# IL-TUR: Benchmark for Indian Legal Text Understanding and Reasoning

**Abhinav Joshi**[¶][*]  **Shounak Paul**[◇][*]
**Akshat Sharma**[¶]  **Pawan Goyal**[◇]  **Saptarshi Ghosh**[◇]
**Ashutosh Modi**[¶][†]
[¶]IIT Kanpur, [◇]IIT Kharagpur

shounakpaul95@kgpian.iitkgp.ac.in, {pawang,saptarshi}@cse.iitkgp.ac.in,
{ajoshi, akshatsh, ashutoshm}@cse.iitk.ac.in

## Abstract

Legal systems worldwide are inundated with exponential growth in cases and documents. There is an imminent need to develop NLP and ML techniques for automatically processing and understanding legal documents to streamline the legal system. However, evaluating and comparing various NLP models designed specifically for the legal domain is challenging. This paper addresses this challenge by proposing **IL-TUR**: Benchmark for Indian Legal Text Understanding and Reasoning. **IL-TUR** contains monolingual (English, Hindi) and multi-lingual (9 Indian languages) domain-specific tasks that address different aspects of the legal system from the point of view of understanding and reasoning over Indian legal documents. We present baseline models (including LLM-based) for each task, outlining the gap between models and the ground truth. To foster further research in the legal domain, we create a leaderboard (available at: https://exploration-lab.github.io/IL-TUR/) where the research community can upload and compare legal text understanding systems.
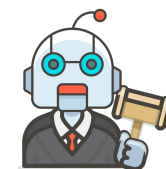
*"Justice delayed is justice denied"* - Legal Maxim

## 1 Introduction

Besides several other purposes, legal systems have been established in various countries to ensure, at the very minimum, order and fairness in society and to safeguard fundamental human rights. However, legal systems worldwide struggle with exponentially growing legal cases in various courts. It is even more pronounced in populous countries; e.g., in India, there are about 50 million pending cases in multiple courts at various levels (district, state, federal) (National Judicial Data Grid, 2023). Such a massive backlog of cases goes against the



Figure 1: **IL-TUR**: A consolidated benchmark covering a wide range of legal text understanding and reasoning tasks with a publically available leaderboard.

fundamental human right of fair access to justice. Documents in different natural languages are the backbone of various legal processes. Natural Language Processing (NLP) based techniques could be helpful in various legal processes involving fundamental tasks related to information extraction, document understanding, and prediction. This paper introduces **IL-TUR**, a benchmark for *Indian Legal Text Understanding and Reasoning*. The purpose of **IL-TUR** is twofold. First, it aims to foster research in the Legal-NLP (L-NLP) domain and plans to address the pain points associated with processing legal texts (see below); second, it provides a platform for comparing different models and further advancing the L-NLP domain.

**Why a separate benchmark for the legal domain?** The legal text involves natural language but differs from the regular text used to train NLP models. 1) Many of the terms used in legal documents are domain-specific. For example, some words used in everyday language have specialized meanings in legal parlance. The presence of a different lexicon posits a need for specialized NLP tools to handle legal texts. 2) Legal documents are typically very long compared to regular texts. For example, the average length of a legal document from the Supreme Court of India (SCI) is 4000

---

[*] Equal Contributions
[†] Corresponding Author

11460

words (Malik et al., 2021). It poses a challenge for existing NLP models (e.g., LLMs) as the information is spread throughout the document and must be linked together for reasoning. Moreover, many of the existing language models (e.g., BERT (Devlin et al., 2019)) have limitations on the length (512 tokens) of the input. It requires developing specialized models for processing and handling long legal documents. 3) Legal documents are highly unstructured and sometimes noisy (for example, in the Indian setting, most documents are typed manually in the courts and prone to grammatical mistakes and typos). The absence of structure in the documents makes extracting semantically relevant information from large chunks of text difficult. 4) The legal domain is further subdivided into specialized subdomains; for example, criminal law differs from civil law, and both differ from banking and insurance law. Even though some fundamental legal principles are shared across various laws, models trained on a particular law (e.g., civil law) may not work on another (e.g., banking and insurance law). Hence, domain adaptation is a challenge. 5) Lastly, many existing state-of-the-art (SOTA) NLP models are black boxes; however, explainability is not a second-class citizen for the legal domain. For models to be widely usable by legal practitioners, these need to be explainable. Due to the above reasons, a separate set of models/systems is required to process and understand legal documents. Given the huge backlog of cases, NLP-based technologies could come to our rescue and help streamline the legal workflow. Even a small technical intervention can have a considerable impact. Hence, a benchmark is needed to promote the development of models in this area. In a nutshell, we make the following contributions:

- We introduce **IL-TUR**: a benchmark for Indian Legal Text Understanding and Reasoning. The benchmark has eight tasks (in English and 9 Indian languages) requiring different types of legal knowledge and skills to solve. Moreover, the list of tasks is not exhaustive, and we plan to keep adding more tasks to **IL-TUR**. Currently, there are various L-NLP-specific tasks; however, these occur in isolation, making it difficult to keep track of progress made in the field. Similar to existing NLP benchmarks (e.g., GLUE (Wang et al., 2018a)), we consolidate and harmonize some of the existing L-NLP tasks and create new tasks resulting in a unified benchmark.

- We report baseline model results on each of the tasks. We also experiment with various LLMs (§4), and results show that LLMs are far from solving the tasks and hence point towards the need to develop better models.
- We release the dataset and baseline models associated with each task. Further, we create a leaderboard where anyone can upload their model and test against the baselines and other proposed systems (e.g., Fig. 1). The datasets, models, and the leaderboard are available via the following website: `https://exploration-lab.github.io/IL-TUR/`.

## 2 Related Work

Over the past few years, L-NLP has been a fertile area for research. Researchers have explored different aspects of the legal domain via various tasks such as Prior Case Retrieval (Joshi et al., 2023; Jackson et al., 2003a), Case Prediction (Malik et al., 2021; Chalkidis et al., 2019; Strickson and De La Iglesia, 2020; Kapoor et al., 2022), Summarization (Moens et al., 1999), Semantic Segmentation of Legal Documents (Malik et al., 2022; Kalamkar et al., 2022b; Bhattacharya et al., 2019), and Information Extraction and Retrieval (Tran et al., 2019; Lagos et al., 2010). On the modeling side, various techniques have been proposed, ranging from classical ML-based methods such as SVM (Al-Kofahi et al., 2001; Jackson et al., 2003b) to recent transformer-based models (Chalkidis et al., 2019; Malik et al., 2021). Researchers have also proposed legal domain-specific language models such as LegalBERT (Chalkidis et al., 2020), CaseLawBERT (Zheng et al., 2021) and InLegalBERT and InCaseLawBERT (Paul et al., 2023). However, legal LLMs have shown limited success and have not demonstrated generalization and transfer learning capabilities (Chalkidis, 2023; Malik et al., 2021; Joshi et al., 2023).

**Comparison with Existing Benchmarks:** Benchmarks have played a crucial role in the development of better techniques and models in almost every domain, such as computer vision (Deng et al., 2009; Guo et al., 2014; Wu et al., 2013) and reinforcement learning (Laskin et al., 2021; Cobbe et al., 2020; Zhang et al., 2018). Similarly, in the NLP domain, various benchmarks have been proposed, for example, GLUE (Wang et al., 2018a), Super-GLUE (Wang et al., 2019a), XTREME (Hu et al., 2020), CLUE (Xu et al., 2020), GLGE (Liu

et al., 2020), and IndicNLPSuite (Kakwani et al., 2020). However, these benchmarks focus on the general NLP domain, and models developed for the generic domains do not perform well for the legal domain (Malik et al., 2022; Joshi et al., 2023). Similar attempts have thus been made for the legal domain; for example, Chalkidis et al. (2022a) developed LexGLUE, a specialized English language benchmark (restricted to EU and US legal systems) for evaluating legal NLP models, by consolidating existing datasets for various tasks. LexGLUE introduces six main (all classification-based) tasks: violated article identification, case issue classification, concept identification, contract topic prediction, unfair contractual terms identification, and case holding identification. Niklaus et al. (2023) have proposed LEXTREME, a multi-lingual (24 EU languages) legal NLP benchmark (all tasks classification-based) restricted to EU and Brazilian jurisdictions. Chalkidis et al. (2022b) have introduced FAIRLEX, a multi-lingual benchmark consisting of cases from 5 languages and 4 jurisdictions, to test the fairness of different models on legal judgment and topic prediction. Hwang et al. (2022) have introduced LBOX benchmark for the Korean legal system. The benchmark targets tasks related to classification and summarization; the documents are in Korean. Recently, Guha et al. (2023) released LegalBench, a large, collaborative legal benchmark (restricted to US legal system) consisting of 162 tasks (in English) to test the reasoning abilities of LLMs. The tasks belong to six different categories of legal reasoning and address various stages in the pipeline of the litigation process. LegalBench is primarily focused on testing the ability of LLMs to handle legal processes at various stages of litigation; consequently, the tasks involve shorter texts (avg. length $\sim$ 200 words). To benchmark LLMs for Chinese law, Fei et al. (2023) released LawBENCH, a benchmark consisting of 20 tasks (in Chinese) to evaluate the capability of LLMs to memorize and understand legal knowledge. Most of these tasks consist of longer texts than LegalBench (avg. length $\sim$ 300 words).

**IL-TUR differs from the existing benchmarks (see Table 1)**. First, **IL-TUR** focuses on multiple tasks that are not restricted to classification but also involve information retrieval, generation, and explanation. Second, via **IL-TUR**, we introduce tasks that are grounded in the actual legal workflow and, consequently, are more complex and involve actual long legal documents (average length 4000

| Dataset | Jurisdictions | System | Task types | Languages |
|---|---|---|---|---|
| LexGLUE | U.S., E.U. | Predominantly Civil Law | Classification | English |
| LEXTREME | E.U., Brazil | Predominantly Civil Law | Classification | E.U. |
| FAIRLEX | E.U., U.S., China, Switzerland | Predominantly Civil Law | Fairness evaluation | E.U., Chinese |
| LBOX | Korea | Civil Law | Classification, Generation | Korean |
| LEGALBENCH | Multiple | Common & Civil Law | Generation | English |
| LAWBENCH | China | Civil Law | Classification, Generation, Extraction | Chinese |
| IL-TUR (ours) | India | Common & Civil Law | Classification, Retrieval, Generation, Extraction | English, Indian |

Table 1: Comparison of different L-NLP benchmarks.

words). In contrast to some of the popular benchmarks, **IL-TUR** is not introduced to test the law understanding capability of LLMs but rather to address the problems plaguing the judiciary. In the future, if LLMs are replaced by some other class of machine learning models, **IL-TUR** would still be relevant. In fact, as shown in our experiments, we observe that long legal documents are challenging for LLMs. Third, **IL-TUR** is based on Indian legal documents. Given that India is the most populous country in the world (population of $\sim$ 1.4 billion (United Nations, 2023)) and there is a backlog of almost 43 million cases, it is imperative to develop benchmarks and datasets for the Indian legal system. From the language perspective, **IL-TUR** benchmark covers English and 9 major Indian languages. Although **IL-TUR** is India-specific, the models developed for **IL-TUR** could provide inspiration (or possibly adapted) for developing models for the legal systems of other countries. Lastly and most importantly, **IL-TUR** covers tasks related to the common-law system as well as the civil law system. India has a predominantly common-law system, which implies that a judge in a higher court can overrule existing precedents, so the decision may not always be as per the rule book (written statutes and laws). It introduces some subjectivity into the decision-making process and must be backed by solid reasoning, making the tasks in **IL-TUR** much more difficult. Additionally, India has a civil law system for certain matters (e.g., banking and insurance). In the proposed benchmark, we cover both settings. Moreover, the legal domain has various areas (following common or civil systems) of laws such as criminal, civil, and banking; via the benchmark, we want to test the cross-area generalization capabilities of the models, i.e., how well the models developed on data from one area

| Task | Dataset (Language) | Avg. #Words | Task Type | Key Skills Required |
|------|---------------------|-------------|-----------|----------------------|
| L-NER | 105 docs 650k words (English) | 6,180 | Sequence Classification | Foundational task, legal understanding |
| RR | 21,184 sentences (English) | 25,796 | Multi-Class Classification | Foundational task, legal knowledge and legal semantics understanding |
| CJPE | ILDC 34k Docs (English) | 3,336 | Classification, Extraction | Legal understanding and reasoning |
| BAIL | HLDC 176k Docs (Hindi) | 86 | Classification | Legal understanding (in Hindi) and reasoning |
| LSI | ILSI 65k samples (English) | 2,406 | Multi-Label Classification | Understanding of the statutes and their applicability in various factual situations, commonsense knowledge and reasoning |
| PCR | IL-PCR 7,070 Docs (English) | 8,096 | Retrieval | Understanding of facts (commonsense + legal knowledge) and statutes, concept of legal relevance |
| SUMM | In-Abs 7,130 Docs (English) | 4,376 | Generation | Legal understanding and generation |
| L-MT | MILPaC 17,853 text pairs (English and Indian Langs.) | 49 | Generation | Parallel understanding of the legal text in English and 9 Indian languages |

Table 2: Summary of Tasks introduced in **IL-TUR**.

generalize across other areas. In contrast, Korea, China, (and, to a large extent, the EU) mainly follow civil law where a decision is as per the rule book. **IL-TUR** aims to fill the voids in the Legal NLP for the Indian setting by introducing some of the foundational tasks that can be useful for various legal applications.

