei Aharoni, Jonathan Herzig, Hagai	889
834	Shawn Curran, Sam Lansley, and Oliver Bethell.		Taitelbaum, Doron Kukliansy, Vered Cohen, Thomas	890
835	2023. <a href="#">Hallucination is the last thing you need</a> .		Scialom, Idan Szpektor, Avinatan Hassidim, and	891
836	ArXiv:2306.11520 [cs].		Yossi Matias. 2022. <a href="#">TRUE: Re-evaluating Factual</a>	892
			<a href="#">Consistency Evaluation</a> . In <i>Proceedings of the 2022</i>	893
837	Marina Danilevsky, Kun Qian, Ranit Aharonov, Yan-		<i>Conference of the North American Chapter of the</i>	894
838	nis Katsis, Ban Kawas, and Prithviraj Sen. 2020. <a href="#">A</a>		<i>Association for Computational Linguistics: Human</i>	895
839	<a href="#">Survey of the State of Explainable AI for Natural</a>		<i>Language Technologies</i> , pages 3905–3920, Seattle,	896
840	<a href="#">Language Processing</a> . In <i>Proceedings of the 1st Con-</i>		United States. Association for Computational Lin-	897
841	<i>ference of the Asia-Pacific Chapter of the Association</i>		guistics.	898
842	<i>for Computational Linguistics and the 10th Interna-</i>			
843	<i>tional Joint Conference on Natural Language Pro-</i>		Yichong Huang, Xiachong Feng, Xiaocheng Feng, and	899
844	<i>cessing</i> , pages 447–459, Suzhou, China. Association		Bing Qin. 2023. <a href="#">The Factual Inconsistency Prob-</a>	900
845	for Computational Linguistics.		<a href="#">lem in Abstractive Text Summarization: A Survey</a> .	901
			ArXiv:2104.14839 [cs].	902
846	Yue Dong, John Wieting, and Pat Verga. 2022. <a href="#">Faith-</a>			
847	<a href="#">ful to the Document or to the World? Mitigating</a>		Deepali Jain, Malaya Dutta Borah, and Anupam Biswas.	903
848	<a href="#">Hallucinations via Entity-Linked Knowledge in Ab-</a>		2023. <a href="#">Bayesian Optimization based Score Fusion</a>	904
849	<a href="#">stractive Summarization</a> . In <i>Findings of the Associa-</i>		<a href="#">of Linguistic Approaches for Improving Legal Doc-</a>	905
850	<i>tion for Computational Linguistics: EMNLP 2022</i> ,		<a href="#">ument Summarization</a> . <i>Knowledge-Based Systems</i> ,	906
851	pages 1067–1082, Abu Dhabi, United Arab Emirates.		264:110336.	907
852	Association for Computational Linguistics.			
			Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan	908
853	Mohamed Elaraby and Diane Litman. 2022. <a href="#">ArgLegal-</a>		Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea	909
854	<a href="#">Summ: Improving Abstractive Summarization of Le-</a>		Madotto, and Pascale Fung. 2023. <a href="#">Survey of Hal-</a>	910
855	<a href="#">gal Documents with Argument Mining</a> . In <i>Proce-</i>		<a href="#">lucination in Natural Language Generation</a> . <i>ACM</i>	911
			<i>Computing Surveys</i> , 55(12):248:1–248:38.	912

