# Towards the Development of a LegalNLP Dataset for Bodo and the Evaluation of Abstractive Summarization Models

**Anonymous ACL submission**

## Abstract

Natural Language Processing (NLP) has become a transformative tool for analyzing large volumes of unstructured legal text, enabling tasks such as document summarization, judgment prediction, and legal information retrieval. However, most advancements in Legal NLP have been focused on high-resource languages like English, leaving low-resource languages such as Bodo significantly underrepresented. To address this gap, this paper presents the development of a legal training and test dataset for Bodo, a language spoken in Northeast India. Legal case documents and their summaries were sourced from publicly available platforms and translated into Bodo using the IndicTrans2 model, followed by preprocessing and standardization to ensure linguistic consistency and data quality using BLEU score and manual human evaluation. The dataset was also used to evaluate several state-of-the-art abstractive summarization models, including sequence-to-sequence architectures, pretrained transformers, and large language models, with performance assessed using ROUGE and CHRF scores. The findings emphasize the importance of building language-specific resources and provide a foundational benchmark for advancing Legal NLP research in Bodo and other low-resource languages.

## 1 Introduction

India has a complex judicial structure, characterized by a hierarchical system comprising the Supreme Court, High Courts, and District Courts. Among these, the Supreme Court holds the highest jurisdiction. The country follows a common law system, similar to that of the UK and the USA. Thousands of cases are registered within this system, and each case may take months or even years to reach a resolution. From the moment a case is filed to its final disposal, numerous documents and thousands of pages are generated and recorded. Legal practitioners are required to thoroughly read these documents to understand the case details. However, legal documents are often lengthy, complex, unstructured, and noisy making the task time consuming and tedious. Automatic summarization of lengthy legal documents has therefore become an essential tool within the judicial system (Bhattacharya et al., 2019). Despite its importance, legal document summarization remains a significant challenge due to the documents' complexity and volume, often necessitating the involvement of legal professionals for accurate interpretation.

India is characterized by significant linguistic diversity, with 22 officially recognized scheduled languages and numerous regional dialects (Mallikarjun, 2021). While English remains the primary language for legal proceedings at higher levels of the judiciary, regional courts frequently conduct their proceedings in the respective regional languages. Furthermore, the Government of India has initiated efforts to ensure that legal documentation and court processes are increasingly made available in regional languages to enhance accessibility and inclusivity. Given the diverse linguistic landscape, each regional language exhibits unique structural and characteristic features. Bodo is one such language, recognized as a scheduled language under the Indian Constitution. It is predominantly spoken in the northeastern region of India, particularly in the state of Assam, with an estimated 1.5 million (Narzary et al., 2021) speakers across the country. In the northeastern region of India, particularly in Assam, regional court proceedings are predominantly conducted in the Assamese language, which is the most widely spoken language in the state. Recently, Bodo has also been declared an official state language of Assam [1], a development that may influence the linguistic practices within judicial proceedings in the near future. However, there remains a limited availability of legal documents in regional
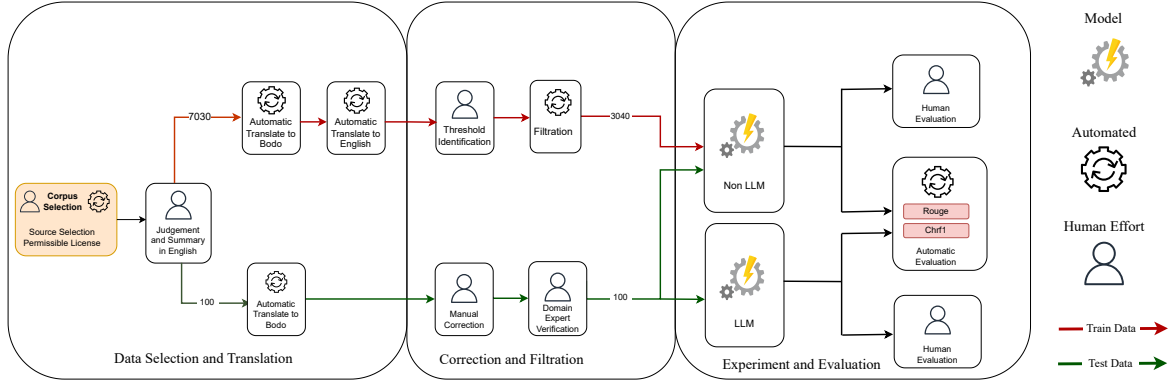
---

[1]https://prsindia.org

1

Figure 1: Proposed Architecture

languages such as Assamese and Bodo. Even when legal texts are translated into regional languages, their inherently lengthy, unstructured, and complex nature continues to pose significant challenges for legal practitioners, who must invest substantial time and effort to comprehend the material. This situation gives rise to two primary challenges: (1) The need for a comprehensive legal corpus in regional languages, and (2) The necessity of processing extensive legal documents into concise and accessible formats.

Our contributions are summarized as follows:

1. The development of a legal Bodo dataset derived from Indian legal case judgments and their summaries.

2. The evaluation of various abstractive summarization models on this dataset. To the best of our knowledge, this represents the first systematic effort to construct a Bodo legal dataset and to investigate the performance of different summarization techniques on legal texts in the Bodo language.