## 3 IL-TUR: Legal-NLP Benchmark

Table 2 summarizes various tasks proposed in **IL-TUR**. The tasks cover multiple aspects of the legal domain and require specialized skills and knowledge to solve them.

### 3.1 Design Philosophy

We want to develop technology that enables automated semantic and legal understanding of legal documents and processes. We created **IL-TUR** with the following principles in mind.

**1) Legal Understanding and World Knowledge:** The tasks should cater exclusively to the legal domain. Solving a task should require in-depth knowledge and understanding of the law and its associated areas. Further, the tasks should not be restricted to only classification but should also involve retrieval, generation, and explanation. The proposed tasks address the pain points of processing legal texts (§1). Moreover, solving legal tasks should require knowledge about the law as well as

commonsense knowledge and societal norms about the world (e.g., facts in conjunction with socio-economic conditions in a particular case). **2) Difficulty Level:** The difficulty level should be such that these are not solvable by a layperson (having minimal knowledge and expertise in legal matters). It ensures that general language learners cannot easily solve the tasks, and the tasks would be sufficiently challenging for the current state-of-the-art models (e.g., LLMs). **3) Language:** Since India is a multi-lingual society, the tasks should cater to the most frequent languages used in the courts. We cover tasks in English and 9 other Indian languages. **4) Evaluation:** The tasks should be automatically evaluable, and the metrics used should align with human judgments. **5) Public Availability:** The data used for the tasks should be publicly available so anyone can use it for research purposes without licensing or copyright restrictions. Further, a leaderboard should be available to compare different systems and models. We release the data via a Creative Common Attribution-NonCommercial-ShareAlike (CC BY-NC-SA) license and create a public leaderboard.

### 3.2 IL-TUR Tasks

Based on the design philosophy, in this version of **IL-TUR**, we selected eight different tasks. Table 2 provides a summary of the tasks. We briefly describe the tasks here; details about the dataset

and evaluation metrics are provided in App. A.

- **Legal Named Entity Recognition (L-NER):** This is a newly created task in **IL-TUR**. Formally, given a legal document, the task of Legal Named Entity Recognition is to identify entities (set of 12 entity types), namely, Appellant, Respondent, Judge, Appellant Counsel, Respondent Counsel, Court, Authority, Witness, Statute, Precedent, Date, and Case Number. L-NER is different from the standard NER task; if one were to run a standard NER system on a legal document, the judge, petitioner, and respondent would all be labeled with a "PERSON" tag. Hence, a separate task is needed to identify the legal named entities in the documents. The standard NER (identifying person/organization/location names) can be done by any non-legal professional/person, but identifying the roles of entities involved in a legal case (L-NER) requires an in-depth understanding of the legal terminologies and the law. Hence, we develop a gold-standard dataset for L-NER with the help of law students (details in A.1). Moreover, the set of legal entities and corresponding definitions are formulated with the help of legal academicians (experts).

- **Rhetorical Role Prediction (RR):** As pointed out earlier, legal documents are typically long (avg. length 4000 words) and highly unstructured, with the legal information spread throughout the document. Segmenting the long documents into topically coherent units (such as facts, arguments, precedent, statute, etc.) helps highlight the relevant information and reduces human effort. These topically coherent units are termed as *Rhetorical Roles* (RR). Given a legal document, the task of RR prediction involves assigning RR label(s) to each sentence. We focus on 13 RR labels: *Fact, Issue, Arguments (Respondent), Argument (Petitioner), Statute, Dissent, Precedent Relied Upon, Precedent Not Relied Upon, Precedent Overruled, Ruling By Lower Court, Ratio Of The Decision, Ruling By Present Court, None.* Details about RR labels, definitions, and the dataset are provided in the App. A.2.

- **Court Judgment Prediction with Explanation (CJPE):** Formally, the task of Court Judgment Prediction with Explanation (CJPE) involves predicting the final judgment (appeal accepted or denied, i.e., the binary outcome of 0 or 1) for a given judgment document (having facts and other details) and providing the *explanation for the decision*. In this case, the explanations are in the form of the salient sentences that lead to the decision. Note that the idea behind this task is *not* to replace human judges but to augment them in decision-making. Furthermore, the task requires the system to explain its decision so that it is interpretable for a human (details in App. A.3).

- **Bail Prediction (BAIL):** A large fraction of the pending cases in India are from the district-level courts and have to do with bail applications (https://en.wikipedia.org/wiki/Bail) (Kapoor et al., 2022). Many of the district courts in India use Hindi as their official language (also refer to the Limitations section). Given a legal document in the Hindi language (having the facts of the case), the task of Bail Prediction involves predicting if the accused should be granted bail or not (i.e., a binary decision of 0/1) (details in App. A.4).

- **Legal Statue Identification (LSI):** The task of Legal Statute Identification (LSI) is formally defined to automatically identify the relevant statutes given the facts of a case. One of the first steps in the judicial process is finding the applicable statutes/laws based on the facts of the current situation. Manually rummaging through multiple legislation and laws to find out the relevant statutes can be time-consuming, making the LSI task important for reducing the workload and improving efficiency (more details in App. A.5).

- **Prior Case Retrieval (PCR):** When framing a legal document, legal experts (judges and lawyers) use their expertise to cite previous cases to support their arguments/reasoning. Legal experts have relied on their expertise to cite previous cases; however, with an exponentially growing number of cases, it becomes practically impossible to recall all possible cases. Given a query document (without citations), the task of Prior Case Retrieval (PCR) is to retrieve the legal documents from the candidate pool that are relevant (and hence can be cited) in the given query document (details in App. A.6).

- **Summarization (SUMM):** Summarization is a standard task in NLP; however, as mentioned in §1, summarizing legal documents

requires legal language understanding and reasoning. The task of summarization involves generating a gist (of a legal document) that captures the critical aspects of the case. We focus on abstractive summarization (more details in App. A.7).

- **Legal Machine Translation (L-MT):** In the Indian legal setting, when a case is transferred (due to re-appeal) from a district court to a High court, the corresponding document (typically in a regional language) needs to be translated to English. Additionally, since a large majority of the Indian population is not proficient in English, High Court / Supreme Court documents often need to be translated from English to Indian languages. In both scenarios, such translations, if done by humans, become a primary reason for delay in administering justice. Machine translation (MT) can augment human translators who can post-edit the translated document rather than translating from scratch. India is a diverse country with multiple languages across different states; the task of Legal Machine Translation (L-MT) attempts to close the language barrier by encouraging the development of systems for translating legal documents from English to Indian languages and vice-versa. Given that many Indian languages are low-resource, MT becomes even more challenging, requiring specialized models for translating legal documents in low-resource Indian languages. We focus on 9 Indian languages, namely, Bengali (BN), Hindi (HI), Gujarati (GU), Malayalam (ML), Marathi (MR), Telugu (TE), Tamil (TA), Punjabi (PA), and Oriya (OR) (details in App. A.8).

The tasks in **IL-TUR** require quite varied skills to solve the problem (Table 2). The skills include a deep understanding of language, the ability to generate legal language, foundational knowledge of law and statutes, application of law to social settings (e.g., decision-making in CJPE and BAIL), and the ability to reason using legal principles. The requirement of such a rich set of skills makes **IL-TUR** quite challenging; a single model struggles to solve all these tasks, as we observed in our experiments with BERT, LegalBERT, InLegalBERT, GPT3.5 and GPT-4 models (§4).

**Harmonization of Tasks:** This resource paper introduces a new benchmark for promoting research and development in the Indian legal system. Since it is a benchmark paper, the aim is to bring domain-specific tasks and datasets under one umbrella so that researchers can compare their models across tasks and with respect to each other. Earlier, no such effort was made for the Indian legal NLP domain. Some of the tasks included in the benchmark already exist; however, there is a lack of standardization across these, e.g., each task and dataset follows its file format, evaluation metric, etc. We have collated all these datasets and converted them to a uniform, JSON-based format so that the community can easily understand and use them. We have also collated all the training scripts for these different tasks together and devised a standard evaluation setup for all these tasks. Further, we have created a website (`https://exploration-lab.github.io/IL-TUR/`) and a public leaderboard that brings all relevant tasks together. The public leaderboard will further promote transparent and fair comparisons of techniques for each task. Moreover, the leaderboard will lead to the development of more sophisticated models (e.g., GLUE (Wang et al., 2018a) and SuperGLUE (Wang et al., 2019a) benchmarks promoted further research in NLP). Furthermore, to harmonize these tasks, we also conducted experiments with GPT-3.5 and GPT-4 (see §4) on all the tasks (except PCR), which involved converting the data to the desired format for GPT and formulating the prompts and verbalizers. Also, we plan to grow IL-TUR by introducing more new tasks in the future. We would also like to point out that many existing popular NLP benchmarks such as GLUE (Wang et al., 2018a), SuperGLUE (Wang et al., 2019a), as well as legal benchmarks like LEXGLUE (Chalkidis et al., 2022a) mostly comprised of datasets released by prior works. GLUE and LEXGLUE introduced only one new dataset each, whereas SuperGLUE did not have any new datasets.

**Anonymization of datasets:** In order to address ethical concerns (also see Ethical Considerations section) and to prevent the model from developing any bias, we anonymized named entities in the dataset of the relevant tasks, namely RR, CJPE, BAIL, LSI and PCR (details in App. A.10).

### 3.3 Relevance of Tasks to Litigation Process

In general, considering the pipeline of a litigation process for a case, all the tasks in the IL-TUR benchmark help formulate various ways in which automatic legal language processing can augment legal practitioners. Among the tasks, LSI is con-

| Task | SOTA | Metric | Model Details |
|------|------|--------|---------------|
| L-NER | 48.58% | strict mF1 | InLegalBERT + CRF |
| RR | 69.01% | mF1 | MTL-BERT |
| CJPE | 81.31%<br>0.56<br>0.32 | mF1<br>ROUGE-L<br>BLEU | InLegalBERT + BiLSTM |
| BAIL | 81% | mF1 | TF-IDF + IndicBERT |
| LSI | 28.08% | mF1 | LeSICiN (Graph-based Model) |
| PCR | 39.15% | $\mu$F1@$K$ | Event-Based |
| SUMM | 0.33<br>0.86 | ROUGE-L<br>BERTScore | Legal-LED |
| L-MT | 0.28<br>0.32<br>0.57 | BLEU<br>GLEU<br>chrF++ | MSFT (Microsoft Translation) |

Table 3: Summary of the best result for each task, along with the model that achieved the best result.

| Task | Arch | BERT | | | |
|------|------|------|------|------|------|
| | | V | L | InL | Ind |
| L-NER | Flat(CS) | 39.59% | 45.58% | 48.58% | - |
| RR | Flat(S) | 58% | 54% | 58% | - |
| CJPE | Hier | 71.14% | 78.21% | 81.31% | - |
| BAIL | Flat(L) | - | - | - | 76% |
| LSI | Hier | 18.44% | 21.74% | 26.23% | - |
| PCR | Flat(CS) | 9.24% | 8.67% | 7.57% | - |

Table 4: Results of different BERT-based models on tasks of **IL-TUR**. V, L, InL, and Ind refer to Vanilla BERT, LegalBERT, InLegalBERT, and IndicBERT, respectively. All metrics are in terms of macro-F1 (strict mF1 for L-NER). All the BERT-based models are implemented with the default architectures: either in a flat setup by taking individual sentences (S), segmenting long texts (CS), choosing the last 512 tokens for encoding (L), or in a hierarchical (Hier) setup with a BiLSTM on top of BERT.

sidered one of the first steps in the judicial process – right after identifying the facts, legal personnel must find out the statutes of the law that are violated. Since India follows a mixture of civil and common law systems, identifying the statutes is not the sole basis of legal reasoning; precedent cases must also be considered (PCR task). Subsequently, the final step in the litigation process is to decide the outcome of the case; the CJPE and BAIL tasks are relevant in this case, and human judges can use corresponding models to get suggestions/recommendations. The tasks of L-NER, RR, and SUMM, though not directly required for the judicial process, significantly help legal practitioners (e.g., lawyers conducting legal research to argue an ongoing case) get a quick understanding of the documents. Sometimes, a case gets re-appealed in a higher court, and consequently, the case document (in a regional language) in the lower court needs to be translated into English (L-MT Task).