913	Hongyan Jing and Kathleen R. McKeown. 1999. <a href="#">The decomposition of human-written summary sentences</a> . In <i>Proceedings of the 22nd annual international ACM SIGIR conference on Research and development in information retrieval</i> , SIGIR '99, pages 129–136, New York, NY, USA. Association for Computing Machinery.	Jayawant N. Mandrekar. 2010. <a href="#">Receiver Operating Characteristic Curve in Diagnostic Test Assessment</a> . <i>Journal of Thoracic Oncology</i> , 5(9):1315–1316.	969
914			970
915			971
916			
917		Kazuki Matsumaru, Sho Takase, and Naoaki Okazaki. 2020. <a href="#">Improving Truthfulness of Headline Generation</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1335–1346, Online. Association for Computational Linguistics.	972
918			973
919			974
920	Ambedkar Kanapala, Sukomal Pal, and Rajendra Pamula. 2019. <a href="#">Text summarization from legal documents: a survey</a> . <i>Artificial Intelligence Review</i> , 51(3):371–402.		975
921			976
922			977
923			
924	Yuta Kikuchi, Tsutomu Hirao, Hiroya Takamura, Manabu Okumura, and Masaaki Nagata. 2014. <a href="#">Single Document Summarization based on Nested Tree Structure</a> . In <i>Proceedings of the 52nd Annual Meeting of the Association for Computational Linguistics (Volume 2: Short Papers)</i> , pages 315–320, Baltimore, Maryland. Association for Computational Linguistics.	Joshua Maynez, Shashi Narayan, Bernd Bohnet, and Ryan McDonald. 2020. <a href="#">On Faithfulness and Factuality in Abstractive Summarization</a> . In <i>Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics</i> , pages 1906–1919, Online. Association for Computational Linguistics.	978
925			979
926			980
927			981
928			982
929			983
930		Gianluca Moro and Luca Ragazzi. 2022. <a href="#">Semantic Self-Segmentation for Abstractive Summarization of Long Documents in Low-Resource Regimes</a> . <i>Proceedings of the AAI Conference on Artificial Intelligence</i> , 36(10):11085–11093. Number: 10.	984
931			985
932	Daniel King, Zejiang Shen, Nishant Subramani, Daniel S. Weld, Iz Beltagy, and Doug Downey. 2022. <a href="#">Don't Say What You Don't Know: Improving the Consistency of Abstractive Summarization by Constraining Beam Search</a> . In <i>Proceedings of the 2nd Workshop on Natural Language Generation, Evaluation, and Metrics (GEM)</i> , pages 555–571, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.		986
933			987
934		Ankan Mullick, Abhilash Nandy, Manav Kapadnis, Sohan Patnaik, Raghav R, and Roshni Kar. 2022. <a href="#">An Evaluation Framework for Legal Document Summarization</a> . In <i>Proceedings of the Thirteenth Language Resources and Evaluation Conference</i> , pages 4747–4753, Marseille, France. European Language Resources Association.	988
935			989
936			990
937			991
938			992
939			993
940			994
941	Svea Klaus, Ria Van Hecke, Kaweh Djafari Naini, Ismail Sengor Altinogvde, Juan Bernabé-Moreno, and Enrique Herrera-Viedma. 2022. <a href="#">Summarizing Legal Regulatory Documents using Transformers</a> . In <i>Proceedings of the 45th International ACM SIGIR Conference on Research and Development in Information Retrieval</i> , SIGIR '22, pages 2426–2430, New York, NY, USA. Association for Computing Machinery.	Emre Mumcuoğlu, Ceyhan E. Öztürk, Haldun M. Ozaktas, and Aykut Koç. 2021. <a href="#">Natural language processing in law: Prediction of outcomes in the higher courts of Turkey</a> . <i>Information Processing and Management: an International Journal</i> , 58(5).	995
942			996
943			997
944			998
945			999
946			1000
947		Feng Nan, Ramesh Nallapati, Zhiguo Wang, Cicero Nogueira dos Santos, Henghui Zhu, Dejiao Zhang, Kathleen McKeown, and Bing Xiang. 2021a. <a href="#">Entity-level Factual Consistency of Abstractive Text Summarization</a> . In <i>Proceedings of the 16th Conference of the European Chapter of the Association for Computational Linguistics: Main Volume</i> , pages 2727–2733, Online. Association for Computational Linguistics.	1001
948			1002
949			1003
950	Anastassia Kornilova and Vladimir Eidelman. 2019. <a href="#">BillSum: A Corpus for Automatic Summarization of US Legislation</a> . In <i>Proceedings of the 2nd Workshop on New Frontiers in Summarization</i> , pages 48–56, Hong Kong, China. Association for Computational Linguistics.		1004
951			1005
952			1006
953			1007
954			1008
955			
956	Wojciech Kryscinski, Bryan McCann, Caiming Xiong, and Richard Socher. 2020. <a href="#">Evaluating the Factual Consistency of Abstractive Text Summarization</a> . In <i>Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)</i> , pages 9332–9346, Online. Association for Computational Linguistics.	Feng Nan, Cicero Nogueira dos Santos, Henghui Zhu, Patrick Ng, Kathleen McKeown, Ramesh Nallapati, Dejiao Zhang, Zhiguo Wang, Andrew O. Arnold, and Bing Xiang. 2021b. <a href="#">Improving Factual Consistency of Abstractive Summarization via Question Answering</a> . In <i>Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)</i> , pages 6881–6894, Online. Association for Computational Linguistics.	1009
957			1010
958			1011
959			1012
960			1013
961			1014
962			1015
963	LexisNexis. 2024. <a href="#">Lawyers cross into the new era of generative AI</a> . Technical report, LexisNexis.		1016
964			1017
965	Peter J. Liu, Mohammad Saleh, Etienne Pot, Ben Goodrich, Ryan Sepassi, Lukasz Kaiser, and Noam Shazeer. 2018. <a href="#">Generating Wikipedia by Summarizing Long Sequences</a> . ArXiv:1801.10198 [cs].	Shashi Narayan, Joshua Maynez, Reinald Kim Amplayo, Kuzman Ganchev, Annie Louis, Fantine Huot, Dipanjan Das, and Mirella Lapata. 2022a. <a href="#">Conditional Generation with a Question-Answering Blueprint</a> . ArXiv:2207.00397 [cs].	1018
966			1019
967			1020
968			1021
			1022
			1023
			1024