The organization of the paper is as follows: Section 2 reviews related work, Section 3 describes the methodology, and Section 4 presents model experiments and evaluation. Section 5 discusses the results and analysis, followed by the conclusion in Section 6 and future work in Section 7. The paper concludes with a discussion of limitations and ethical considerations.

## 2 Related Works

The summarizations are classified in two part: Extractive and Abstractive, In extractive summarization it selects and compiles direct sentences or phrases from the original text to form a summary. Whereas in abstractive summarization it rephrases and rewrites the content in a more concise form using natural language understanding.

The Supervised approaches (Liu and Chen, 2019) specifically tailored for the summarization of legal case documents. In the paper (Gong and Liu, 2001), the author employs matrix factorization to identify the most representative sentences. while (He et al., 2012) utilizes a data reconstruction framework for summarization, the graph-based LexRank algorithm (Erkan and Radev, 2004) both of which rely on sentence salience and connectivity. Models such as SummaRuNNer (Nallapati et al., 2017) have gained prominence in the extractive summarization domain. These approaches formulate summarization as a binary classification task. Similarly, in the paper (Zhong et al., 2019) the author introduces a template-based summarization framework that incorporates a two-stage classifier, the work (Polsley et al., 2016) enhances sentence ranking by combining TF-IDF weights with legal-domain-specific features, in the paper (Farzindar, 2004), the author employ term distribution-based models to rank sentences—utilizing TF-IDF and a k-mixture model.

A divide-and-conquer approach (Gidiotis and Tsoumakas, 2020) for long document summarization, wherein both the documents and their corresponding summaries are segmented based on sentence similarity. The paper (Bajaj et al., 2021) presents the development of a two-step extractive-abstractive framework for long document summarization, wherein salient sentences are first identified to create a compressed version of the document, which is then processed using a pre-trained

2

BART model to generate the final abstractive summary. The paper (Beltagy et al., 2020) presented transformer based architectures which have more efficient attention mechanisms, enabling these models to effectively handle and summarize long documents. BART (Lewis et al., 2019), a pre-trained using a denoising autoencoder objective, where the input text is corrupted through various noising strategies. The paper (Zhang et al., 2020), introduces a novel pretraining objective tailored to the summarization task. It removes important sentences (gap-sentences) from a document and trains the model to generate those sentences based on the remaining context, simulating the summarization process during pretraining. In the context of abstractive summarization for long documents, presenting deep communicating agents (Celikyilmaz et al., 2018) collaboratively process the input by dividing the text into manageable segments, with each agent responsible for encoding a specific subsection. These agents then exchange information to generate a coherent and concise summary of the entire document. RNN-based encoder-decoder models have shown strong performance in abstractive summarization tasks involving short input and output sequences. A neural network based model (Paulus et al., 2017) that incorporates a novel intra-attention mechanism, which separately attends to both the input and the previously generated output. Proposed a new training approach that combines traditional supervised learning for word prediction with reinforcement learning (RL) to optimize summary quality. The paper (Narayan et al., 2018) propose a novel abstractive model which is conditioned on the article's topics and based entirely on convolutional neural networks.

## 3 Methodology

In this section, we present a comprehensive overview of the data collection, translation, and evaluation processes. As shown in the Figure: 1, the details overflow of work resulting the creating of: 1. Corpus Selection and Automatic Translation, 2. Quality filteration, 3. Manual Correction, 4. Domain expert evaluation and 5. Inter Annotator Agreement.

### 3.1 Corpus Selection and Automatic Translation

In this section, we detail the process of constructing the Bodo legal judgement and its summary corpus.

Given the absence of an existing Bodo legal summarization dataset, we curated a total of 7,130 pairs of legal judgments and their corresponding summaries in English Collected Indian Supreme Court judgments and its summary dataset from a open source and publicly available platform [2], providing access to Indian legal databases. The database is having two sets of documents 7,030 set for train and 100 set for test set. Each document is having multiple pages with an average of 10 pages. The translation of dataset to low resource such as Bodo is done using existing open source machine translation IndicTrans2 (Gala et al., 2023) model which provides robust machine translation capabilities, especially for Indian languages. The IndicTrans2 platform allow us to translate from English to Bodo and vice versa. We have used IndicTrans2 having eng-indic-1B parameters from Huggingface Library [3]. The translation is performed in a NVIDIA TESLA 16GB machine with chuck size of 1000 in each translation. To translate all 7,130 documents from English to Bodo it took approximately 3 days. It may not always produce accurate or coherent outputs particularly when dealing with complex legal texts. Such translations can often result in lengthy, unstructured, or noisy data.

### 3.2 Quality filteration

Due to the large size of the training translated document, manually verifying each file was impractical. Instead, we applied a quality filtration technique using a Self-BLEU (Zhu et al., 2018) score threshold on a train set of 7,030 legal judgment and summary pairs. The original English data was first translated to Bodo using the IndicTrans2 indic-en-1B model from Hugging Face library [4], and then translated back to English using the same model. An average Self-BLEU score of 30.06 was observed between the original and back translated English texts. After filtering out files with below average BLUE scores, we retained 3,819 Bodo judgment summary pairs, representing 45% of the original training set.