## 4 Models, Experiments and Results

We extensively experimented with various models for each proposed task, including transformer-based language models. Table 3 summarizes baseline models and results for all tasks. Due to space limitations, we provide only the top-performing models here; details of experiments (e.g., hyper-parameters) and other models are in App. B. In general, results indicate that the tasks are far from being solved, and more research is required. In par-

ticular, we experimented with both generic BERT model (Devlin et al., 2019) and legal domain-specific BERT models: LegalBERT (Chalkidis et al., 2020) (BERT pre-trained on EU legal documents), CaseLawBERT (Zheng et al., 2021) (BERT pre-trained on US legal documents), and InLegal-BERT (Paul et al., 2023) (BERT pre-trained on Indian legal documents). For L-NER, InLegal-BERT (with CRF on top) shows the best performance, possibly because of in-domain data pre-training. For the RR task, vanilla BERT (or other transformers) and Legal-BERT do not work well; hence, RR prediction is posed as a sequence prediction problem (at the sentence level), and the Multi-Task Learning (MTL) model based on BERT developed by Malik et al. (2022) shows the best performance. Since legal documents are long, and BERT has a limitation of 512 tokens in the input, for the CJPE task, hierarchical InLegalBERT (InLe-galBERT and BiLSTM on top of that) (Paul et al., 2023) works best. For explanations, we use the occlusion method for finding the sentences leading to the final decision (Malik et al., 2021). But these fall short of expert-annotated important sentences in terms of ROUGE-L and BLEU scores. For BAIL prediction, since the documents are in Hindi, IndicBERT (Kakwani et al., 2020), a BERT model trained on Indian languages, was used. A pre-filtering of salient sentences, followed by In-dicBERT, works best (Kapoor et al., 2022). For the LSI task, we conduct experiments with hierarchical LegalBERT and InLegalBERT, along with LeSICiN, a graph-based method proposed by Paul et al. (2022). We observe that LeSICIN outper-

| Model | Trained On | ROUGE | | | BERTScore |
|---|---|---|---|---|---|
| | | R-1 | R-2 | R-L | |
| Summa-RuNNer | In-Abs | 0.604 | 0.264 | 0.225 | 0.828 |
| | UK-Abs | 0.493 | 0.255 | 0.274 | 0.849 |
| Legal LED | In-Abs | 0.557 | 0.244 | 0.242 | 0.844 |
| | UK-Abs | 0.482 | 0.186 | 0.264 | 0.851 |

Table 5: Performance of SummaRuNNer and Legal LED on UK-Abs test set after being trained on In-Abs (part of **IL-TUR**) and UK-Abs train sets.

forms the BERT-based methods. For the PCR task, an event-based model works the best (Joshi et al., 2023). An event refers to an action/activity (in the form of a predicate (typically a verb) and corresponding arguments) mentioned in the document. For SUMM, Legal-LED (HuggingFace, a) performs the best, and the commercially available Microsoft Azure Cognitive Services Translation API works best for the L-MT task. In general, across all tasks except LSI, PCR, SUMM, and L-MT, BERT (or its variant) performs the best. In order to compare the same model type across all tasks, we also experimented with BERT and its variants across all tasks that do not require text generation (i.e., SUMM and L-MT). Specifically, we took BERT (`bert-base-uncased`), LegalBERT, and In-LegalBERT as the encoders and ran it either in the flat text setup (either over the last 512 tokens or individual sentences/chunks) or hierarchical setup (full document), as per the task requirement. For BAIL, we use IndicBERT since the documents are in Hindi. These results are reported in Table 4. In general, except for PCR, we observe that performance increases going from BERT to LegalBERT to InLegalBERT, which correlates with the degree of in-domain pre-training.

**Model generalization beyond Indian jurisdiction:** Law is country/region specific. The laws of one country cannot be directly applied to another country. Hence, the legal NLP models developed for one region are less likely to generalize across countries. A similar pattern is also observed among human lawyers, i.e., an Indian lawyer cannot practice directly in EU/US jurisdictions. Moreover, even many of the tasks are jurisdiction-specific, e.g., tasks like BAIL, CJPE (since the processes used for deciding bail are different across countries), and LSI (since statutes are country-specific by nature), PCR, and L-MT require a deep understanding of the Indian legal system. For tasks like L-NER and RR, one could

test the generalization capabilities of models across jurisdictions; however, we could not find equivalent datasets (with the same labels) in other jurisdictions. For the summarization task, one can easily check the generalization capability of models across legal systems. So we conducted cross-jurisdiction experiments on abstractive summarization of *UK Supreme Court documents* using the UK-Abs dataset (Shukla et al., 2022). Uk-Abs consists of gold standard summaries released by the UK Supreme Court as press summaries. We experimented with two best-performing models: SummaRuNNer and Legal LED (trained on the In-Abs dataset, which is part of **IL-TUR**). These are used to generate summaries on the test set of UK-Abs (100 documents). Results are reported in Table 5. For comparison, we also report the results of these models on the UK-Abs test set when trained on the train set of UK-Abs itself. The results show that the summarization models trained on In-Abs perform decently when tested on UK-Abs (in a zero-shot setting). Both SummaRuNNer and Legal LED trained on In-Abs outperform their counterparts trained over UK-Abs in terms of ROUGE-1 and ROUGE-2 and achieve comparable ROUGE-L and BERT Scores. This experiment further shows the utility of our **IL-TUR** benchmark. Nevertheless, the generalization of models across jurisdictions requires more research in cross-jurisdiction domain adaptation techniques; we leave this for future work.

**Experiments with LLMs:** We also conducted experiments with LLMs. In particular, we experimented with large models (in terms of the number of parameters) like Open-AI GPT-3.5 (`gpt-3.5-turbo-16k`) and GPT-4 (`gpt-4-turbo`) and smaller models like GPT-Neo (Black et al., 2021) family of three models (GPT-Neo-125M, GPT-Neo-1.3B, GPT-Neo-2.7B), GPT-J-6B (Wang and Komatsuzaki, 2021), Llama-2-7b-chat-hf (Touvron et al., 2023), and Mistral-7B-v0.1 (Jiang et al., 2023). We experimented with zero-shot settings and In-Context Learning (ICL)-based settings (one-shot and two-shots). Table 6 shows the results for Open-AI GPT-3.5 and GPT-4 models (details about prompts and other settings in App. C). We could not experiment with ICL for PCR since it requires a comparison between the query document and the pool of all candidate documents, and passing the content of all the documents to GPT (or other LLMs) exceeds the token length limit even for GPT-4 (having context length of 16,000

| Task | GPT-3.5 | | | GPT-4 | | | SOTA | Metric |
|---|---|---|---|---|---|---|---|---|
| | 0-Shot | 1-Shot | 2-Shot | 0-Shot | 1-Shot | 2-Shot | | |
| L-NER | 30.59% | 23.68% | <u>32.84%</u>* | 13.65% | 10.51% | 24.03% | **48.58**% | strict mF1 |
| RR | 30.95% | 30.05% | 30.31% | 37.37% | 37.43% | <u>38.18%</u> | **69.01**% | mF1 |
| CJPE | 54.17% | 51.46% | 56.74% | <u>68.29%</u> | 47.26% | 60.44% | **81.31**% | mF1 |
| | 0.30 | 0.29 | 0.30 | 0.40 | 0.39 | <u>0.43</u> | **0.56** | ROUGE-L |
| | 0.08 | 0.15 | 0.113 | 0.14 | 0.16 | <u>0.18</u> | **0.32** | BLEU |
| BAIL | 51.04% | 46.35% | 61.0% | 51.46% | 56.90% | <u>66.67%</u> | **81**% | mF1 |
| LSI | 21.55% | 22.61% | 21.40% | <u>23.99</u> | 22.26 | 20.53 | **28.08**% | mF1 |
| SUMM | 0.21 | 0.20 | 0.22 | <u>0.23</u> | 0.16 | 0.17 | **0.33** | ROUGE-L |
| | <u>0.85</u> | 0.84 | 0.84 | <u>0.85</u> | 0.81 | 0.81 | **0.86** | BERTScore |
| L-MT | 0.23 | 0.25 | 0.26 | 0.33 | 0.35 | **0.36** | 0.28 | BLEU |
| | 0.28 | 0.28 | 0.29 | 0.36 | 0.38 | <u>**0.39**</u> | 0.32 | GLEU |
| | 0.42 | 0.43 | 0.43 | 0.50 | 0.52 | <u>0.53</u> | **0.57** | chrF++ |

Table 6: Performance of OpenAI GPT-3.5 (`gpt-3.5-turbo-16k`) and GPT-4 (`gpt-4-turbo`) model on various tasks for zero-shot, one-shot and two-shot settings. The SOTA corresponds to the best-performing model as given in Table 3. The best result for each task is marked in **boldface**. The best GPT-based result for each task is <u>underlined</u>.

tokens). In the future, we would like to experiment with a two-stage retrieval process to feed a subset of candidates to GPT instead of all of them. As observed, the GPT models perform worse than the SOTA models for each task, except for GPT-4 on MiLPAC for L-MT. It may be because the tasks are quite complex and require reasoning across long contexts. Also, for some tasks like L-NER and RR, it can be hard to come up with output formats that the models can understand in a zero-shot setting. Results for one-shot and two-shot show a similar trend in most cases. In some cases, one-shot performance is worse than zero-shot performance (also observed in other works (Brown et al., 2020)), while ICL has no effect on some tasks. Overall, in most cases and under most settings, GPT-4 outperforms GPT-3.5 by significant margins. Experiments with smaller models (GPT-Neo-125M, GPT-Neo-1.3B, GPT-Neo-2.7B, GPT-J-6B, Llama-2-7b-chat-hf, and Mistral-7B-v0.1) showed similar trends (details in App. C).

**Discussion:** Tasks in **IL-TUR** are quite varied, requiring different types of knowledge and skills. Developing systems for the legal domain is not easy due to the inherent challenges (§1). Moreover, legal datasets are expensive to annotate; consequently, annotated legal datasets are relatively small in size, and hence, learning in a low-resource setting is challenging. Experiments indicate that transformers fine-tuned on legal texts have shown limited success in the legal domain. Further, LLMs like Chat-GPT, which have demonstrated SOTA results in other domains (and have been shown to pass the bar exam (Chalkidis, 2023)), have not performed well on the IL-TUR benchmark, indicative of further research required in the legal domain.