1025	Shashi Narayan, Gonçalo Simões, Yao Zhao, Joshua	Alex Wang, Kyunghyun Cho, and Mike Lewis. 2020.	1083
1026	Maynez, Dipanjan Das, Michael Collins, and Mirella	Asking and Answering Questions to Evaluate the	1084
1027	Lapata. 2022b. A Well-Composed Text is Half Done!	Factual Consistency of Summaries. In <i>Proceedings</i>	1085
1028	Composition Sampling for Diverse Conditional Gen-	of the 58th Annual Meeting of the Association for	1086
1029	eration. In <i>Proceedings of the 60th Annual Meet-</i>	Computational Linguistics, pages 5008–5020, Online.	1087
1030	ing of the Association for Computational Linguistics	Association for Computational Linguistics.	1088
1031	(Volume 1: Long Papers), pages 1319–1339, Dublin,		
1032	Ireland. Association for Computational Linguistics.		
1033	Shashi Narayan, Yao Zhao, Joshua Maynez, Gonçalo	Thomas Wolf, Lysandre Debut, Victor Sanh, Julien	1089
1034	Simões, Vitaly Nikolaev, and Ryan McDonald. 2021.	Chaumond, Clement Delangue, Anthony Moi, Pier-	1090
1035	Planning with Learned Entity Prompts for Abstrac-	ric Cistac, Tim Rault, Remi Louf, Morgan Funtowicz,	1091
1036	tive Summarization. <i>Transactions of the Association</i>	Joe Davison, Sam Shleifer, Patrick von Platen, Clara	1092
1037	for Computational Linguistics	Ma, Yacine Jernite, Julien Plu, Canwen Xu, Teven	1093
	, 9:1475–1492.	Le Scao, Sylvain Gugger, Mariama Drame, Quentin	1094
1038	Joel Niklaus and Daniele Giofre. 2023. Can we Pretrain	Lhoest, and Alexander Rush. 2020. Transformers:	1095
1039	a SotA Legal Language Model on a Budget From	State-of-the-Art Natural Language Processing. In	1096
1040	Scratch? In <i>Proceedings of The Fourth Workshop</i>	<i>Proceedings of the 2020 Conference on Empirical</i>	1097
1041	on Simple and Efficient Natural Language Process-	Methods in Natural Language Processing: System	1098
1042	ing (SustainLP), pages 158–182, Toronto, Canada	Demonstrations, pages 38–45, Online. Association	1099
1043	(Hybrid). Association for Computational Linguistics.	for Computational Linguistics.	1100
1044	Milda Norkute, Nadja Herger, Leszek Michalak, An-	Huihui Xu, Jaromir Savelka, and Kevin D. Ashley. 2021.	1101
1045	drew Mulder, and Sally Gao. 2021. Towards Ex-	Toward summarizing case decisions via extracting	1102
1046	plainable AI: Assessing the Usefulness and Impact	argument issues, reasons, and conclusions. In <i>Pro-</i>	1103
1047	of Added Explainability Features in Legal Document	ceedings of the Eighteenth International Conference	1104
1048	Summarization. In <i>Extended Abstracts of the 2021</i>	on Artificial Intelligence and Law, ICAIL ’21, pages	1105
1049	CHI Conference on Human Factors in Computing	250–254, New York, NY, USA. Association for Com-	1106
1050	Systems, CHI EA ’21, pages 1–7, New York, NY,	puting Machinery.	1107
1051	USA. Association for Computing Machinery.		
1052	Vedant Parikh, Vidit Mathur, Parth Mehta, Namita Mit-	Yumo Xu and Mirella Lapata. 2022. Text Summariza-	1108
1053	tal, and Prasenjit Majumder. 2021. LawSum: A	tion with Oracle Expectation. ArXiv:2209.12714	1109
1054	weakly supervised approach for Indian Legal Docu-	[cs].	1110
1055	ment Summarization. ArXiv:2110.01188 [cs].		
1056	Thiago Dal Pont, Federico Galli, Andrea Loreggia,	Manzil Zaheer, Guru Guruganesh, Kumar Avinava	1111
1057	Giuseppe Pisano, Riccardo Rovatti, and Giovanni	Dubey, Joshua Ainslie, Chris Alberti, Santiago On-	1112
1058	Sartor. 2023. Legal Summarisation through LLMs:	tanon, Philip Pham, Anirudh Ravula, Qifan Wang,	1113
1059	The PRODIGIT Project. ArXiv:2308.04416 [cs].	Li Yang, and Amr Ahmed. 2020. Big Bird: Trans-	1114
		formers for Longer Sequences. In <i>Advances in</i>	1115
1060	Zejiang Shen, Kyle Lo, Lauren Yu, Nathan Dahlberg,	<i>Neural Information Processing Systems</i> , volume 33,	1116
1061	Margo Schlanger, and Doug Downey. 2022. Multi-	pages 17283–17297. Curran Associates, Inc.	1117
1062	LexSum: Real-world Summaries of Civil Rights	Jingqing Zhang, Yao Zhao, Mohammad Saleh, and	1118
1063	Lawsuits at Multiple Granularities. In <i>Thirty-sixth</i>	Peter J. Liu. 2020a. PEGASUS: pre-training with	1119
1064	Conference on Neural Information Processing Sys-	extracted gap-sentences for abstractive summariza-	1120
1065	tems Datasets and Benchmarks Track.	tion. In <i>Proceedings of the 37th International Confer-</i>	1121
		<i>ence on Machine Learning</i> , volume 119 of <i>ICML’20</i> ,	1122
1066	Abhay Shukla, Paheli Bhattacharya, Soham Poddar, Ra-	pages 11328–11339. JMLR.org.	1123
1067	ajdeep Mukherjee, Kripabandhu Ghosh, Pawan Goyal,	Tianyi Zhang, Varsha Kishore, Felix Wu, Kil-	1124
1068	and Saptarshi Ghosh. 2022. Legal Case Document	ian Q. Weinberger, and Yoav Artzi. 2020b.	1125
1069	Summarization: Extractive and Abstractive Meth-	BERTScore: Evaluating Text Generation with BERT.	1126
1070	ods and their Evaluation. In <i>Proceedings of the 2nd</i>	ArXiv:1904.09675 [cs].	1127
1071	Conference of the Asia-Pacific Chapter of the Asso-		
1072	ciation for Computational Linguistics and the 12th	Yusen Zhang, Ansong Ni, Ziming Mao, Chen Henry Wu,	1128
1073	International Joint Conference on Natural Language	Chenguang Zhu, Budhaditya Deb, Ahmed Awadallah,	1129
1074	Processing (Volume 1: Long Papers), pages 1048–	Dragomir Radev, and Rui Zhang. 2022. Summ^N:	1130
1075	1064, Online only. Association for Computational	A Multi-Stage Summarization Framework for Long	1131
1076	Linguistics.	Input Dialogues and Documents. In <i>Proceedings</i>	1132
		of the 60th Annual Meeting of the Association for	1133
1077	Arvind Krishna Sridhar and Erik Visser. 2022. Im-	Computational Linguistics (Volume 1: Long Papers),	1134
1078	proved Beam Search for Hallucination Mitigation in	pages 1592–1604, Dublin, Ireland. Association for	1135
1079	Abstractive Summarization. ArXiv:2212.02712 [cs].	Computational Linguistics.	1136
1080	Kazem Taghva, Tom Nartker, Allen Condit, and Julie	Zheng Zhao, Shay B. Cohen, and Bonnie Webber.	1137
1081	Borsack. Automatic Removal of “Garbage Strings”	2020. Reducing Quantity Hallucinations in Abstrac-	1138
1082	in OCR Text: An Implementation.	tive Summarization. In <i>Findings of the Association</i>	1139



1140 *for Computational Linguistics: EMNLP 2020*, pages  
1141 2237–2249, Online. Association for Computational  
1142 Linguistics.

1143 Lucia Zheng, Neel Guha, Brandon R. Anderson, Peter  
1144 Henderson, and Daniel E. Ho. 2021. [When does  
1145 pretraining help? assessing self-supervised learning  
1146 for law and the CaseHOLD dataset of 53,000+ legal  
1147 holdings](#). In *Proceedings of the Eighteenth Interna-  
1148 tional Conference on Artificial Intelligence and Law,  
1149 ICAIL '21*, pages 159–168, New York, NY, USA.  
1150 Association for Computing Machinery.

1151 Haoxi Zhong, Chaojun Xiao, Cunchao Tu, Tianyang  
1152 Zhang, Zhiyuan Liu, and Maosong Sun. 2020. [How  
1153 Does NLP Benefit Legal System: A Summary of Le-  
1154 gal Artificial Intelligence](#). In *Proceedings of the 58th  
1155 Annual Meeting of the Association for Computational  
1156 Linguistics*, pages 5218–5230, Online. Association  
1157 for Computational Linguistics.

1158 Linwu Zhong, Ziyi Zhong, Zinian Zhao, Siyuan Wang,  
1159 Kevin D. Ashley, and Matthias Grabmair. 2019. [Au-  
1160 tomatic Summarization of Legal Decisions using It-  
1161 erative Masking of Predictive Sentences](#). In *Proceeed-  
1162 ings of the Seventeenth International Conference on  
1163 Artificial Intelligence and Law, ICAIL '19*, pages  
1164 163–172, New York, NY, USA. Association for Com-  
1165 puting Machinery.