### 3.3 Manual Correction

This section outlines the process of creating a high quality test set comprising 100 legal judgment and summary pairs. Each pair consists of a full-length legal judgment and its corresponding summary. To support the development and evaluation of Bodo

---

[2]https://zenodo.org/records/7152317#.Yz6mJ9JByC0
[3]ai4bharat/indictrans2-en-indic-1B
[4]ai4bharat/indictrans2-indic-en-1B

3

| Dataset | Type | Total Set | Total Words | Unique Words | Average Seq. Length |
|---------|------|-----------|-------------|--------------|---------------------|
| Train | Judgement | 7,030 | 12,261,564 | 287,133 | 3210.67 |
| | Summary | 7,030 | 2,234,167 | 102,393 | 585.01 |
| Test | Judgement | 100 | 3,43,463 | 20,987 | 3434.63 |
| | Summary | 100 | 65,532 | 8,456 | 655.32 |

Table 1: Dataset statistics for newly created dataset for Bodo.

language models, the English version of this test set has been translated into Bodo using the Indic-Trans2 eng-indic-1B translation model available on Huggingface. To ensure the quality and linguistic integrity of the dataset and hence the dataset is also limited only 100 set therefore the machine-translated Bodo texts were carefully reviewed and manually corrected by both language experts and legal domain specialists. This approach helps maintain the semantic fidelity of the original judgments and summaries while ensuring that the Bodo translations are fluent, contextually appropriate, and legally accurate.

### 3.4 Domain expert evaluation

This section describe expert annotation. Before sharing it to domain expert the translated text are manually check and corrected by the language expert. Afterward the corrected files are shared with 3 legal domain experts. The experts are from Legal background and are having knowledge of Bodo language. The experts were asked to read the translated text if it is generated as per the original english text and asked to rate each document file from 1 to 5 depending on the quality of the translated text. The rating details are mentioned in Table 2.

After annotation of all three experts we calculated the average Likert scores of each expert. As mentioned in Table 3 the average Likert scores given by three human annotators: A, B, and C. Annotator A assigned an average score of 4.10, while Annotator B gave the highest average score of 4.46. Annotator C provided the lowest average score of 3.97. These scores indicate that Annotator B rated the items most favorably, whereas Annotator C was comparatively more critical in their evaluations.

### 3.5 Inter Annotator Agreement

The Table 4 shows the inter-annotator agreement scores between three annotators: A, B, and C. The agreement score between Annotator A and Annotator B is 0.49, indicating a moderate level of consistency in their annotations. The agreement between Annotator A and Annotator C is 0.53, which is the highest among all pairs, suggesting that these

two annotators had the most consistent judgments. On the other hand, the agreement score between Annotator B and Annotator C is 0.34, the lowest among the three pairs, indicating relatively less consistency in their annotations. These scores help assess how reliably the annotators performed the task.

## 4 Experiments

We experimented with various pretrained and non-pretrained models to evaluate their effectiveness on the Bodo legal dataset. Pretrained models while powerful struggle with low resource languages, since most pretrained models are primarily trained on high-resource languages, they often underperform on underrepresented languages due to the lack of linguistic and contextual data. On the other hand, non-pretrained models, although more flexible in learning from scratch, require large amounts of high-quality annotated data to achieve competitive performance, something that is scarce in the low-resource Bodo language.

### 4.1 Models

Several models have been used to experiment with Bodo dataset. We trained models on the newly created dataset namely: sequence to sequence models (Seq2Seq), pre-trained models, and large language models.

**Sequence to Sequence Models (Seq2Seq)**: We used two models Pointer Generator and LSTM with encoding decoding. The pointer generator model (See et al., 2017) consists of two sections, the baseline sequence to sequence model and the pointer generator network. Generates words from a fixed vocabulary (like a standard Seq2Seq model) and copies words from the source text, useful for handling out-of-vocabulary (OOV) and rare words. The LSTM (Staudemeyer and Morris, 2019) encoder's role is to process the input sequence and summarise the information into a context vector (also known as a thought vector). The decoder uses the context vector produced by the encoder to generate the output sequence.

**Pretrained Models**: The three pretrained models

4

| Score | Rating | Description |
|---|---|---|
| 1 | Poor quality | meaningless translation or incomprehensible |
| 2 | Fair quality | many errors, hard to understand |
| 3 | Average quality | understandable but needs significant revision |
| 4 | Good quality | minor errors, but generally fluent and accurate |
| 5 | Excellent quality | accurate and fluent, almost no errors |

Table 2: Rating Score and Description

| Annotator | Average Likert Score |
|---|---|
| A | 4.10 |
| B | 4.46 |
| C | 3.97 |

Table 3: Human Annotator Scores

have been used in the experiment: Longformer Encoder Decoder, Pegasus and Bart model. The Longformer Encoder and Decoder (LED) (Beltagy et al., 2020) model has designed to handle long documents which includes attention mechanism. The attention pattern scales linearly with the input sequence, making it efficient for longer sequences. The Transformer encoder-decoder architecture in PEGASUS is pre-trained using the Gap Sentence Generation (GSG) objective (Zhang et al., 2020), where key sentences are masked and then predicted, effectively simulating the process of abstractive summarization. BART (Lewis et al., 2019) is a multilingual, pre-trained sequence-to-sequence model designed for natural language generation tasks involving for low-resource languages. With a compact design, BART contains 244 million parameters, making it suitable for low-resource environments.