## 5 Conclusion and Future Directions

This paper presented **IL-TUR**, a benchmark for Indian Legal Text Understanding and Reasoning. The benchmark has eight tasks requiring different types of legal skills to solve. Results indicate that the tasks are far from solved using state-of-the-art transformer-based models and LLMs. The list of tasks in **IL-TUR** is not exhaustive, and we plan to expand the list of tasks in the future; for example, we are working on developing foundational tasks like **Legal Coreference Resolution (L-Coref)** that are required for various applications such as information extraction and knowledge graph creation. Although such tasks have been addressed well in general NLP, our initial experiments show that using SOTA transformer models (which have become part of standard NLP toolkits) do not perform well on legal texts. Due to the usage of specialized terms, new models are needed for the legal domain. On the modeling side, in the future, we plan to develop one model that generalizes and works across all the tasks (e.g., mT5 (Xue et al., 2020) and Multi-task Adapters (Pfeiffer et al., 2020)). Overall, we hope that **IL-TUR** (along with its leaderboard) and its successive versions would create excitement in the Legal-NLP community and lead to the development of new technologies that could benefit society immensely and facilitate fair access to justice, a fundamental human right.

## Limitations

**IL-TUR** is a first step towards creating a benchmark for the Indian legal domain, which desperately needs technological solutions. The benchmark is not perfect and has certain limitations. Given the dynamic nature of the legal domain, new cases and precedents keep getting added. Hence, we plan to keep updating **IL-TUR** in the future. The legal domain is vast and covers various areas such as criminal law, civil law, banking, insurance, etc. In **IL-TUR**, we could not cover each of the subdomains in each task as it is a time-consuming and expensive affair to annotate many documents. One of our goals for **IL-TUR** is to test the cross-area generalization abilities of models; nevertheless, we would expand the datasets of each task in the future. **IL-TUR** is multi-lingual only concerning the L-MT task. Additionally, the BAIL task is in Hindi. All the High Courts and the Supreme Court in India use English as the official language. Hindi is the prominent language used in the district courts in most north Indian states. Nevertheless, India is a multi-lingual society, and legal models for other languages should also be developed for more tasks in the legal domain. We plan to extend the benchmark in the future and include some more tasks in Indian languages. The main challenge in doing so is a scarcity of legal data in regional languages in digitized formats from lower courts. Datasets of some of the tasks (e.g., LSI) use ML-based models (that may not be perfect) in the dataset creation process (e.g., fact extraction in the case of LSI). Extracting facts manually at a large scale is an expensive and time-consuming effort; in the future, we plan to employ legal professionals and create a more refined dataset. Regarding explainability, at present, we mainly address model explainability in the context of the CJPE task. For discussion regarding other tasks, please refer to App. A.9. Regarding LLM experiments, some of the tasks, such as BAIL and CJPE, require the entire document to be a part of the model's input. Obtaining LLM predictions overall test set samples is challenging in terms of expense and computation. Hence, we evaluated over a small subset, assuming that it is a good proxy of LLM performance. Lastly, the benchmark has only eight tasks. Creating legal tasks is time-consuming and expensive since it requires the help of legal experts. Nevertheless, as explained earlier, **IL-TUR** is a work in progress, and we will keep growing by adding more tasks. In this work,

we presented different models for various tasks; although many of the models (e.g., BERT, GPT) are common across all tasks, in the future, we plan to develop a single model that could solve all the tasks (e.g., mT5) with reasonable accuracy.

## Ethical Considerations

We use publicly available and open-source datasets for the tasks; no copyright is infringed. To the best of our knowledge, five of the proposed tasks (L-NER, RR, LSI, PCR, and Summ) do not have any direct ethical consequences since the proposed tasks are mainly related to information retrieval and summarization. Moreover, the tasks are meant to encourage the development of systems that would lead to streamlining the legal workflow and will not directly affect the life of any personnel.

For the LSI task, to prevent any bias in the model, named entities in the dataset were anonymized (details in App. A.10). Similarly, the named entities were anonymized in the RR and PCR datasets. App. A.10 provides more details about various measures and potential risks associated with failure to anonymize legal data. The documents are selected randomly for all tasks to avoid bias towards any entity, organization, or law.

Two tasks (CJPE and BAIL) have ethical considerations. Given a large quantum of pending cases in Indian courts, these tasks aim to develop systems that augment judges and *not* replace them; consequently, the systems are meant to provide recommendations, and a human judge takes the final decision. We follow all the steps as done by Malik et al. (2021); Kapoor et al. (2022) to avoid any bias in the data for these two tasks. For example, we removed cases (documents) related to sensitive issues like rape and sexual violence, and named entities were anonymized.

**Note that we do not endorse the use of the benchmark data for non-research (commercial and real-life) applications, and the primary motivation for creating the IL-TUR benchmark is to consolidate all the research happening in parallel for the Indian Legal domain.** Hence, we will release the benchmark and datasets under the Creative Common Attribution-NonCommercial-Share-Alike (CC BY-NC-SA) license. Moreover, we believe providing a platform by maintaining a common leaderboard for multiple tasks will advance the field with more transparency and reproducibility.

# References

Khalid Al-Kofahi, Alex Tyrrell, Arun Vachher, and Peter Jackson. 2001. A Machine Learning Approach to Prior Case Retrieval. In *Proceedings of the Eighth International Conference on Artificial Intelligence and Law (ICAIL)*.

Paheli Bhattacharya, Shounak Paul, Kripabandhu Ghosh, Saptarshi Ghosh, and Adam Wyner. 2019. Identification of Rhetorical Roles of Sentences in Indian Legal Judgments. *CoRR*, abs/1911.05405.

Sid Black, Gao Leo, Phil Wang, Connor Leahy, and Stella Biderman. 2021. GPT-Neo: Large Scale Autoregressive Language Modeling with Mesh-Tensorflow. If you use this software, please cite it using these metadata.

Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.

Ilias Chalkidis. 2023. Chatgpt may pass the bar exam soon, but has a long way to go for the lexglue benchmark. *arXiv preprint arXiv:2304.12202*.

Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. 2019. Neural Legal Judgment Prediction in English. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics(ACL)*.

Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. LEGAL-BERT: The muppets straight out of law school. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 2898–2904, Online. Association for Computational Linguistics.

Ilias Chalkidis, Manos Fergadiotis, Dimitrios Tsarapatsanis, Nikolaos Aletras, Ion Androutsopoulos, and Prodromos Malakasiotis. 2021. Paragraph-level rationale extraction through regularization: A case study on European court of human rights cases. In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 226–241, Online. Association for Computational Linguistics.

Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. 2022a. LexGLUE: A benchmark dataset for legal language understanding in English. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4310–4330, Dublin, Ireland. Association for Computational Linguistics.

Ilias Chalkidis, Tommaso Pasini, Sheng Zhang, Letizia Tomada, Sebastian Schwemer, and Anders Søgaard. 2022b. FairLex: A multilingual benchmark for evaluating fairness in legal text processing. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4389–4406, Dublin, Ireland. Association for Computational Linguistics.

Karl Cobbe, Chris Hesse, Jacob Hilton, and John Schulman. 2020. Leveraging procedural generation to benchmark reinforcement learning. In *International conference on machine learning*, pages 2048–2056. PMLR.

Jia Deng, Wei Dong, Richard Socher, Li-Jia Li, Kai Li, and Li Fei-Fei. 2009. Imagenet: A large-scale hierarchical image database. In *2009 IEEE conference on computer vision and pattern recognition*, pages 248–255. Ieee.

Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies (NAACL-HLT)*.

Zhiwei Fei, Xiaoyu Shen, Dawei Zhu, Fengzhe Zhou, Zhuo Han, Songyang Zhang, Kai Chen, Zongwen Shen, and Jidong Ge. 2023. Lawbench: Benchmarking legal knowledge of large language models. *arXiv preprint arXiv:2309.16289*.

Joseph L Fleiss, Bruce Levin, and Myunghee Cho Paik. 2013. *Statistical methods for rates and proportions*. john wiley & sons.

Leo Gao, Stella Biderman, Sid Black, Laurence Golding, Travis Hoppe, Charles Foster, Jason Phang, Horace He, Anish Thite, Noa Nabeshima, Shawn Presser, and Connor Leahy. 2020. The Pile: An 800gb dataset of diverse text for language modeling. *arXiv preprint arXiv:2101.00027*.

Alexios Gidiotis and Grigorios Tsoumakas. 2020. A divide-and-conquer approach to the summarization of academic articles. *CoRR*, abs/2004.06190.

Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi, Tom Zur, Varun Iyer, and Zehua Li. 2023. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models.

Yulan Guo, Jun Zhang, Min Lu, Jianwei Wan, and Yanxin Ma. 2014. Benchmark datasets for 3d computer vision. In *2014 9th IEEE Conference on Industrial Electronics and Applications*, pages 1846–1851. IEEE.

Prakhar Gupta, Matteo Pagliardini, and Martin Jaggi. 2019. Better word embeddings by disentangling contextual n-gram information. In *NAACL-HLT (1)*, pages 933–939. Association for Computational Linguistics.

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

Montani Ines Honnibal Matthew and Boyd Adriane Van Landeghem Sofie. 2020. spaCy: Industrial-Strength Natural Language Processing in Python.

Junjie Hu, Sebastian Ruder, Aditya Siddhant, Graham Neubig, Orhan Firat, and Melvin Johnson. 2020. Xtreme: A massively multilingual multi-task benchmark for evaluating cross-lingual generalisation. In *International Conference on Machine Learning*, pages 4411–4421. PMLR.

HuggingFace. a. Legal LED. https://huggingface.co/nsi319/legal-led-base-16384. [Online].

HuggingFace. b. Legal Pegasus. https://huggingface.co/nsi319/legal-pegasus. Online.

Wonseok Hwang, Dongjun Lee, Kyoungyeon Cho, Hanuhl Lee, and Minjoon Seo. 2022. A multi-task benchmark for korean legal language understanding and judgement prediction. *arXiv preprint arXiv:2206.05224*.

Peter Jackson, Khalid Al-Kofahi, Alex Tyrrell, and Arun Vachher. 2003a. Information extraction from case law and retrieval of prior cases. *Artificial Intelligence*, 150(1-2):239–290.

Peter Jackson, Khalid Al-Kofahi, Alex Tyrrell, and Arun Vachher. 2003b. Information Extraction from Case Law and Retrieval of Prior Cases. *Artificial Intelligence, Elsevier*.

Albert Q. Jiang, Alexandre Sablayrolles, Arthur Mensch, Chris Bamford, Devendra Singh Chaplot, Diego de las Casas, Florian Bressand, Gianna Lengyel, Guillaume Lample, Lucile Saulnier, Lélio Renard Lavaud, Marie-Anne Lachaux, Pierre Stock, Teven Le Scao, Thibaut Lavril, Thomas Wang, Timothée Lacroix, and William El Sayed. 2023. Mistral 7b.

Abhinav Joshi, Akshat Sharma, Sai Kiran Tanikella, and Ashutosh Modi. 2023. U-CREAT: Unsupervised Case Retrieval using Events extrAcTion. In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing*

*(Volume 1: Long Papers)*. Association for Computational Linguistics.

Divyanshu Kakwani, Anoop Kunchukuttan, Satish Golla, NC Gokul, Avik Bhattacharyya, Mitesh M Khapra, and Pratyush Kumar. 2020. Indicnlpsuite: Monolingual corpora, evaluation benchmarks and pre-trained multilingual language models for indian languages. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pages 4948–4961.

Prathamesh Kalamkar, Astha Agarwal, Aman Tiwari, Smita Gupta, Saurabh Karn, and Vivek Raghavan. 2022a. Named entity recognition in Indian court judgments. In *Proceedings of the Natural Legal Language Processing Workshop 2022*, pages 184–193, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics.

Prathamesh Kalamkar, Aman Tiwari, Astha Agarwal, Saurabh Karn, Smita Gupta, Vivek Raghavan, and Ashutosh Modi. 2022b. Corpus for Automatic Structuring of Legal Documents. In *Proceedings of the 13th Language Resources and Evaluation Conference -Association for Computational Linguistics (ACL-LREC)*.

Arnav Kapoor, Mudit Dhawan, Anmol Goel, Arjun T H, Akshala Bhatnagar, Vibhu Agrawal, Amul Agrawal, Arnab Bhattacharya, Ponnurangam Kumaraguru, and Ashutosh Modi. 2022. HLDC: Hindi Legal Documents Corpus. In *Findings of the Association for Computational Linguistics (ACL)*.