## 1166 A Summary Granularities in 1167 Multi-LexSum

1168 Here we present examples of the three summary  
1169 granularities in Multi-LexSum: long, short, and  
1170 tiny.

1171 *Source Input Excerpt ...* And, even if the agency  
1172 had made an internal decision to maintain the sta-  
1173 tus quo, the documents at issue would not lose  
1174 their predecisional status because plaintiff has not  
1175 shown that they have been “adopted, formally or  
1176 informally, as the agency position on an issue or  
1177 is used by the agency in its dealings with the pub-  
1178 lic.”<sup>1</sup> *Coastal States Gas Corp.*, 617 F.2d at 866;  
1179 *Sears*, 421 U.S. at 161 (“[I]f an agency chooses  
1180 expressly to adopt or incorporate by reference an  
1181 intraagency memorandum previously covered by  
1182 Exemption 5 in what would otherwise be a final  
1183 opinion” that memorandum may not be withheld  
1184 under Exemption 5). Plaintiff does not point to any  
1185 public statements that OMB has made referencing,  
1186 adopting, or incorporating the records or the sub-  
1187 ject matter at issue, nor has plaintiff provided the  
1188 Court with any evidence that the records were in-  
1189 formally adopted as the agency’s position. Plaintiff  
1190 references a statement made by Karen Battle, chief  
1191 of the Census Bureau’s Population Division, on Jan-  
1192 uary 26, 2018, where she explained that additional

1193 research and testing were necessary before the Cen-  
1194 sus Bureau could proceed to implement a separate  
1195 Middle Eastern or North African category. Pl.’s  
1196 Cross-Mem. at 13. Plaintiff argues that “[t]o the  
1197 extent that Ms. Battle’s explanation about the need  
1198 for more research, and indeed the entire underlying  
1199 decision to maintain the status quo, is evidenced  
1200 in the withheld documents, it has been adopted as  
1201 the agency’s policy.” Id. But, this statement was  
1202 made by a Census Bureau official, not an OMB  
1203 official. And, in any event, the statement 1 Courts  
1204 in this district have held that the plaintiff carries  
1205 the burden to show that the agency has formally  
1206 or informally adopted a record as policy. See, e.g.,  
1207 *Heffernan v. Azar*, 317 F. Supp. 3d 94, 122 (D.D.C.  
1208 2018), citing *Sec. Fin. Life Ins. Co. v. U.S. Dep’t*  
1209 *of Treasury*, No. 03-102, 2005 WL 839543, at \*7  
1210 (D.D.C. Apr. 12, 2005). ...

1211 *Long Summary:* On April 13, 2018, the Arab  
1212 American Institute (“AAI”) sued the Office of Man-  
1213 agement and Budget (“OMB”) under the Freedom  
1214 of Information Act (“FOIA”), 5 U.S.C. § 552, in  
1215 the U.S. District Court for the District of Columbia.  
1216 AAI alleged that OMB violated FOIA by failing  
1217 to disclose requested records pertaining to OMB’s  
1218 decision not to include a combined race and ethnic-  
1219 ity question or a Middle Eastern or North African  
1220 (MENA) category on the 2020 Census. AAI asked  
1221 the court to declare that OMB violated FOIA and  
1222 to issue an injunction ordering the agency to re-  
1223 lease the requested records. This case was assigned  
1224 to Judge Amy Berman Jackson One month later,  
1225 on May 18, 2018, the court ordered OMB to file a  
1226 dispositive motion or a status report setting a sched-  
1227 ule for OMB’s production of documents to AAI.  
1228 OMB chose the latter, filing its first status report  
1229 on June 15, 2018. Over the next two years, the par-  
1230 ties filed several joint status reports detailing which  
1231 documents OMB had disclosed to AAI and which  
1232 documents were still outstanding or disputed. By  
1233 May 13, 2020, OMB had reviewed approximately  
1234 2,000 potentially responsive documents, produc-  
1235 ing “a number” of them to AAI and withholding  
1236 161 of them, claiming they were FOIA exempt.  
1237 AAI objected to the withholding of five of the al-  
1238 legedly exempt documents. OMB filed a motion  
1239 for summary judgment on February 10, 2020, argu-  
1240 ing that the five disputed documents were exempt  
1241 under FOIA Exemption 5, which allows agencies to  
1242 withhold “inter-agency or intra-agency memoran-  
1243 dums or letters that would not be available by law  
1244 to a party other than an agency in litigation with

the agency,” including “predecisional and deliberative” documents that reflect internal Executive Branch deliberations. AAI filed a cross-motion for summary judgment on March 12, 2020, arguing that OMB had not provided a sufficient basis for exempting the documents and that the exemption didn’t apply because the documents were not “predecisional.” On August 13, 2020, after conducting in camera review, the court granted OMB’s motion for summary judgment and denied AAI’s cross-motion, finding that the disputed documents were predecisional and exempt from FOIA. 2020 WL 4698098. As of December 25, 2020, AAI has not appealed the court’s decision.

*Short Summary:* On April 13, 2018, the Arab American Institute sued the Office of Management and Budget under the Freedom of Information Act in the U.S. District Court for the District of Columbia. AAI alleged that OMB violated FOIA by failing to disclose requested records pertaining to OMB’s decision not to include a combined race and ethnicity question or a Middle Eastern or North African (MENA) category on the 2020 Census. In May, the court ordered OMB to file a dispositive motion or a status report setting a schedule for OMB’s production of documents to AAI. Over the next two years, the parties filed several joint status reports detailing which documents OMB had disclosed to AAI and which documents were still outstanding or disputed. OMB produced a number of documents to AAI but withheld some, claiming they were FOIA exempt. AAI objected to five claimed exemptions. The parties both filed motions for summary judgment. After conducting in camera review, on August 13, 2020, the court granted OMB’s motion for summary judgment and denied AAI’s cross-motion, finding that the disputed documents were predecisional and exempt from FOIA. As of December 25, 2020, AAI has not appealed the court’s decision.

*Tiny Summary:* The Office of Management and Budget is forced to disclose documents requested by the Arab American Institute under the Freedom of Information Act. (D.D.C.)