**Large Language Models**: We generally used three LLM models in our experiments Gemma, Deepseek and Llamma. Google Gemma (Team et al., 2024) is a family of lightweight, open-source language models. The Gemma family includes two main model sizes Gemma 2B and Gemma 7B each available in both pretrained and instruction-tuned versions. In our experiments use used 2B size Gemma. DeepSeek (Bi et al., 2024) a family of open-source large language models (LLMs) the goal of advancing high quality, bilingual language models that are both powerful and accessible. The core of the DeepSeek model family is based on a transformer decoder-only architecture, following the structure of GPT-like models most notably DeepSeek 7B and DeepSeek 67B. The 7B variant is used in our experiments. LLaMA (Roumeliotis et al., 2023) is a series of open-source large language models designed to provide high-

performance natural language processing capabilities in a compact, accessible format. LLaMA 2 models were available in three sizes: LLaMA 2-7B, LLaMA 2-13B, and LLaMA 2-70B. Out of those LLaMA 2-7B used in our experiment. This is to be noted that, all these three LLM is experimented in prompting techniques with 100 test set data.

## 4.2 Evaluation Metrics

**ROUGE**: ROUGE (Recall-Oriented Understudy for Gisting Evaluation) (Lin, 2004) measures the overlap between a generated text and a reference text, primarily focusing on recall. It evaluates how much of the reference content is captured in the generated output. The most commonly used variants include ROUGE-1 (unigram overlap), ROUGE-2 (bigram overlap), and ROUGE-L (longest common subsequence). ROUGE is particularly popular in summarization tasks, where capturing key information from the original text is essential.

**ChrF-1**: ChrF (Character n-gram F-score) (Popović, 2015) evaluates the similarity between generated and reference texts based on character-level n-grams. It calculates a balanced F1-score considering both precision and recall over sequences of characters rather than words. ChrF is widely used in machine translation tasks where such nuances are critical, and it is particularly effective when traditional word-level tokenization is unreliable or inconsistent.

## 5 Results and Analysis

### 5.1 Automatic evaluation

The Table 5 presents a comparative evaluation of various algorithms used for a text summarization task using standard evaluation matrics such as Rouge and Chrf1 score. Among the traditional models, the Pointer Generator network achieved the highest performance, with an R-1 of 0.230, R-2 of 0.038, R-L of 0.092 and ChrF score of 42.136, indicating its strong ability to produce accurate and fluent summaries. In contrast, the LSTM with Encoder-Decoder and Legal LED models per-

5

| Annotator | Score |
|---|---|
| Between Annotator A and B | 0.49 |
| Between Annotator A and C | 0.53 |
| Between Annotator B and C | 0.34 |

Table 4: Inter Annotator Agreement Score

| Models | Algorithms | R-1 | R-2 | R-L | ChrF1 Score | BertScore |
|---|---|---|---|---|---|---|
| Seq2Seq | LSTM with Encoder and Decoder | 0.10 | 0.01 | 0.10 | 5.49 | 0.87 |
| | Pointer Generator | 0.23 | 0.03 | 0.09 | 42.13 | 0.90 |
| | Legal LED | 0.03 | 0 | 0.03 | 10.45 | 0.91 |
| Pretrained | BART | 0.23 | 0.03 | 0.16 | 28.65 | 0.90 |

Table 5: Performance comparison of summarization algorithms using ROUGE, ChrF and Bart Scores.

| 1 shot | | | | |
|---|---|---|---|---|
| Models | Rouge-1 | Rouge-2 | Rouge-L | ChrF1 score |
| Gemma-2b-it | 0.50 | 0.36 | 0.50 | 0.38 |
| Gemma-2b | 0.34 | 0.30 | 0.34 | 23.77 |
| Deepseek-llm-7b-chat | 0.23 | 0.08 | 0.16 | 34.37 |
| Llama-2-7b-chat-hf | 0.02 | 0 | 0.02 | 2.34 |
| Llama-3-8B-Instruct | 0.36 | 0.30 | 0.36 | 19.94 |

Table 6: Evaluation metrics of different LLM-based models on 1-shot

| 100 shot | | | | |
|---|---|---|---|---|
| Models | Rouge-1 | Rouge-2 | Rouge-L | ChrF1 score |
| Gemma-2b-it | 0.06 | 0.04 | 0.06 | 3.80 |
| Gemma-2b | 0.10 | 0.06 | 0.10 | 4.75 |
| Deepseek-llm-7b-chat | 0.23 | 0.08 | 0.16 | 34.37 |
| Llama-2-7b-chat-hf | 0.06 | 0 | 0.05 | 6.65 |
| Llama-3-8B-Instruct | 0.11 | 0.05 | 0.10 | 6.69 |

Table 7: Evaluation metrics of different LLM-based models on 100-shot

formed poorly across all metrics, suggesting limitations in their ability to handle the summarization task effectively. LED considered to be maximum BertScore of 0.91 indicating strong alignment between its generated summaries and the reference texts and LSTM having lowest BertScore of 0.87 suggesting comparatively less semantic closeness, Additionally, PEGASUS achieved a score of zero across all evaluation metrics, highlighting its inability to effectively perform the summarization task in Bodo languages. Consequently, its results were excluded from the table for clarity and relevance. BART demonstrated a competitive performance, particularly with a strong R-L score of 0.168 and a moderate ChrF score of 28.65 highlighting its potential for summarization in multilingual or low-resource languages.