Nikolaos Lagos, Frederique Segond, Stefania Castellani, and Jacki O'Neill. 2010. Event extraction for legal case building and reasoning. In *International Conference on Intelligent Information Processing*, pages 92–101. Springer.

Michael Laskin, Denis Yarats, Hao Liu, Kimin Lee, Albert Zhan, Kevin Lu, Catherine Cang, Lerrel Pinto, and Pieter Abbeel. 2021. Urlb: Unsupervised reinforcement learning benchmark. *arXiv preprint arXiv:2110.15191*.

Mike Lewis, Yinhan Liu, Naman Goyal, Marjan Ghazvininejad, Abdelrahman Mohamed, Omer Levy, Veselin Stoyanov, and Luke Zettlemoyer. 2019. BART: denoising sequence-to-sequence pre-training for natural language generation, translation, and comprehension. *CoRR*, abs/1910.13461.

Chao-Lin Liu and Kuan-Chun Chen. 2019. Extracting the gist of chinese judgments of the supreme court. In *Proc. International Conference on Artificial Intelligence and Law (ICAIL)*.

Dayiheng Liu, Yu Yan, Yeyun Gong, Weizhen Qi, Hang Zhang, Jian Jiao, Weizhu Chen, Jie Fu, Linjun Shou, Ming Gong, et al. 2020. Glge: A new general language generation evaluation benchmark. *arXiv preprint arXiv:2011.11928*.

Sayan Mahapatra, Debtanu Datta, Shubham Soni, Adrijit Goswami, and Saptarshi Ghosh. 2023. Improving access to justice for the indian population: A benchmark for evaluating translation of legal text to indian languages.

Vijit Malik, Rishabh Sanjay, Shouvik Kumar Guha, Angshuman Hazarika, Shubham Nigam, Arnab Bhattacharya, and Ashutosh Modi. 2022. Semantic Segmentation of Legal Documents via Rhetorical Roles. In *Proceedings of the Natural Legal Language Processing Workshop (NLLP) EMNLP*.

Vijit Malik, Rishabh Sanjay, Shubham Kumar Nigam, Kripabandhu Ghosh, Shouvik Kumar Guha, Arnab Bhattacharya, and Ashutosh Modi. 2021. ILDC for CJPE: Indian Legal Documents Corpus for Court Judgment Prediction and Explanation. In *Proceedings of the 59th Annual Meeting of the Association for Computational Linguistics and the 11th International Joint Conference on Natural Language Processing (ACL-IJCNLP)*.

Marie-Francine Moens, Caroline Uyttendaele, and Jos Dumortier. 1999. Abstracting of legal cases: the potential of clustering based on the selection of representative objects. *Journal of the American Society for Information Science*, 50(2):151–161.

Ramesh Nallapati, Feifei Zhai, and Bowen Zhou. 2017. Summarunner: A recurrent neural network based sequence model for extractive summarization of documents. In *Proc. AAAI Conference on Artificial Intelligence*.

Shashi Narayan, Shay B Cohen, and Mirella Lapata. 2018. Ranking sentences for extractive summarization with reinforcement learning. In *Proc. NAACL-HLT*, pages 1747–1759.

National Judicial Data Grid. 2023. National judicial data grid statistics. https://www.njdg.ecourts.gov.in/njdgnew/index.php.

Joel Niklaus, Veton Matoshi, Pooja Rani, Andrea Galassi, Matthias Stürmer, and Ilias Chalkidis. 2023. Lextreme: A multi-lingual and multi-task benchmark for the legal domain. *arXiv preprint arXiv:2301.13126*.

Shounak Paul, Pawan Goyal, and Saptarshi Ghosh. 2022. Lesicin: A heterogeneous graph-based approach for automatic legal statute identification from indian legal documents. *Proceedings of the AAAI Conference on Artificial Intelligence*, 36(10):11139–11146.

Shounak Paul, Arpan Mandal, Pawan Goyal, and Saptarshi Ghosh. 2023. Pre-trained language models for the legal domain: A case study on indian law.

Jonas Pfeiffer, Ivan Vulić, Iryna Gurevych, and Sebastian Ruder. 2020. Mad-x: An adapter-based framework for multi-task cross-lingual transfer. *arXiv preprint arXiv:2005.00052*.

Seth Polsley, Pooja Jhunjhunwala, and Ruihong Huang. 2016. Casesummarizer: A system for automated summarization of legal texts. In *Proc. Iinternational conference on Computational Linguistics (COLING) System Demonstrations*.

Gowtham Ramesh, Sumanth Doddapaneni, Aravinth Bheemaraj, Mayank Jobanputra, Raghavan AK, Ajitesh Sharma, Sujit Sahoo, Harshita Diddee, Mahalakshmi J, Divyanshu Kakwani, Navneet Kumar, Aswin Pradeep, Srihari Nagaraj, Kumar Deepak, Vivek Raghavan, Anoop Kunchukuttan, Pratyush Kumar, and Mitesh Shantadevi Khapra. 2022. Samanantar: The largest publicly available parallel corpora collection for 11 Indic languages. *Transactions of the Association for Computational Linguistics*, 10:145–162.

Joshua Robinson and David Wingate. 2023. Leveraging large language models for multiple choice question answering. In *The Eleventh International Conference on Learning Representations*.

Isabel Segura-Bedmar, Paloma Martínez, and María Herrero-Zazo. 2013. SemEval-2013 task 9 : Extraction of drug-drug interactions from biomedical texts (DDIExtraction 2013). In *Second Joint Conference on Lexical and Computational Semantics (*SEM), Volume 2: Proceedings of the Seventh International Workshop on Semantic Evaluation (SemEval 2013)*, pages 341–350, Atlanta, Georgia, USA. Association for Computational Linguistics.

Abhay Shukla, Paheli Bhattacharya, Soham Poddar, Rajdeep Mukherjee, Kripabandhu Ghosh, Pawan Goyal, and Saptarshi Ghosh. 2022. Legal case document summarization: Extractive and abstractive methods and their evaluation. In *Proceedings of the 2nd Conference of the Asia-Pacific Chapter of the Association for Computational Linguistics and the 12th International Joint Conference on Natural Language Processing (Volume 1: Long Papers)*, pages 1048–1064, Online only. Association for Computational Linguistics.

Benjamin Strickson and Beatriz De La Iglesia. 2020. Legal judgement prediction for uk courts. In *Proceedings of the 2020 The 3rd International Conference on Information Science and System*, pages 204–209.

Hugo Touvron, Louis Martin, Kevin Stone, Peter Albert, Amjad Almahairi, Yasmine Babaei, Nikolay Bashlykov, Soumya Batra, Prajjwal Bhargava, Shruti Bhosale, Dan Bikel, Lukas Blecher, Cristian Canton Ferrer, Moya Chen, Guillem Cucurull, David Esiobu, Jude Fernandes, Jeremy Fu, Wenyin Fu, Brian Fuller, Cynthia Gao, Vedanuj Goswami, Naman Goyal, Anthony Hartshorn, Saghar Hosseini, Rui Hou, Hakan Inan, Marcin Kardas, Viktor Kerkez, Madian Khabsa, Isabel Kloumann, Artem Korenev, Punit Singh Koura, Marie-Anne Lachaux, Thibaut Lavril, Jenya Lee, Diana Liskovich, Yinghai Lu, Yuning Mao, Xavier Martinet, Todor Mihaylov, Pushkar Mishra, Igor Molybog, Yixin Nie, Andrew Poulton, Jeremy Reizenstein, Rashi Rungta, Kalyan Saladi, Alan Schelten,

Ruan Silva, Eric Michael Smith, Ranjan Subramanian, Xiaoqing Ellen Tan, Binh Tang, Ross Taylor, Adina Williams, Jian Xiang Kuan, Puxin Xu, Zheng Yan, Iliyan Zarov, Yuchen Zhang, Angela Fan, Melanie Kambadur, Sharan Narang, Aurelien Rodriguez, Robert Stojnic, Sergey Edunov, and Thomas Scialom. 2023. Llama 2: Open foundation and fine-tuned chat models.

Vu Tran, Minh Le Nguyen, and Ken Satoh. 2019. Building legal case retrieval systems with lexical matching and summarization using a pre-trained phrase scoring model. In *Proceedings of the Seventeenth International Conference on Artificial Intelligence and Law*, pages 275–282.

United Nations. 2023. India overtakes china as the world's most populous country. https://www.un.org/development/desa/dpad/wp-content/uploads/sites/45/PB153.pdf.

Alex Wang, Yada Pruksachatkun, Nikita Nangia, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel Bowman. 2019a. Superglue: A stickier benchmark for general-purpose language understanding systems. *Advances in neural information processing systems*, 32.

Alex Wang, Amanpreet Singh, Julian Michael, Felix Hill, Omer Levy, and Samuel R Bowman. 2018a. Glue: A multi-task benchmark and analysis platform for natural language understanding. *arXiv preprint arXiv:1804.07461*.

Ben Wang and Aran Komatsuzaki. 2021. GPT-J-6B: A 6 Billion Parameter Autoregressive Language Model. https://github.com/kingoflolz/mesh-transformer-jax.

Pengfei Wang, Yu Fan, Shuzi Niu, Ze Yang, Yongfeng Zhang, and Jiafeng Guo. 2019b. Hierarchical Matching Network for Crime Classification. In *Proceedings of the 42nd International ACM Conference on Research and Development in Information Retrieval, (SIGIR)*.

Pengfei Wang, Ze Yang, Shuzi Niu, Yongfeng Zhang, Lei Zhang, and ShaoZhang Niu. 2018b. Modeling Dynamic Pairwise Attention for Crime Classification over Legal Articles. In *The 41st International ACM Conference on Research & Development in Information Retrieval (SIGIR)* .

Yi Wu, Jongwoo Lim, and Ming-Hsuan Yang. 2013. Online object tracking: A benchmark. In *Proceedings of the IEEE conference on computer vision and pattern recognition*, pages 2411–2418.

Liang Xu, Hai Hu, Xuanwei Zhang, Lu Li, Chenjie Cao, Yudong Li, Yechen Xu, Kai Sun, Dian Yu, Cong Yu, et al. 2020. Clue: A chinese language understanding evaluation benchmark. *arXiv preprint arXiv:2004.05986*.

Linting Xue, Noah Constant, Adam Roberts, Mihir Kale, Rami Al-Rfou, Aditya Siddhant, Aditya Barua, and

Colin Raffel. 2020. mt5: A massively multilingual pre-trained text-to-text transformer. *arXiv preprint arXiv:2010.11934*.

Vikas Yadav and Steven Bethard. 2019. A survey on recent advances in named entity recognition from deep learning models. *arXiv preprint arXiv:1910.11470*.

Amy Zhang, Yuxin Wu, and Joelle Pineau. 2018. Natural environment benchmarks for reinforcement learning. *arXiv preprint arXiv:1811.06032*.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q Weinberger, and Yoav Artzi. 2019. Bertscore: Evaluating text generation with bert. *arXiv preprint arXiv:1904.09675*.

Lucia Zheng, Neel Guha, Brandon R. Anderson, Peter Henderson, and Daniel E. Ho. 2021. When does pre-training help? assessing self-supervised learning for law and the casehold dataset of 53,000+ legal holdings. In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law*, ICAIL '21, page 159–168, New York, NY, USA. Association for Computing Machinery.

## Appendix

## Table of Contents

## List of Figures

# A  Tasks and Dataset Details

We will release the baseline codes along with a compiled list of task-specific datasets and evaluation scripts with the camera-ready version of the paper. The consolidated leaderboard website for the benchmark will be made public with the camera-ready release. We now describe the tasks in the benchmark in detail.

## A.1  Legal Named Entity Recognition (L-NER)

**Task Motivation and Description:** Named Entity Recognition (NER) is a foundational task in NLP (Yadav and Bethard, 2019). However, in the legal domain, the types of named entities one may be interested in differ (e.g., judge, petitioner (appellant), and respondent), which may not be identified by a standard NER system. Hence, a separate task is needed to identify the legal named entities in the documents. Note that L-NER is very different from the standard NER; the standard NER (identifying person/organization/location names) requires a language understanding; in contrast, identifying the roles of entities involved in a legal case (L-NER) requires an understanding of the legal terminologies. Hence, we develop a gold-standard dataset for L-NER annotated with the help of law students (details in App. A). Moreover, the set of legal entities and corresponding definitions are formulated with the help of legal academicians (experts). **Formally, given a legal document, the task of Legal Named Entity Recognition is to identify entities (set of 12 entity types), namely, Appellant, Respondent, Judge, Appellant Counsel, Respondent Counsel, Court, Authority, Witness, Statute, Precedent, Date, and Case Number.** Fig. 2 shows an example. Table 7 shows the definition for 12 NE types/classes.