## B Document Cleaning

The use of OCR (as required in real-world scenarios) to obtain plain text data from PDF court documents (Shen et al., 2022) of variable legibility containing formatting such as headers, footnotes, citations, and tables results in the source text in

the Multi-LexSum dataset containing errors and noise. Therefore, despite the underlying quality of the judicial documents, we first conducted dataset cleaning to allow for subsequent steps such as segmentation to be meaningfully applied, as in many cases we find ‘junk’ in the middle of paragraphs or sentences, and erroneous line breaks.

The overall cleaning pipeline for each source document is illustrated in Figure B. To define the rules for cleaning, we studied the text in the Multi-LexSum dataset and the corresponding original documents available on the CRLC website for cases in the validation set. For each newly implemented rule, we tested their validity on subsequent documents in the validation set, and ensured that previously considered documents were not adversely affected. This process continued until a stable set of rules was reached, which was then applied to all source documents.

- Removal of footers: We removed document footers containing irrelevant entities.
- Removal of headers: We only keep lines meeting at least one of the following conditions:
  - The line stripped of numerals only occurs once in the document - headers occur multiple times in the document, but may contain page numbers; thus, when stripped of numerals, this stripped line occurs multiple times.
  - The length of the stripped line is less than 20 characters - headers are long, we do not want to remove other information which may be repeated throughout the document, such as names, or terms such as ‘v.’ or ‘and’.
  - The line does not contain any numerals or hyperlinks - headers usually contain one or both, and we do not want to remove useful information.
- Removal of dirty lines: Dirty lines include page numbers, hyperlinks, lines not containing alphabetical characters, timestamps, and ‘junk’ resulting from OCR. Timestamp lines were identified using the *dateutils* parser. To remove ‘junk’ lines resulting from the OCR process, we edited *garbage\_detector*<sup>4</sup> (based on Taghva et al.), which identifies a line of text as ‘garbage’ if any one of several given conditions holds. We

<sup>4</sup><https://github.com/foodoh/rmgarbage>



Figure 1: Summary of main stages of the cleaning pipeline.

removed two of the conditions originally provided, as these gave many false positives in the legal domain: uppercase between lowercase; two distinct punctuation marks in the same line. We kept the remaining three original conditions, relating to a string’s ratio of alphanumeric characters to total characters, ratio of consonants to vowels, and if a punctuation mark repeats consecutively (this condition was edited to reflect the fact that while periods and brackets can legitimately repeat consecutively, punctuation marks such as commas, colons, semicolons, and dollars cannot). Finally, we added a condition to capture the fact that certain punctuation marks appearing between lower-case letters is indicative of junk text.

- Line breaks: This includes removing blank lines, removing newlines in the middle of sentences or paragraphs, and correctly ensuring a newline before each new legal paragraph. We kept existing line breaks only after colons (used to precede legal lists), after periods where the previous character was not a capital letter or ‘v’ (to avoid line breaks after abbreviations such as v. or U.S.), or if the whole line consisted of upper case letters (indicative of a section title). To insert the correct line breaks between legal paragraphs, in judicial documents of ‘standard’ format a new legal paragraph can be identified by a numeral or letter (in the case of lists) followed by a period. At this phase, we had to consider a number of special cases. For example, we do not insert a newline after a colon if the colon is not followed by whitespace, so as not to insert a newline in the middle of a hyperlink.
- Clean remaining lines: We remove footnotes and floating punctuation.
- Additional docket processing: Docket documents have a distinct format to judicial documents of other types. In particular, dockets contain tables with two columns giving the date (left), and the action taking place (right), which are not well represented in plain text format. For dockets, we remove lines consisting solely of dates (the

left column of the table), and numbers at the start of lines, as this is noise from attempting to linearise the table. In the vast majority of cases, no information is lost as the corresponding date is included in the main column entry.

- Address line breaks: Removing junk information often allows us to retrieve the correct line breaks. For docket documents, this phase is different, as due to the text originally being table cells, newline characters separate sentences.

An annotated representative excerpt from a case document before and after cleaning is given in Figures 2 and 3 respectively. This displays the effectiveness of our cleaning pipeline, however we note that cleaning cannot be perfect in all cases since documents have different formats and levels of OCR noise, and we do not want to erroneously remove valid text.

The cleaning process allows the source text to be correctly segmented into sentences and paragraphs, which is vital for subsequent stages in our methodology. The newline stages of the cleaning process allow for correct paragraph segmentation. For sentence segmentation, we use LexNLP (Bommarito et al., 2018) as this is specifically designed for legal text. Despite this, we still found that some post-processing was required to achieve the best results as certain cases were not well handled. Following (Parikh et al., 2021), we merge a sentence with the previous sentence if the previous sentence ends in an acronym (such as ‘v.’), or if the current sentence begins with ‘Section’ (to address incorrect segmentation within legal articles). We also introduce a sentence boundary between ‘;’ and ‘(’ to segment long legal lists. For docket type documents, as there is no period at the end of entries in table cells, we must first divide the text into paragraphs, which correspond to each cell of the table, before applying sentence segmentation to each paragraph.

With respect to data filtering, we filter out training examples with low entity extractivity to discourage hallucination, as training examples which are unfaithful to the source text can encourage generative models to produce hallucinations (Nan et al.,

1386  
1387  
1388  
1389  
1390  
1391  
1392  
1393  
1394  
1395  
1396  
1397  
1398  
1399  
1400  
1401  
1402  
1403  
1404  
1405  
1406  
1407  
1408  
1409  
1410  
1411  
1412  
1413  
1414  
1415  
1416  
1417  
1418  
1419  
1420  
1421  
1422  
1423  
1424  
1425  
1426  
1427  
1428  
1429



UNITED STATES DISTRICT COURT WESTERN DISTRICT OF MICHIGAN  
SOUTHERN DIVISION

" R-,:P ";9 ← Floating Junk From OCR

EQUAL EMPLOYMENT OPPORTUNITY COMMISSION

Plaintiff,  
Vo

-4 Honorable Swene~io.dr,eJUL .AS..NDillest~rict ,Judge

ROBERT BOSCH CORPORATION

COMPLAINT AND JURY DEMAND

Defendant. / ← Floating Punctuation

NATURE OF TITLE ACTION

This is an action under Title VII of the Civil Rights Act of 1964 and Title I of the Civil ← Line Breaks Within Sentence / Paragraph

Rights Act of 1991 to correct unlawful employment practices on the bases of religion, and to

provide appropriate relief to Jeff Carter who was adversely affected by such practices. The