Large Language Models (LLMs), including Gemma-2b, Deepseek-llm-7b-chat, Llama-2-7b-chat-hf, and Meta-Llama-3-8B-Instruct, were evaluated using 1-shot prompting techniques as mentioned in Table 6. Among the models, Gemma-2b-it achieved the highest performance across all ROUGE metrics, with a R-1 score of 0.50, R-2 of 0.36, and R-L of 0.50. It also recorded the highest ChrF1 score of 0.38, Gemma-2b also performed reasonably well, with R-1, R-2, and R-L scores of 0.34, 0.30, and 0.34 respectively, and a ChrF1 score of 23.77. Llama-3-8B-Instruct showed similar ROUGE performance (0.36 for R-1 and R-L, and 0.30 for R-2), but had a lower ChrF1 score of 19.94. Deepseek-llm-7b-chat had lower ROUGE

449
450
451
452
453
454
455
456
457
458
459
460
461
462
463
464
465
466
467
468
469
470
471
472
473
474
475
476
477
478
479
480
481
482
483
484
485
486
487
488
489
490
491
492
493
494
495
496
497

scores (R-1: 0.23, R-2: 0.08, R-L: 0.16) but a notably high ChrF1 score of 34.37, In contrast, Llama-2-7b-chat-hf performed poorly across all metrics, with negligible ROUGE scores and a ChrF1 score of only 2.34, indicating limited effectiveness in generating relevant summaries in a 1-shot setting.

Table 7 shows experiments with 100-shot prompting setup. Among these, Deepseek-llm-7b-chat performed particularly well, achieving ROUGE scores on par with the Pointer Generator and a high ChrF score of 34.372, showcasing its adaptability to summarization tasks even with limited examples. Other LLMs such as Gemma variants and Llama-based models yielded lower scores, with ChrF values ranging from 3.805 to 6.699, indicating comparatively weaker summarization performance under the same 100-shot setting. The results also note the absence of output for Pegasus, likely due to evaluation constraints.

Overall, the findings underscore that while traditional models like the Pointer Generator remain strong baselines, newer LLMs especially when guided through few-shot prompting demonstrate promising capabilities, with potential for further enhancement through fine-tuning and domain adaptation.

### 5.2 Human evaluation

The Table A.1 presents the results of a human evaluation conducted on three summarization models: BART, Pointer Generator, and LED. Two annotators independently rated the output of each model using a Likert scale, where higher scores indicate better performance in terms of summary quality. Among the models, BART received the highest scores, with Annotator A assigning a score of 2.91 and Annotator B giving 3.08. This suggests that BART generated more coherent and relevant summaries compared to the other models. The Pointer Generator model received moderate scores—2.45 from Annotator A and 2.15 from Annotator 2—indicating average performance. In contrast, the LED model was rated the lowest, with Annotator A giving a score of 1.93 and Annotator B assigning 1.52. These results reflect that LED's summaries were perceived as less effective. Overall, BART outperformed the other models in human judgment, while the variation in scores between the two annotators was relatively small, indicating consistency in evaluation.

## 6 Conclusion & Future Works

The study presents the first comprehensive benchmark dataset for legal document summarization in the Bodo language, a low resource language. We have collected and translated over 7,000 legal judgment-summary pairs from Indian Supreme Court cases using IndicTrans2 and conducting rigorous quality filtration and human annotation. Our experiments with both traditional Seq2Seq and state-of-the-art pretrained and large language models (LLMs) reveal that while classical models like the Pointer Generator perform strongly, LLMs such as DeepSeek-7B also show promising results even in a few-shot setting.

Future research can be focus on several areas-First, Fine-tuning domain-specific model on the Bodo legal dataset may yield significant improvements in performance, particularly by capturing the nuances of legal language in a low-resource setting. Second, The inclusion of more annotated data and incorporating multi-domain legal texts from district and high courts could enrich the diversity and applicability of the dataset. Third, Incorporating syntactic and semantic features specific to the Bodo language could lead to more linguistically informed summarization, potentially improving fluency and coherence and Finally, Integrating this summarization framework into practical tools for legal professionals and citizens could improve legal accessibility, promote transparency and bridge the gap between complex legal language and public understanding.

## Limitations

The translation of complex legal texts into Bodo using machine translation tools like IndicTrans2, while efficient, may introduce semantic inconsistencies or syntactic inaccuracies. Manual correction and expert verification were applied only to the test set, whereas the training data was filtered using BLEU-based back-translation, which may not capture nuanced translation errors. Furthermore, due to resource constraints, only 100 samples were manually evaluated, which may limit the robustness of evaluation insights.