**Dataset:** We collected a total of 105 case documents in English (a total of 650K words and 12.5K entities). Table 8 lists some important statistics about the NER dataset. The NE type label statistics are displayed along with the class descriptions in Table 7.

**Annotation Details:** For the L-NER task, we collected a total of 105 cases publicly available from the Supreme Court and a few High Courts of India by scrapping the website: `https://www.indiankanoon.org`. Please note that the IndianKanoon website allows free downloads of public documents. In discussion with legal experts, we

decided on a comprehensive set of 12 NE (Named Entity) classes suited for the legal domain (Table 7). Two law students from a reputed law college in India were tasked with annotating the case documents. The annotation procedure involved the following steps:

- To ensure that entity spans are marked consistently, we discussed with both annotators how to mark every label. Such decisions involved leaving out prefixes/salutations such as 'Shri' (a polite way to address Mr. in the Indian context) and 'Smt.' (a polite way to address Ms. in the Indian context), 'Justice' (Honorific for a Judge), etc., from the entity names, including the (optional) precedent citations that follow case titles as part of the precedent (PREC) entities, and so on.
- We randomly chose a set of 25 documents, and each annotator worked on all 25 documents independently based on the rules devised in the previous step.
- We observed a high degree of agreement between the annotators for these 25 documents (Cohen's Kappa: 0.82, Krippendorff's Alpha: 0.85).
- Both annotators worked together to resolve the disagreements to arrive at one single consolidated set of annotations for these 25 documents.
- Above steps performed over 25 documents calibrated the annotators and led to a high degree of agreement among them. Since annotation is an expensive and time-consuming process, the remaining 80 documents were split equally between the two annotators for annotation.

**Task Evaluation:** NER can be formulated as a sequence prediction task, where each word receives either of the labels {B-X, I-X, O} as per the popular 'B-I-O' scheme (Yadav and Bethard, 2019) ('X' represents any of the legal classes we are interested in). We use standard metrics of *strict* macro-averaged precision, recall, and F1 score for evaluation. The *strict* score assumes a correct match only if *both* the entity boundary and entity type are correctly predicted. L-NER evaluation F1 score metric is computed using `https://pypi.org/project/nervaluate/`. We use strict macro-averaged scores in our setup. The *strict* scoring mechanism ensures that a match is considered correct if the entity span and entity type are the same. In other words, if either the span is incorrect (the model predicts more/fewer tokens as

| Broad Category | Label | Frequency | Description |
|---|---|---|---|
| Party | APPELLANT (APP) | 660 | Party filing an appeal to the court |
| | RESPONDENT (RESP) | 516 | Party against whom appeal has been filed |
| Legal Professional | JUDGE (JUD) | 366 | Judge of the current or prior/cited cases |
| | A.COUNSEL (AC) | 288 | Lawyer(s) on behalf of the appellant(s) |
| | R.COUNSEL (RC) | 255 | Lawyer(s) on behalf of the respondent(s) |
| Organizations | COURT (CRT) | 1,572 | Any court occurring in the document |
| | AUTHORITY (AUTH) | 1,342 | Any organization/body having administrative/legal authority |
| Other Person(s) | WITNESS (WIT) | 312 | Witness(es) who are testifying in the case |
| Legal References | STATUTE (STAT) | 2,055 | Citation to legal acts |
| | PRECEDENT (PREC) | 1,804 | Citation to prior cases |
| Legal Artefacts | DATE | 2,316 | Mention of any date in the case |
| | CASE NO. (CN) | 1,102 | Mention of any case number, including that of the current case |

Table 7: Named Entity (NE) types used in the L-NER dataset

| | |
|---|---|
| # Documents | 105 |
| # Labels | 12 |
| Total no. of words | 648,937 |
| Avg. Document Size (in #words) | 6180.35 |
| Total no. of entities (All occurrences) | 12,588 |
| Total no. of entities (Unique occ.) | 5,658 |
| Avg. no. of entities per doc (All occ.) | 119.89 |
| Avg. no. of entities per doc (Unique occ.) | 53.89 |

Table 8: The dataset statistics for the L-NER task

part of the entity) or the predicted label type does not match the ground truth, the match is considered incorrect.

**Comparison with existing L-NER datasets:** Recently, Kalamkar et al. (2022a) released a dataset for L-NER over Indian legal documents. However, unlike our dataset, which comprises of *full-length documents* annotated with every occurrence of every NE, the dataset by Kalamkar et al. (2022a) consists of segments of documents and not full documents. This is a crucial difference since models trained on our data will be able to detect NEs even when provided with a snippet of a case document. There is also a slight variation in the set of NEs considered in our dataset as compared to those considered by Kalamkar et al. (2022b), although most common entity types have been covered in both datasets.



Figure 2: Example of L-NER

## A.2 Rhetorical Role Prediction (RR)

**Task Motivation and Description:** As pointed out earlier, legal documents are typically long (avg. length 4000 words) and highly unstructured, with legal information spread throughout the document. Segmenting the long documents into topically coherent units (such as facts, arguments, precedent, statute, etc.) not only helps highlight the relevant information but also reduces human effort when going through a long list of documents. These topically coherent units are termed as *Rhetorical Roles* (RR). **Given a legal document, the task of RR prediction involves assigning RR label(s) to each sentence.** The sentences are annotated with 13 RRs by as many as six legal experts (from a reputed Indian law school). The 13 RR labels are: *Fact, Issue, Arguments (Respondent), Argument (Petitioner), Statute, Dissent, Precedent Relied Upon, Precedent Not Relied Upon, Precedent Overruled, Ruling By Lower Court, Ratio Of The Decision, Ruling By Present Court, None.* The definition of each RR label is given in Table 9. We utilize the dataset and role definitions provided by prior work on structuring Indian legal documents (Malik et al., 2022). Fig. 3 shows an excerpt from a legal document annotated with RR labels. RR Prediction is a foundational task that helps structure the information and thus aids downstream applications related to document understanding, information extraction, summarization, and retrieval.

**Dataset:** For this task, we use the dataset developed by Malik et al. (2022) primarily due to the large number of annotations by several Law academicians and public availability. The dataset consists of 21,184 sentences from legal documents (in English) about banking and competition law. The RR dataset was created by scrapping (from IndianKanoon website: https://indiankanoon.org/) publicly available documents from the Supreme Court of India, High Courts, and Tribunal courts. The documents pertain to Banking/Income Tax law (IT) and Competition Law (CL) (also called as Anti-Trust Law in the US). The dataset consists of 21,184 sentences annotated with 13 RRs. Figure 4 shows the distribution of RR labels. The dataset is split randomly (at document level) into 80% train, 10% validation, and 10% test set.

**Annotation Details:** The dataset was annotated by six legal experts (graduate law student researchers), three annotated CL documents, and the remaining three annotated IT documents (Malik et al., 2022). The annotators showed a high degree of agreement. The Fleiss kappa (Fleiss et al., 2013) between the annotators is 0.65 for the IT domain and 0.87 for the CL domain, indicating a substantial agreement between annotators. Annotating RR is not a trivial task, and annotators can have disagreements. Several strategies were employed to resolve these disagreements. More details about annotation case studies can be found in Malik et al. (2022).

**Evaluation:** RR Prediction is evaluated using standard Macro F1 metric. Macro F1 is the average F1 score calculated per class.

## A.3 Court Judgment Prediction with Explanation (CJPE)

**Task Motivation and Description:** The task of Court Judgment Prediction with Explanation (CJPE) aims to augment a judge in the judicial decision-making process by predicting the final outcome of the case. Note that the idea behind this task is *not* to replace human judges but to aid them. Furthermore, the task requires the system to explain its decision so that it is interpretable for a human using it. **Formally, the task of Court Judgment Prediction with Explanation (CJPE) involves predicting the final judgment (appeal accepted or denied, i.e., the binary outcome of 0 or 1) for a given judgment document (having facts and other details) and providing the explanation for the decision.** The explanations, in this case, are in the form of the crucial sentences appearing in the input text that lead to the decision.

**Dataset Details:** For the CJPE task, we use the Indian Legal Document Corpus (ILDC) (Malik et al., 2021). ILDC is a corpus of 34k legal judgment documents (in English) from the Supreme Court of India. Each document is annotated with the ground truth (actual decision given by the judge); further, a small subset of the documents are annotated with explanations by legal experts. This makes it a suitable dataset to consider for a legal understanding benchmark as it covers both judgment as well as relevant explanations annotated by human experts. Table 10 provides dataset statistics. Some cases consist of multiple appeals, which can contain corresponding decisions for each appeal. However, since the task has been posed as binary text classification, the final decision is considered as ACCEPT if *at least one* appeal is accepted, otherwise REJECT. The documents are stripped of the final decision given by the Judge with the help of

Figure 3: Example of the Rhetorical Role Prediction Task (Kalamkar et al., 2022b)



Figure 4: Distribution of RR labels in IT and CL documents (Malik et al., 2022).

regex-based matching. Table 10 provides details of the dataset.

Regarding ethical concerns, we follow Malik et al. (2021) who took various steps, such as normalizing the dataset concerning named entities to remove any biases in the data (also check the Ethical Considerations section).

**Annotation Details:** The explanation aspect of the CJPE task was annotated with the help of 5 legal experts (Malik et al., 2021). The annotators were graduate students and a law professor from a reputed law school. The annotators were not shown the final decision of the case. They were asked to predict the final decision and annotate the sentences (explanations) in the document that led to the final decision. More details about agreement among the annotators are provided in (Malik et al., 2021). In a nutshell, the average prediction F1 score of annotators w.r.t. to the ground truth judgment was 94.32%. This points towards the

challenging nature of the CJPE task; as pointed out earlier, India has a common-law system, and hence, judges could override existing precedents. Disagreements among the annotators were mainly due to differences in the linguistic interpretation of the case and law. For the explanation part, similar trends are reported with the average agreement in terms of the BLEU score to be around 0.4.

**Evaluation:** The prediction part of the CJPE task is evaluated using standard F1 score metric, and the explanation part is evaluated using BLEU and ROUGE scores.