United States Equal Employment Opportunity Commission (hereinafter "EEOC") alleges that Robert Bosch Corporation (hereinafter "Defendant") failed to provide a reasonable

accommodation to the known religious practices of Carter, who is a member of the International Old Path Church of God Inc. The Defendant unlawfully terminated Carter because his religious practices conflicted with an employment requirement. ← Missing Line Break After Section Title

JURISDICTION AND VENUE 1. Jurisdiction of this Court is invoked pursuant to 28 U.S.C. § 451, 1331, 1337, 1343, and 1345. This action is authorized and instituted pursuant to Section 706(0)(1) and (3) and 707(e) of Title VII of the Civil Rights Act of 1964, as amended, 42 U.S.C. § 2000e-5(f)(1) and (3) and § 2000e-6(e) ("Title VII"), and Section 102 of the Civil Rights Act of 1991, 42

1 ← Floating Page Number Dividing Paragraph

U.S.C. § 1981a. 2. The employment practices alleged to be unlawful were committed within the jurisdiction of the United States District Court for the Western District of Michigan, Southern Division.

PARTIES 3. Plaintiff, EEOC is the agency of the United States of

Figure 2: Annotated representative excerpt, before cleaning process.

UNITED STATES DISTRICT COURT WESTERN DISTRICT OF MICHIGAN  
SOUTHERN DIVISION  
EQUAL EMPLOYMENT OPPORTUNITY COMMISSION  
Plaintiff, Vo  
ROBERT BOSCH CORPORATION  
COMPLAINT AND JURY DEMAND ← Some Noise From OCR Remains

Defendant.

NATURE OF TITLE ACTION This is an action under Title VII of the Civil Rights Act of 1964 and Title I of the Civil Rights Act of 1991 to correct unlawful employment practices on the bases of religion, and to provide appropriate relief to Jeff Carter who was adversely affected by such practices. The United States Equal Employment Opportunity Commission (hereinafter "EEOC") alleges that Robert Bosch Corporation (hereinafter "Defendant") failed to provide a reasonable accommodation to the known religious practices of Carter, who is a member of the International Old Path Church of God Inc. The Defendant unlawfully terminated Carter because his religious practices conflicted with an employment requirement.

JURISDICTION AND VENUE

1. Jurisdiction of this Court is invoked pursuant to 28 U.S.C. § 451, 1331, 1337, 1343, and 1345. This action is authorized and instituted pursuant to Section 706(0)(1) and (3) and 707(e) of Title VII of the Civil Rights Act of 1964, as amended, 42 U.S.C. § 2000e-5(f)(1) and (3) and § 2000e-6(e) ("Title VII"), and Section 102 of the Civil Rights Act of 1991, 42 U.S.C. § 1981a.

2. The employment practices alleged to be unlawful were committed within the jurisdiction of the United States District Court for the Western District of Michigan, Southern Division.

PARTIES

3. Plaintiff, EEOC is the agency of the United States of America charged with the administration, interpretation and enforcement of Title VII, and is expressly authorized to bring this action by Section 706(f)(1) and (3) and 707(e) of Title VII, 42 U.S.C. § 2000e-5(f)(1) and (3) and § 2000e-6(e).

4. At all relevant times, Defendant has continuously been a corporation doing business in the State of Michigan, and has continuously had at least 15 employees.

5. At all relevant times, Defendant has continuously been an employer engaged in an industry affecting commerce within the meaning of Sections 701(b), (g) and (h) of Title VII, 42 U.S.C. §§ 2000e(b), (g) and (h).

STATEMENT OF CLAIMS

6. More than thirty days before the institution of this lawsuit, Carter filed a Charge of Discrimination with the Commission alleging violations of Title VII by the Defendant. All conditions precedent to the institution of this lawsuit have been fulfilled.

7. Since at least February, 2002, Defendant Employer has engaged in unlawful employment practices at its Saint Joseph, Michigan facility, in violation of Section 703(a), 42 U.S.C. § 2000e-2(a), and Section 704(a), 42 U.S.C. § 2000e-3(a). The Defendant's unlawful employment practices include the unlawful failure to provide a reasonable accommodation to the known sincerely-held religious beliefs of Carter, to wit: the belief that he should not work on

Figure 3: Annotated representative excerpt, after cleaning process.

2021a; Ji et al., 2023; Dong et al., 2022; Chaudhury et al., 2022; Narayan et al., 2021). While the summaries in Multi-LexSum are expertly constructed and faithful to the source documents as on the CRLC website, the OCR process means that not all documents are adequately represented by the plain text format in Multi-LexSum - for example, the dataset contains handwritten source documents for which the OCR software struggles to extract any text. Therefore, the Multi-LexSum dataset contains cases where the source text does not contain key information in the summary; these cases should be removed.

We based our filtering on verifying if the named entities in the summary occur in the source text. Firstly, in order to conduct the named entity recognition (NER), we use a state-of-the-art NER system (Barale et al., 2023) developed specifically for the legal domain in collaboration with legal professionals. The NER model was trained on human-annotated Canadian refugee law cases, fine-tuning LEGAL-BERT (Chalkidis et al., 2020). We include standard NER categories (DATE, PERSON, GPE, ORG, NORP, LAW) and the CLAIMANT\_INFO legal-specific category. Additionally, we added the MONEY category from LexNLP (Bommarito

et al., 2018). We manually evaluated results of the LEGAL-BERT NER systems on a subset of the validation set, studying the performance and relevance of all categories. While overall the NER system performed well, we found one common error for the GPE and ORG categories - the system included additional words between two true entities, resulting in one false entity (eg 'AT&T employee against AT&T Corp.') being returned. To solve this, we used the NLTK (Bird and Loper, 2004) part-of-speech tagger to postprocess these categories, removing words which were not nouns, adjectives, 'in', or 'of' from the entity and segmenting at the newly created boundaries.

Our filtering was based on verifying if the entities extracted from the gold standard summary appeared in the source text. However, matching named entities is nontrivial (Nan et al., 2021a), with several recurring scenarios causing difficulty:

- Dates - the same date can occur in different formats. We adopted a very optimistic approach to filter out obvious errors, however we note that this may give false positives by indicating entities are extractive when they are not. To deal with



1480 generalisations, such as ‘September 2003’ occur-  
 1481 ring in the summary while the source documents  
 1482 may only contain specific dates (i.e. - the day  
 1483 of the month is also specified), we parsed such  
 1484 expressions into multiple date formats and at-  
 1485 tempted to find a match in the source text for any  
 1486 of these formats, for any day of the month. Sim-  
 1487 ilarly, for expressions such as ‘early 2003’ we  
 1488 solely attempted to verify the year. For relative  
 1489 expressions such as ‘the next day’, we optimisti-  
 1490 cally assumed these were valid.