## Ethics

The newly created dataset in the Bodo language, which includes translations of legal documents and their summaries, follows same licensing and permission terms as the original source data. Since

498
499
500
501
502
503
504
505
506
507
508
509
510
511
512
513
514
515
516
517
518
519
520
521
522
523
524
525
526
527
528
529
530
531
532
533
534
535
536
537
538
539
540
541
542
543
544
545
546

| Model | Annotator | Average Likert Score |
|-------|-----------|---------------------|
| BART | A | 2.91 |
| | B | 3.08 |
| Pointer Generator | A | 2.45 |
| | B | 2.15 |
| LED | A | 1.93 |
| | B | 1.52 |

Table 8: Results of Human Evaluation

the Bodo dataset is a derivative work based on the publicly available legal dataset from Zenodo [5], under the Creative Commons Attribution 4.0 International license. Permission to use, distribute, and reproduce the translated data aligns with the terms under which the source data was released.

The use of this data aligns with its intended purpose, as the resulting dataset is exclusively developed for research in the domain of Legal NLP, particularly in low-resource languages like Bodo. The translated and annotated Bodo Legal Summarization Dataset is intended solely for academic and research applications, including model development, benchmarking, and linguistic analysis. All derivatives of the original data, including translations and annotations, adhere to the access conditions of the source to ensure compliance with ethical and legal standards.

# References

Ahsaas Bajaj, Pavitra Dangati, Kalpesh Krishna, Pradhiksha Ashok Kumar, Rheeya Uppaal, Bradford Windsor, Eliot Brenner, Dominic Dotterrer, Rajarshi Das, and Andrew McCallum. 2021. Long document summarization in a low resource setting using pretrained language models. *arXiv preprint arXiv:2103.00751*.

Iz Beltagy, Matthew E Peters, and Arman Cohan. 2020. Longformer: The long-document transformer. *arXiv preprint arXiv:2004.05150*.

Paheli Bhattacharya, Kaustubh Hiware, Subham Rajgaria, Nilay Pochhi, Kripabandhu Ghosh, and Saptarshi Ghosh. 2019. A comparative study of summarization algorithms applied to legal case judgments. In *Advances in Information Retrieval: 41st European Conference on IR Research, ECIR 2019, Cologne, Germany, April 14–18, 2019, Proceedings, Part I 41*, pages 413–428. Springer.

Xiao Bi, Deli Chen, Guanting Chen, Shanhuang Chen, Damai Dai, Chengqi Deng, Honghui Ding, Kai Dong, Qiushi Du, Zhe Fu, and 1 others. 2024. Deepseek llm: Scaling open-source language models with longtermism. *arXiv preprint arXiv:2401.02954*.

Asli Celikyilmaz, Antoine Bosselut, Xiaodong He, and Yejin Choi. 2018. Deep communicating agents for abstractive summarization. *arXiv preprint arXiv:1803.10357*.

Günes Erkan and Dragomir R Radev. 2004. Lexrank: Graph-based lexical centrality as salience in text summarization. *Journal of artificial intelligence research*, 22:457–479.

Atefeh Farzindar. 2004. Atefeh farzindar and guy lapalme,'letsum, an automatic legal text summarizing system'in t. gordon (ed.), legal knowledge and information systems. jurix 2004: The seventeenth annual conference. amsterdam: Ios press, 2004, pp. 11-18. In *Legal knowledge and information systems: JURIX 2004, the seventeenth annual conference*, volume 120, page 11. IOS Press.

Jay Gala, Pranjal A Chitale, Raghavan AK, Varun Gumma, Sumanth Doddapaneni, Aswanth Kumar, Janki Nawale, Anupama Sujatha, Ratish Puduppully, Vivek Raghavan, and 1 others. 2023. Indictrans2: Towards high-quality and accessible machine translation models for all 22 scheduled indian languages. *arXiv preprint arXiv:2305.16307*.

Alexios Gidiotis and Grigorios Tsoumakas. 2020. A divide-and-conquer approach to the summarization of academic articles. *arXiv preprint arXiv:2004.06190*.

Yihong Gong and Xin Liu. 2001. Generic text summarization using relevance measure and latent semantic analysis. In *Proceedings of the 24th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 19–25.

Zhanying He, Chun Chen, Jiajun Bu, Can Wang, Lijun Zhang, Deng Cai, and Xiaofei He. 2012. Document summarization based on data reconstruction. In *Proceedings of the AAAI Conference on Artificial Intelligence*, volume 26, pages 620–626.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Ves Stoyanov, and Luke Zettlemoyer. 2019. Bart: Denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *arXiv preprint arXiv:1910.13461*.

---

[5]https://zenodo.org/records/7152317#.Yz6mJ9JByC0

Chin-Yew Lin. 2004. Rouge: A package for automatic evaluation of summaries. In *Text summarization branches out*, pages 74–81.

Chao-Lin Liu and Kuan-Chun Chen. 2019. Extracting the gist of chinese judgments of the supreme court. In *proceedings of the seventeenth international conference on artificial intelligence and law*, pages 73–82.

B Mallikarjun. 2021. The eighth schedule languages* a critical appraisal. *Language in India*, 21(1).

Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. Summarunner: A recurrent neural network based sequence model for extractive summarization of documents. In *Proceedings of the AAAI conference on artificial intelligence*, volume 31.

Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Don't give me the details, just the summary! topic-aware convolutional neural networks for extreme summarization. *arXiv preprint arXiv:1808.08745*.