### A.4 Bail Prediction (BAIL)

**Task Motivation and Description:** A large fraction of the pending cases in India are from the district-level courts, and have to do with bail applications (https://en.wikipedia.org/wiki/Bail) (Kapoor et al., 2022). Many of the district courts in India use Hindi as their official language (also refer to the Limitations section). Given the importance of Hindi (the most frequently spoken/written language in India), the task of Bail Prediction for Hindi legal documents is of immense importance, incorporating both language diversity and wider applicability in the Indian legal system. **Formally, given a legal document (having the facts of the case), the task of Bail Prediction involves predicting if the accused should be granted bail or not (i.e., a binary decision of 0 and 1).**

**Dataset:** For the task of BAIL prediction, Kapoor et al. (2022) created a corpus of 900k Hindi Legal Documents (referred to as HLDC (Hindi Legal Document Copus)). The corpus is created

| Rhetorical Role Label | Definition |
|---|---|
| Fact (FAC) | These are the facts specific to the case based on which the arguments have been made and judgment has been issued. In addition to Fact, we also have the fine-grained label |
| Issues (ISS) | The issues which have been framed/accepted by the present court for adjudication. |
| Argument Petitioner (ARG-P) | Arguments which have been put forward by the petitioner/appellant in the case before the present court and by the same party in lower courts (where it may have been petitioner/respondent) |
| Argument Respondent (ARG-R) | Arguments which have been put forward by the respondent in the case before the present court and by the same party in lower courts (where it may have been petitioner/respondent) |
| Statute (STA) | The laws referred to in the case. |
| Dissent (DIS) | Any dissenting opinion expressed by a judge in the present judgment/decision. |
| Precedent Relied Upon (PRE-R) | The precedents which have been relied upon by the present court for adjudication. These may or may not have been raised by the advocates of the parties and amicus curiae. |
| Precedent Not Relied Upon (PRE-NR) | The precedents which have not been relied upon by the present court for adjudication. These may have been raised by the advocates of the parties and amicus curiae. |
| Precedent Overruled (PRE-O) | Any precedents (past cases) on the same issue that have been overruled through the current judgment. |
| Ruling By Lower Court (RLC) | Decisions of the lower courts which dealt with the same case. |
| Ratio Of The Decision (ROD) | The principle that has been established by the current judgment/decision which can be used in future cases. Does not include the obiter dicta which is based on observations applicable to the specific case only. |
| Ruling By Present Court (RPC) | The decision of the court on the issues that have been framed/accepted by the present court for adjudication. |
| None (NON) | any other matter in the judgment which does not fall in any of the above-mentioned categories. |

Table 9: Definitions for different Rhetorical Roles

| Corpus (Avg. tokens) | Number of docs (Accepted Class %) | | |
|---|---|---|---|
| | Train | Validation | Test |
| **ILDC-multi (3231)** | 32305 (41.43%) | 994 (50%) | 1517 (50.23%) |
| **ILDC-single (3884)** | 5082 (38.08%) | | |
| **ILDC-expert (2894)** | 56 (51.78%) | | |

Table 10: Statistics for the CJPE dataset (ILDC) (Malik et al., 2021)

by scrapping publicly available documents on the eCourts website (https://ecourts.gov.in/ecourts_home/). The documents are scrapped from district courts of the state of Uttar Pradesh (a Hindi-speaking state in northern India). The data is anonymized to take care of biases and ethical aspects; please refer to (Kapoor et al., 2022) for more details. Bail cases in HLDC are pre-processed to remove the final decision (using regex) since we aim to predict this automatically. For the task of Bail prediction we selected only the documents related Bail cases from HLDC, this resulted in 176K documents, having 86 words per document on average. More details about the dataset are discussed in (Kapoor et al., 2022). For model training and evaluation, we divide the data into train, validation, and test split in the ratio of 70:10:20.

**Evaluation:** The BAIL prediction is a binary task; it is evaluated using the standard macro-F1 score metric.

## A.5 Legal Statute Identification (LSI)

**Task Motivation and Description:** One of the first steps in the judicial process is finding the applicable statutes/laws based on the facts of the current situation. Manually rummaging through multiple legislation and laws to find out the relevant statutes can be time-consuming, making the LSI task important for reducing the workload, helping improve the efficiency of the judicial system. **The task of Legal Statute Identification (LSI) is formally defined to automatically identify the relevant statutes given the facts of a case.** An example of the LSI task is presented in the Table 11. We utilize the ILSI dataset for this task, which comprises of 100 target statutes from the Indian Penal Code (IPC), the main legislation codifying criminal laws in India.

**Dataset:** For LSI, we use the Indian Legal

| Facts of the case | "On the fateful day at about 9.30 a.m. deceased accompanied by [PERSON1] (PW 4) and [PERSON2] (PW 7) was going from his village Talod to Alote. The accused persons were hiding behind bushes on the road near village Gharola. They were armed with lathies and farsies. When the deceased and the aforesaid two persons reached near the Khakhra, the respondents surrounded them and started attacking the deceased with weapons with which they were armed. His nose was cut. PWs. 4 and 7 tried to intervene, but they were also attacked by the accused persons as a result of which they also received injuries. The two witness rushed to the police station where PW 4 lodged the FIR (Exhibit P-10). The deceased in injured condition was taken to the hospital, and later he succumbed to the injuries. Post-mortem was conducted and large number of injuries were found on his body. During investigation the alleged weapons of the assailants were seized. After investigation charge sheet was placed." |
|---|---|
| IPC S.324 | *Voluntarily causing hurt by dangerous weapons or means* |
| IPC S.302 | *Punishment for murder* |

Table 11: Example of the LSI task, fact section taken a High Court Document **"State Of Madhya Pradesh vs. Mansingh And Ors. on 13 August, 2003"**, along with the IPC Sections (324 and 302) that the case cites.

Statute Identification (ILSI) dataset (Paul et al., 2022). The dataset consists of fact portions of 65k court case documents (derived from criminal court cases from the Supreme Court of India (SCI) and 6 High Courts of India). The Indian Penal Code (IPC) comprises most criminal statutes and procedures in India; the 100 most frequently occurring statutes in the IPC were chosen as the target statutes. The original ILSI dataset released by Paul et al. (2022) contains named entities. In line with recent works in legal NLP (Malik et al., 2021), we anonymize the dataset by masking entities of types 'PERSON' and 'ORGANIZATION' to remove any possible bias. Table 12 lists some important statistics about the ILSI dataset. In addition to the facts extracted from case documents and their corresponding statute mappings, Paul et al. (2022) also provided the statute descriptions as part of the dataset.

**Dataset Preprocessing:** The LSI task requires the input to be *only* the facts of the case, and thus, an automated RR method (Bhattacharya et al., 2019) was employed to extract the facts. Since this method is not foolproof, some sentences containing statute citations may get mislabeled as facts. The version of the dataset released by Paul et al. (2022) contains some unmasked statute citations. Thus, we used an existing automated Legal NER method (Kalamkar et al., 2022a), which can identify both the act/law names and the statute/section references in the text, to mask all possible statute and act references (statutes from all acts were masked, not just IPC). To prevent model biases, we also masked all entities identified by the Legal NER method.

**Evaluation:** LSI is formulated as a multi-label text classification task. The facts, a functional segment of the entire case document, are provided as in-

put. The expected output is one or more statutes from a list of target statutes relevant to the given fact portion. Standard classification metrics such as macro-averaged precision, recall, and F1 score are used for evaluation. In principle, the LSI task can also be considered a retrieval task instead of a multi-label classification task, i.e., the task is to retrieve from a dynamic set of statutes and provide a bigger pool of relevant candidates to be retrieved for a particular query document having facts. However, as it is an initial phase of establishing the benchmark, we followed the classification setting proposed in previous works (Wang et al., 2018b, 2019b; Chalkidis et al., 2019, 2021).

| Dataset | ILSI |
|---|---|
| # Documents | 66,090 |
| # Labels | 100 |
| Train/Dev/Test Split | 42,835/ 10,200/ 13,039 |
| Avg. Document Size (in #words) | 2406 |
| Avg. no. of citations (#labels per doc) | 3.78 |

Table 12: The table shows the dataset statistics for the ILSI dataset (Paul et al., 2022).

### A.6 Prior Case Retrieval (PCR)

**Task Motivation and Description:** When framing a legal document, legal experts (judges and lawyers) use their expertise to cite previous cases to support their arguments/reasoning. Legal experts have relied on their expertise to cite previous cases; however, with an exponentially growing number of cases, it becomes practically impossible to recall all possible cases. **Given a query document (without citations), the task of Prior Case Retrieval (PCR) is to retrieve the legal documents from**

the candidate pool that are relevant (and hence can be cited) in the given query document. Automating this process directly impacts the justice delivery logistics. Moreover, including this task in the benchmark incorporates the retrieval aspects and understanding of legal similarity (as opposed to semantic similarity), opening research directions for retrieval systems in the legal domain.

**Dataset:** For the task of PCR we use the `Indian Legal Prior Case Retrieval (IL-PCR)` corpus (Joshi et al., 2023). To the best of our knowledge, IL-PCR is the largest publicly available retrieval dataset for the Indian judicial system, making it a suitable candidate to be added to the benchmark. The IL-PCR corpus was created by scraping legal documents (available in the public domain) from the website IndianKanoon (https://indiankanoon.org/). The pool of documents is expanded by scraping documents cited by documents scraped previously. It was done to ensure sufficient citation links from the query to the candidate pool in the final dataset. Names of individuals and organizations were anonymized to the <NAME> and <ORG> tags, respectively, using a NER model (Honnibal Matthew and Van Landeghem Sofie, 2020) and a manually compiled gazetteer. This anonymization step is especially pertinent to the PCR task as it removes any biases in the judgment based on entity names. The ground truth labels mark all the candidate's cases relevant to each query case. Statistics for the IL-PCR corpus are shown in Table 13.

**Evaluation:** The PCR task uses micro-averaged F1@K score as the evaluation metric (as done in previous work: https://sites.ualberta.ca/~rabelo/COLIEE2021/). Prediction models predict a relevance score for each candidate for a given query. Top-K-ranked candidates are considered for prediction (i.e., whether a candidate is cited or not).

| Dataset | IL-PCR |
|---|---|
| # Documents | 7070 |
| Avg. Document Size | 8093.19 |
| # query Documents | 1182 |
| Vocab Size | 113340 |
| Total Citation Links | 8008 |
| Avg. Citation Links per query | 6.775 |
| Language | English |

Table 13: The table shows the statistics for the IL-PCR dataset (Joshi et al., 2023)

## A.7 Summarization (SUMM)

**Task Motivation and Description:** Summarization is a standard task in NLP; however, as mentioned in §1, summarizing legal documents requires legal language understanding and reasoning. **The task of summarization involves generating a gist (of a legal document) that captures the critical aspects of the case.** Summarization could be extractive (selecting the important sentences) or abstractive (generating the gist). In our setting, summarization is an abstractive generation task.

**Dataset:** For the summarization task, it is necessary to have a large dataset with gold summaries. Consequently, we use the `In-Abs` dataset (Shukla et al., 2022), created from judgment documents from the Supreme Court of India. The dataset consists of 7130 case documents with abstractive summaries (also called "headnotes"). These documents were collected from the website of the Legal Information Institute of India (http://www.liiofindia.org/in/cases/cen/INSC/), which provides free and non-profit access to databases of Indian law. These documents are accompanied by additional notes called "headnotes", which enumerate the important issues and aspects of the case. Legal experts write these headnotes and can be considered abstractive summaries of the entire case document. Headnotes usually occur in the top part of the document, just below the document header (which contains party names, date, bench, etc.), and just above the main judgment. They are also usually preceded by the heading "HEADNOTE:". Shukla et al. (2022) used these cues, and additionally employed regular expression matching to extract the headnotes from the judgment. Table 14 provides some statistics of this dataset (more details in Shukla et al. (2022)).

| Dataset | In-Abs |
|---|---|
| # Documents | 7,130 |
| Type of Summary | Abstractive |
| Language | English |
| Train/Test Split | 7,030/100 |
| Avg. Document size (in #words) | 4376.98 |
| Avg. Summary size (in #words) | 842.52 |
| Avg. Compression Ratio | 0.235 |

Table 14: The table shows the statistics of the In-Abs dataset (Shukla et al., 2022)

**Evaluation:** Following Shukla et al. (2022), we use standard summarization metrics such as ROUGE-1, ROUGE-2, and ROUGE-L F1-scores (computed

using `https://pypi.org/project/py-rouge/`, with *max_n* set to 2, parameters *limit_length* and *length_limit* not used, and other parameters kept as default), and BertScore (Zhang et al., 2019) (computed using `https://pypi.org/project/bert-score/`, version 0.3.4) that calculates the semantic similarity scores using the pre-trained BERT model.