- 1491 • Paraphrases - for example, ‘AT&T employee’ and  
 1492 ‘employed by AT&T Corp.’.
- 1493 • Expansion and contraction of abbreviations - for  
 1494 example, ‘Corporation’ and ‘Corp.’. Creating a  
 1495 dictionary to match all such abbreviations would  
 1496 be infeasible.
- 1497 • Minor errors such as inconsistent spacing and  
 1498 punctuation.

1499 We note that many of these issues occur due to  
 1500 basing our matching on an exact match of surface  
 1501 forms. While we considered strategies such as  
 1502 fuzzy string matching, we found this to lead to  
 1503 worse results, as for example, changing one letter  
 1504 is very important when referring to legal articles,  
 1505 but could still lead to a fuzzy string match with  
 1506 high confidence. Overall, while our method is not  
 1507 reliable at the level of individual entities, we found  
 1508 through manual inspection that our method suffices  
 1509 to filter out obviously low-quality sources. From  
 1510 inspection of the percentage of entities verified,  
 1511 summaries, court documents on the CRLC website,  
 1512 and source text in Multi-LexSum for a sample of  
 1513 cases, we removed cases where less than 75% of  
 1514 entities could be verified.

1515 We found one legitimate case where summaries  
 1516 contained non-extractive entities: where the final  
 1517 sentence of the summary indicated whether the  
 1518 case was closed ‘as of’ the date of writing. In  
 1519 such cases, the date of writing was evidently not  
 1520 contained in the source documents. Therefore, if  
 1521 the last sentence of summaries in the training set  
 1522 contained ‘as of’, we removed this sentence so as  
 1523 not to encourage hallucination.

### 1524 C OREO: Further Details

1525 Formally, the OREO algorithm defines the  
 1526 summary-worthiness of a sentence  $x_i$  as the ex-  
 1527 pectation of its associated oracle evaluation:

$$\ell'_i := \sum_{Y^*} \mathcal{R}(Y^*, S) p(x_i|Y^*, D) p(Y^*|D, S) =$$

$$Y^* \sim \underbrace{E}_{p(Y^*|D,S)} \left[ \underbrace{\mathcal{R}(Y^*, S)}_{\text{oracle quality}} \quad \underbrace{p(x_i|Y^*, D)}_{\text{oracle membership}} \right]$$

1528 where  $\mathcal{R}$  denotes the mean of ROUGE-1 and  
 1529 ROUGE-2,  $D = \{x_i\}_1^m$  denotes the source text,  
 1530  $S$  is the reference (abstractive) summary, and  $Y^*$   
 1531 is the oracle summary space. The ‘oracle mem-  
 1532 bership’ term refers to if the oracle hypothesis  $Y^*$   
 1533 is in the oracle distribution, which is a uniform  
 1534 distribution over the  $t$  top results of the  $k$  oracle  
 1535 summary hypotheses returned by beam search. The  
 1536 final sentence labels are given by the scaled expect-  
 1537 ation  $\ell(x_i) = (\ell'_i - \bar{\ell}_{min}) / (\bar{\ell}_{max} - \bar{\ell}_{min})$  (Xu and  
 1538 Lapata, 2022). 1539

1540 To obtain the OREO labels for Multi-LexSum,  
 1541 we set the beam size hyperparameter  $k$  to 16, and  
 1542 the oracle distribution hyperparameter  $t$  to 16, as in  
 1543 the hyperparameter search performed in Xu and La-  
 1544 pata (2022), these were the best parameters for the  
 1545 most highly compressive dataset evaluated, Multi-  
 1546 News. We set the summary size hyperparameter  
 1547 to 30 (approx. 1024 / 34) sentences, based on the  
 1548 mean number of tokens (34, very long tail distribu-  
 1549 tion) per source sentence. However, after running  
 1550 OREO, in many cases fewer than 30 sentences were  
 1551 extracted (received a non-zero score) for a given  
 1552 case.

### 1553 D BERT Sentence Salience Classifier: 1554 Further Details

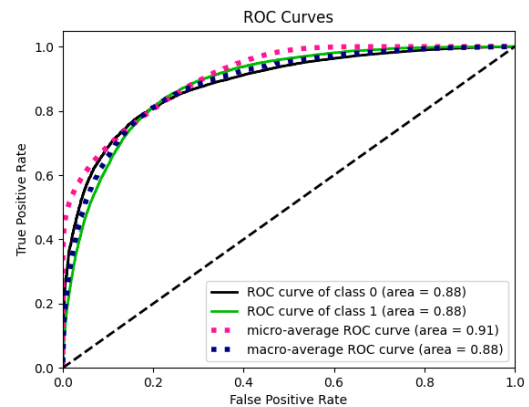


Figure 4: ROC curve for CaseLawBERT classifier.

1555 We train the CaseLawBERT model using its Py-  
 1556 torch implementation in Huggingface (Wolf et al.,

2020) on a single NVIDIA RTX A6000 GPU. The model was trained for 3 epochs (Zheng et al., 2021), with a batch size of 16 and using BertAdam with a learning rate of  $2e-5$  and warmup of 0.01. Inputs were truncated at 128 tokens for feasibility reasons due to the huge number of sentences in the test set; we acknowledge that not truncating may lead to improved results. As output, we obtained the probability of the sentence containing salient information. As we are not working with a threshold (to construct the inputs to PEGASUS, we use a ranked list by probability) and as metrics such as accuracy, precision, and recall are not very informative for highly skewed data, we report the classifier’s ROC-AUC score of 0.884 (Figure 5) - this indicates excellent (Mandrekar, 2010) performance, despite the computational considerations made.

### E Details of Results: Preliminary Content Selection Experiment

Table 5 details the mean number of tokens extracted per input strategy. Figure 5 details the distributions of the ROUGE recall score per input strategy.

	OREO	BERT
Sentences	264.15	1000.78
Windows	821.31	966.10
Paragraphs	679.73	596.47

Table 5: Mean number of tokens extracted for BERT-based and OREO-based input strategies.