Mwnthai Narzary, Gwmsrang Muchahary, Maharaj Brahma, Sanjib Narzary, Pranav Kumar Singh, and Apurbalal Senapati. 2021. Bodo resources for nlp-an overview of existing primary resources for bodo. *AIJR Proceedings*, pages 96–101.

Romain Paulus, Caiming Xiong, and Richard Socher. 2017. A deep reinforced model for abstractive summarization. *arXiv preprint arXiv:1705.04304*.

Seth Polsley, Pooja Jhunjhunwala, and Ruihong Huang. 2016. Casesummarizer: a system for automated summarization of legal texts. In *Proceedings of COLING 2016, the 26th international conference on Computational Linguistics: System Demonstrations*, pages 258–262.

Maja Popović. 2015. chrf: character n-gram f-score for automatic mt evaluation. In *Proceedings of the tenth workshop on statistical machine translation*, pages 392–395.

Konstantinos I Roumeliotis, Nikolaos D Tselikas, and Dimitrios K Nasiopoulos. 2023. Llama 2: Early adopters' utilization of meta's new open-source pretrained model.

Abigail See, Peter J Liu, and Christopher D Manning. 2017. Get to the point: Summarization with pointer-generator networks. *arXiv preprint arXiv:1704.04368*.

Ralf C Staudemeyer and Eric Rothstein Morris. 2019. Understanding lstm–a tutorial into long short-term memory recurrent neural networks. *arXiv preprint arXiv:1909.09586*.

Gemma Team, Thomas Mesnard, Cassidy Hardin, Robert Dadashi, Surya Bhupatiraju, Shreya Pathak, Laurent Sifre, Morgane Rivière, Mihir Sanjay Kale, Juliette Love, and 1 others. 2024. Gemma: Open models based on gemini research and technology. *arXiv preprint arXiv:2403.08295*.

Jingqing Zhang, Yao Zhao, Mohammad Saleh, and Peter Liu. 2020. Pegasus: Pre-training with extracted gap-sentences for abstractive summarization. In *International conference on machine learning*, pages 11328–11339. PMLR.

Linwu Zhong, Ziyi Zhong, Zinian Zhao, Siyuan Wang, Kevin D Ashley, and Matthias Grabmair. 2019. Automatic summarization of legal decisions using iterative masking of predictive sentences. In *Proceedings of the seventeenth international conference on artificial intelligence and law*, pages 163–172.

Yaoming Zhu, Sidi Lu, Lei Zheng, Jiaxian Guo, Weinan Zhang, Jun Wang, and Yong Yu. 2018. Texygen: A benchmarking platform for text generation models. In *The 41st international ACM SIGIR conference on research & development in information retrieval*, pages 1097–1100.

# A  Appendix

## A.1  Expert Annotator Details

All the three annotators, labeled A, B, and C, participated in the annotation. All of them are native Bodo speakers and hold an LLB (Bachelor of Laws) degree. Total 5 annotators we called for an interview for this work, Based on the qualification and experience we found 3 annotators suitable and selected for this work. Annotators A and B are residents of Bongaigaon, India, while Annotator C resides in Baksa, India. The ages of the annotators are 25 for A, 26 for B, and 25 for C. The annotators have mutually agreed to work with a honorarium of Rs.7/- per annotate.

| Annot. | NL | Qual. | Resident | Age |
|--------|------|-------|-------------------|-----|
| A | Bodo | LLB | Bongaigaon, India | 25 |
| B | Bodo | LLB | Bongaigaon, India | 26 |
| C | Bodo | LLB | Baksa, India | 25 |

Details of Human annotators. Here NL represents Native language. Qual. denotes Qualifications, and Annot. as annotators.

## A.2  Hyper parameters used in experiments

The hyper parameters used in the experimental setup of the models are presented in the Table 10

## A.3  Large Language Model Prompts used

The set of prompts used for evaluations on LLM models using Bodo judgement and summary test data are shown Table 11.

## A.4  Generated Results of LED, Pointer Generator and BART model

Here, Fig 2 indicate Bodo judgement and summary generated result of LED, Pointer Generator and BART models.

| Model | Parameters |
|---|---|
| LSTM with encoder and decoder | arch bilstm-lstm Vocab-size = 8200 max-tokens= 8000 max-source-positions=4096 max-target-positions=4096 beam=5 remove-bpe batch-size=32 skip-invalid-size-inputs-valid-test |
| Pointer Generator | Vocab size=10000, Position markers=1000, Epoch=50, Warmup updates=400, Learning Rate (Lr) = 0.0007, Max tokens=2048, Update freq=2, Pointer layer = -2 |
| LED | Learning Rate (lr): 1e-3 Weight Decay: 0.01 Number of Epochs: 4 Per-device Train Batch Size: 8 Warmup Steps: 500 Evaluation Strategy: "epoch" |
| BART | Number of traning epochs=8, warmp up step=500, per device train batch size=4, per gpu eval batch size =8, gradient accumulation steps =16, evauation strategy="epoch", weigth decay=0.01, logging steps=10, eval steps=500, fp16=True, save steps=100, |
| Gemma-2B | Max Tokens=150, Temperature=0.4, Top-p=0.9, Sampling=True |
| Deepseek-7B | MAX INPUT = 3500, MAX OUTPUT = 250, BATCH SIZE = 4, max length = 4096, truncation = True, temperature = 0.7, do sample = True, $max_newtokens = 250$ |
| Llama2 | MAX INPUT = 3500, MAX OUTPUT = 250, max length = 4096, truncation = True, do sample = True, max new tokens = 250, temperature = 0.7 |

Table 10: Hyper parameters used to experiments models

### A.5 Generated Results of Gemma, Llama and DeepSeek LLMs

Here, Fig 3 indicate Bodo judgement and summary generated result of Gemma, Llama and DeepSeek models.

| Model | Shot | Prompts |
|---|---|---|
| Gemma-2b-it | 1-shot | You are a Bodo language expert.Create a concise 3-4 line summary in Bodo language only. Now, based on the judgment above, write only the new summary in Bodo language: |
| | 100-shot | Create a concise 3-5 sentence summary in Bodo language using the example style. New Summary Requirements**: - Use simple Bodo language - Include key facts - 15-20 sentences maximum - Consistent style |
| Deepseek-llm-7b-chat | 1 shot | You are a Bodo language expert. Create a concise 3-4 line summary in Bodo language only. Now, based on the judgment above, write only the new summary in Bodo language: New Summary (Bodo only) |
| | 100-shot | You are a Bodo language expert. Create a concise 3-4 line summary in Bodo language only. Now, based on the judgment above, write only the new summary in Bodo language: New Summary (Bodo only) |
| Llama-2-7b-chat-hf | 1 Shot | Create a concise 3-5 sentence summary in Bodo language using the example style. New Summary Requirements: - Use simple Bodo language - Include key facts - 3-5 sentences maximum - Consistent style |
| | 100 Shot | Create a concise 3-5 sentence summary in Bodo language using the example style. New Summary Requirements: - Use simple Bodo language - Include key facts - 3-5 sentences maximum - Consistent style |
| Llama-3-8B-Instruct | 1 Shot | You are a helpful assistant that summarizes legal judgments. Generate a concise summary in *Bodo language only* that: 1. Captures key information. 2. Maintains original meaning. 3. Is under 100 words. |
| | 100 Shot | Create a concise 3-5 sentence summary in Bodo language using the example style. New Summary Requirements: - Use simple Bodo language - Include key facts - 3-5 sentences maximum - Consistent style |

Table 11: Large Language Model 1 and 100 Shot Prompts

(a)

Output Result of LED Model. Left column indicate original summary and right column indicate generate summary

(b)

Output Result of Pointer Generator Model. Left column indicate original summary and right column indicate generate summary

(c)

Output Result of BART Model. Left column indicate original summary and right column indicate generate summary

Figure 2: Bodo judgement and generated summary pair of LED, Pointer Generator and BART Models

1957 मायथाइनि दायगोनां आर'ज नं 46।

1953 मायथाइनि केस नं 176/एसआव बम्बेनि गाहाय प्रेसीडेन्सी मेजिस्ट्रेट, बम्बेनि 19 जुन, 1954 मायथाइनि बिजिरनाय आरो बिथोननिफ्राय सोमजिनाय 1954 मायथाइनि दायगोनां आर'ज नं 1256 आव बम्बे गोजौ बिजिरसालिनि 4 फेब्रुआरि, 1955 मायथाइनि बिजिरनाय आरो बिथोननिफ्राय जर'खा गनायथिजों आरज गाबनाय।

हरदयाल हार्डी, खावलायगिरिफोरनि थाखाय।

फिन्नाय होगिरिनि थाखाय एच.जे. उमरीगर आरो बि.एच. धेबर।

(a)
Generated summary of DeepSeek model

**Instruction**: Create a concise 3-5 sentence summary in Bodo language using the example style. **Example Judgement**: 1957 मायथाइनि दायगोनां आर'ज नं 46। 1953 मायथाइनि केस नं 176/एसआव बम्बेनि गाहाय प्रेसीडेन्सी मेजिस्ट्रेट, बम्बेनि 19 जुन, 1954 मायथाइनि बिजिरनाय आरो बिथोननिफ्राय सोमजिनाय 1954 मायथाइनि दायगोनां आर'ज नं 1256 आव बम्बे गोजौ बिजिरसालिनि 4 फेब्रुआरि, 1955 मायथाइनि बिजिरनाय आरो बिथोननिफ्राय जर'खा गनायथिजों आरज गाबनाय। हरदयाल हार्डी, खावलायगिरिफोरनि थाखाय। फिन्नाय होगिरिनि थाखाय एच.जे.

(b)
Generated summary of Llama2 model

Sure, here's the summary you requested: गोरबलायनायनि हाबाफारिया जानुवारियाव जागायजेन्दोमोन। 1952 मायथाइयाव खावलायग्रा 1 आरो बिनि खामानि मावग्राफोरनि गेजेराव माखासे नांलायनायफोरनि सोमोन्दै। 9 मे, 1952 मायथाइयाव, जुथाइ आरो 2 जुन, 1952 मायथाइयाव, खावलायग्रा 1 आ गोरोबथा मावख

(c)
Generated summary of Gemma-2b-it model

Figure 3: Bodo generated summary of DeepSeek, Llama-2 and Gemma-2b-it model