### F Experimental Setup - PEGASUS

All experiments were conducted on a single NVIDIA RTX A6000 GPU, using the PyTorch implementations of PEGASUS and Legal-PEGASUS available from the Huggingface (Wolf et al., 2020).

### G Annotated Model Outputs

We include (Figure 6) examples of representative model outputs for two legal cases, compared with the results of our PEGASUS reproduction baseline and the gold standard summaries. Facts inconsistent with the case documents and other errors (such as assimilating information from cited cases) are highlighted in red.

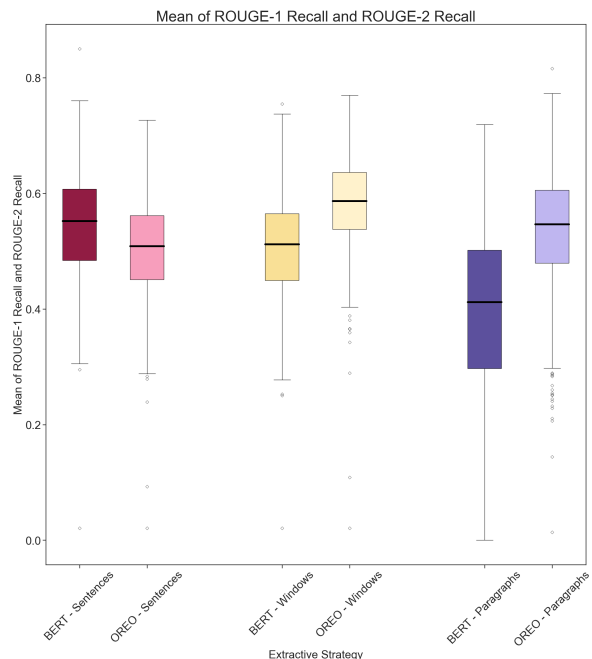


Figure 5: Distributions of ROUGE recall scores against corresponding reference summary for BERT-based and OREO-based strategies, demonstrating the difference in salient information retrieval between BERT-based and OREO-based counterparts.

	<b>Cerda v. Restaurant Associates</b>	<b>City of Los Angeles v. JPMorgan Chase</b>
Gold	On August 9, 2004, the plaintiff filed this class action lawsuit in the Eastern District of New York on behalf of all Hispanic and/or non-white refreshment concession personnel who were hired, employed or offered employment by the defendants, Restaurant Associates and RA Tennis Corp. The plaintiff alleged that the defendants engaged in unlawful employment discrimination on the basis of national origin and race respecting the assignment of positions and locations to food vendors employed at the U.S. Open Tennis Tournament. Specifically, they alleged that the defendants were deliberately and overwhelmingly assigning beer wagons to non-Hispanic white people over other demographics. The parties soon reached a settlement which was approved by U.S. District Court Judge Robert M. Levy, and Restaurant Associates agreed to pay damages. This case is closed.	On May 30, 2014, the City of Los Angeles filed a lawsuit in the U.S. District Court for the Central District of California against JP Morgan Chase Bank, under Equal Credit Opportunity Act (ECOA) and Fair Housing Act/Fair Housing Amendments Act (FHAA). The plaintiff sought damages, declaratory and injunctive relief, alleging that the defendants lending practices were predatory and intentionally discriminated against minority borrowers, and had a disparate impact on minorities. On August 6, 2015, the plaintiff filed an ex parte application to stay the case pending appeal of related action. On September 2, 2015, the case was dismissed as stipulated.
Reproduction Baseline	On August 9, 2004, a restaurant worker filed a class action lawsuit in the U.S. District Court for the Eastern District of New York against the <b>owners and managers of a restaurant in Brooklyn, New York</b> . The plaintiff alleged that the defendants had violated the <b>Fair Employment and Housing Act (FEH) and the Rehabilitation Act of 1973</b> . The court granted the plaintiff's motion for preliminary injunction, and the parties agreed to a settlement. The case is closed.	On May 30, 2014, the City of Los Angeles filed a lawsuit in the U.S. District Court for the Central District of California against JPMorgan Chase & Co. and Chase Manhattan Bank. The City alleged that the banks had <b>violated the Equal Credit Opportunity Act (ECA) and the Fair Credit Reporting Act (FCRA) by denying the City access to credit card information</b> . The court granted the City's <b>motion for summary judgment</b> and granted the defendants' motion to dismiss. <b>The parties agreed to a settlement agreement. The settlement agreement was entered into on September 2, 2015.</b>
BERT-Windows, LegalPEGASUS	On <b>July 20, 2009</b> , a class action lawsuit was filed in the U.S. District Court for the Eastern District of New York on behalf of Hispanic and non-white employees of Restaurant Associates, Inc. and RA Tennis Corp. who worked at the US Open Tennis Tournament in 2003. The complaint alleged that the defendants discriminated against non-Hispanic white employees and shunted them to lower-paying positions in violation of the federal, state and City employment non-discrimination statutes. The parties reached a settlement in which the defendants established a \$90,000 settlement fund to be divided pro rata among non-White employees and employees of non-U.S national origin who file claims. The settlement was approved by the Court on <b>July 20, 2010</b> .	On <b>September 2, 2015</b> , the City of Los Angeles filed this lawsuit in the U.S. District Court for the Central District of California against JPMorgan Chase Co., alleging that the bank violated the Fair Housing Act by targeting minority borrowers for predatory loan terms and increasing the costliness of non-conventional loans at the expense of minority borrowers. The case is one of four related cases brought by the City against large lending institutions, including Wells Fargo, Citigroup Inc., and the <b>California Department of Housing and Urban Development</b> . In 2012, the City entered into a global settlement agreement with the government pursuant to which JPMorgan agreed to pay \$753 million into a settlement fund for distribution 10 to qualified borrowers; and (b) provide an additional \$1.2 billion to foreclosure prevention actions. <b>The City's damages include lost tax revenues and the need to provide 21 increased municipal services.</b>

Restaurant Associates and RA Tennis Corp are food vendors at the U.S. open tennis tournament, not a physical restaurant, and owners and managers were not directly involved - this demonstrates a faithfulness problem not related to entities.

Throughout, dates included occur in the source text, but appear in the incorrect context in the generated summary.

Information following 'In 2012' is related to a cited case, not the current case, and contains artefacts

Figure 6: Annotated examples of representative model outputs for two cases, with facts inconsistent with the case documents and other errors (such as assimilating information from cited cases) highlighted in red.)