## A.8 Machine Translation (MT)

**Task Motivation and Description:** In the Indian legal setting, when a case is transferred (due to re-appeal) from a district court to a High court, the corresponding document (typically in a regional language) needs to be translated to English. Additionally, since a large majority of the Indian population is not proficient in English, High Court / Supreme Court documents often need to be translated from English to Indian languages for a better understanding of the involved parties. In both scenarios, such translations, if done by humans, become a primary reason for delay in administering justice. Machine translation (MT) can augment human translators who could post-edit the translated document rather than translating from scratch. As outlined in §1, legal documents have different lexicons and styles; hence, existing MT systems do not perform well (Mahapatra et al., 2023). Given that many Indian languages are low-resource, MT becomes even more challenging, requiring specialized models for translating legal documents in low-resource Indian languages. **The task of Legal Machine Translation (L-MT) is to translate text in English to Indian languages and vice-versa.**

**Dataset:** For this task, we use the Multilingual Indian Legal Parallel Corpora (MILPaC) (Mahapatra et al., 2023), which comprises of a total of 17,853 parallel text pairs across English and 9 Indian languages, namely, Bengali (BN), Hindi (HI), Gujarati (GU), Malayalam (ML), Marathi (MR), Telugu (TE), Tamil (TA), Punjabi (PA) and Oriya (OR). MILPaC consists of following 3 datasets:

**1) MILPaC-IP:** Developed from a set of primers released by a society of law practitioners, this contains a set of approximately 57 question-answer pairs related to Indian Intellectual Property Laws, developed in EN and 9 Indian languages –BN, HI, MR, TA, GU, TE, ML, PA, OR. The details of the dataset are shown in Table 15

**2) MILPaC-CCI-FAQ:** is developed from a set of QA booklets released by the Competition Commission of India and contains 184 QA pairs on

|      | EN | BN | HI | MR | TA | TE | ML | PA | OR | GU |
|------|----|----|----|----|----|----|----|----|----|----|
| **EN** | × | 110 | 114 | 114 | 114 | 112 | 114 | 114 | 114 | 114 |
| **BN** | *365* | × | 110 | 110 | 110 | 108 | 110 | 110 | 110 | 110 |
| **HI** | *365* | *365* | × | 114 | 114 | 112 | 114 | 114 | 114 | 114 |
| **MR** | *365* | *365* | *365* | × | 114 | 112 | 114 | 114 | 114 | 114 |
| **TA** | *365* | *365* | *365* | *365* | × | 112 | 114 | 114 | 114 | 114 |
| **TE** |  |  |  |  |  | × | 112 | 112 | 112 | 112 |
| **ML** |  |  |  |  |  |  | × | 114 | 114 | 114 |
| **PA** |  |  |  |  |  |  |  | × | 114 | 114 |
| **OR** |  |  |  |  |  |  |  |  | × | 114 |
| **GU** |  |  |  |  |  |  |  |  |  | × |

Table 15: Number of parallel text units per language pair in (1) **MILPaC-IP** - black entries in upper triangular part, and (2) **MILPaC-CCI-FAQ** - blue italicized entries in lower triangular part. For both datasets, text units are QA-pairs, hence not tokenized into sentences (details in text).

|      | EN | BN | HI | MR | TA | TE | ML | PA | OR | GU |
|------|----|----|----|----|----|----|----|----|----|----|
| **EN** | × | 739 | 706 | 578 | 418 | 319 | 443 | 261 | 256 | 316 |
| **BN** |  | × | 439 | 439 | × | 319 | 438 | × | × | × |
| **HI** |  |  | × | 578 | × | 319 | 443 | 262 | 256 | × |
| **MR** |  |  |  | × | × | 319 | 443 | 133 | 128 | × |
| **TA** |  |  |  |  | × | × | × | × | × | × |
| **TE** |  |  |  |  |  | × | 319 | × | × | × |
| **ML** |  |  |  |  |  |  | × | × | × | × |
| **PA** |  |  |  |  |  |  |  | × | 256 | × |
| **OR** |  |  |  |  |  |  |  |  | × | × |
| **GU** |  |  |  |  |  |  |  |  |  | × |

Table 16: Number of Parallel Text units per language pair in **MILPaC-Acts**. Text units are tokenized into sentences for this dataset.

statutory rules based on competition issues in India. The parallel corpus has been developed for EN and 4 Indian languages — BN, HI, MA, and TA (see Table 15).

**3) MILPaC-Acts:** has been developed from 10 popular Indian Acts (statutory documents outlining laws of the country), for which official translations (from the Indian legislature) were available in English and the 9 Indian languages used in MILPaC-IP. For details, see Table 16.

The exact number of pairwise samples are shown in Table 15 (MILPaC-IP and MILPaC-CCI-FAQ) and Table 16. For more details regarding the creation and curation of the dataset, refer to Mahapatra et al. (2023).

**Evaluation:** Following the evaluation strategies proposed by Mahapatra et al. (2023), we use the standard metrics for machine translation, such as BLEU (Bi-Lingual Evaluation Understudy), GLEU (Google BLEU) and chrF++. For all metrics, the IndicNLP tokenizer is first used to tokenize the texts in Indian languages. For BLEU and

chrF++, we use the *SacreBLEU* package (`https://pypi.org/project/sacrebleu/`). In chrF++ calculation, the default order of character and word n-grams are set to 6 and 2 respectively. For GLEU, we use the *Huggingface evaluate* library for computation, and consider subsequences containing 1,2,3 and 4 tokens (`https://huggingface.co/spaces/evaluate-metric/google_bleu`).

### A.9 Limitations: Model Explainability

Model explainability is essential for the legal domain. For some tasks like L-NER, RR, SUMM, and L-MT, based on our discussions with legal practitioners, explainability may not be required as output can easily be observed and interpreted. For tasks like CJPE, BAIL, LSI, and PCR, it is indeed important to know on what basis the model came up with a particular decision. It can be done with the help of various techniques such as occlusion and attention weights (also as done for CJPE); however, for evaluation, such analyses must be verified with legal experts. It requires a significant annotation exercise by legal experts, which is time-consuming and expensive. Nevertheless, in the future, we plan to add more explainability experiments by employing legal experts to annotate explanations for the datasets for tasks like LSI and BAIL.

### A.10 Ethical Considerations

**We do not endorse the use of the benchmark data for non-research (commercial and real-life) applications, and the primary motivation for creating the IL-TUR benchmark is to consolidate all the research happening in parallel for the Indian Legal domain.**

We took various measures to reduce bias in models trained on legal documents. For relevant tasks (RR, LSI, PCR, CJPE, and BAIL), the documents were anonymized for named entities, judge names and organization names. For example, For anonymizing the ILSI dataset, we used the NER method provided by the paper "Named Entity Recognition in Indian Court Judgments" (Kalamkar et al., 2022a). We masked entities belonging to categories like PERSON, such as PETITIONER, RESPONDENT, JUDGE, LAWYER, WITNESS, and OTHER_PERSON. We also masked off entities tagged as PROVISION or STATUTE to remove any mention of statutes from the fact text. According to Kalamkar et al. (2022a), their NER model has a macro-F1 of 91.1%. To further verify the efficacy of the above method

on the ILSI dataset, we manually inspected ten randomly selected documents from the test set. We found that over 95% of entities (belonging to the classes described above) were successfully masked. The model failed in a few cases, e.g., when there was some discrepancy in text formatting, such as no space between a name and a punctuation mark. For CJPE and BAIL tasks, we removed cases (documents) related to sensitive issues like rape and sexual violence, and named entities such as judge names were anonymized.

Legal data is inherently sensitive and requires careful handling. The automated techniques (since doing it manually is not feasible) used to anonymize the data are not perfect and can sometimes fail. This can possibly have adverse effects. For example, if the names of judges are not anonymized then it can lead to model developing certain spurious correlations (or biases) related to specific type of outcomes associated with a particular judge. Failure to anonymize certain person names (or religion names) can lead to a model developing spurious correlations between certain types of crime and certain religious communities (since certain names are more prevalent in some particular religious communities). Similar things have also been observed in COMPAS system[1] for recidivism in the US, where it was biased against certain communities and gender.[2] Since Legal-NLP is a relatively new area, to the best of our abilities, we have taken all steps concerning ethical considerations and privacy. Via these tasks, we want to encourage more research in this area so that any hidden factors that could not have been thought of beforehand are also brought to light.

## B Tasks Models, Experiments and Results

In this section, we provide details for all baseline and SOTA models used for each of the tasks. Apart from these methods, we also conduct inference experiments with LLMs across all of these tasks except PCR, which we discuss in App. C.

### B.1 Legal Named Entity Recognition (L-NER)

We perform NER based on token representations generated by BERT-based models. Since each document in the dataset does not come pre-segmented into sentences or paragraphs, we need to chunk

---

[1] `https://en.wikipedia.org/wiki/COMPAS_(software)`
[2] `https://www.propublica.org/article/how-we-analyzed-the-compas-recidivism-algorithm`

documents before passing them to BERT, as case documents easily exceed the token limits of BERT. However, unlike other tasks like text classification, we need to devise a chunking strategy to avoid splitting true NEs into different chunks. For this, we choose to chunk at the last stopword (based on NLTK's list of English stopwords), which satisfies the chunk size limit. The assumption is that these stopwords are not expected to be part of entity names.

We experiment with five different BERT encoders: (i) `bert-base-uncased` (Devlin et al., 2019), (ii) LegalBERT (Chalkidis et al., 2020), (iii) CaseLawBERT (Zheng et al., 2021), (iv) InLegalBERT (Paul et al., 2023) and (v) InCaseLawBERT (Paul et al., 2023). We applied a Conditional Random Field (CRF) on top of the BERT encoder due to the efficacy of CRFs in sequence labeling tasks.

**Hyper-parameter Settings:** We set the chunk limit to 512 tokens to maximize the input capability of BERT. We trained on a single Nvidia RTX A6000 (48 GB). We used a batch size of 40 during training and 24 during testing. The models were trained for a maximum of 20 epochs with early stopping. We used different learning rates for the different layers, viz., 3e-5 for the BERT layers and 1e-3 for the fully connected and CRF layers. We have used the PyTorch implementation of CRF provided in `https://pypi.org/project/pytorch-crf/`.

**Model Result and Analysis:** Since the dataset is small, we divide the 105 documents into three folds (by trying to maintain the class label frequency distribution across folds as much as possible). We perform 3-fold cross-validation and report the mean across folds. In addition to the strict scores, we also consider another type of scoring, called *ent-type* score (Segura-Bedmar et al., 2013). This scheme considers a match correct if the predicted label type is the same as that of the ground truth, *even if* the predicted span is not correct. Naturally, this scheme is more lenient than the strict mechanism. We report both strict and ent-type scores for all models in Table 17.

In terms of F1 scores, all the models perform relatively poorly. The L-NER dataset contains entire case documents, and evaluation is done over *every occurence* of every named entity. This means that models cannot always rely on the local context to infer the nature of an entity, and all these models are incapable of long range context modeling since the inputs are chunked before feeding

| Method | Strict | | | Ent type | | |
|---|---|---|---|---|---|---|
| | mP | mR | mF1 | mP | mR | mF1 |
| BERT | 38.95 | 41.12 | 39.59 | 47.70 | 49.99 | 48.23 |
| LegalBERT | 43.98 | 48.06 | 45.58 | 53.19 | 58.33 | 55.21 |
| CaseLawBERT | 42.68 | 43.68 | 42.45 | 52.40 | 53.48 | 52.00 |
| InLegalBERT | **47.83** | **50.33** | **48.58** | **57.45** | **60.40** | **58.30** |
| InCaseLawBERT | 45.59 | 44.59 | 44.17 | 56.38 | 54.89 | 54.41 |

Table 17: Performance of BERT-based models over the L-NER dataset. All values are macro-averaged and in terms of percentage.

to them. This could be a possible reason for the low results. For every model, the ent-type scores are around 20% higher than the strict scores, suggesting that these models also struggle to identify the NE boundaries correctly on quite a few occasions. Comparing among the models, we observe increasing performance with greater degree of domain familiarization. BERT performs the poorest, followed by LegalBERT and CaseLawBERT (which have been pre-trained on legal data from other countries). Counterparts for these models pre-trained on Indian legal text, viz., InLegalBERT and InCaseLawBERT, further outperform them.

**Label Analysis:** To further analyze the performance across different labels, we calculate the strict and ent-type F1 scores of every label of the best-performing model, InLegalBERT.

Labels like WITNESS, A.COUNSEL, and R.COUNSEL are straightforward to identify, possibly due to the presence of linguistic cues like "P.W." (abbreviated for "Prosecution Witness") and "learned counsel for the appellant/respondent" close to the entity mentions. Labels like COURT, AUTHORITY, and DATE are slightly more challenging to identify due to the large degree of variations possible in the way these entities are mentioned, e.g., "Delhi High Court" vs. "High Court of Judicature at New Delhi", or "14.06.2023" vs. "14/6/23" vs. "14th June 2023". We also observe very little difference in these classes' strict and ent-type scores.

Labels like APPELLANT, RESPONDENT, and JUDGE are more challenging to identify. There is an apparent confusion between APPELLANT and RESPONDENT roles since the entities belonging to these classes usually occur in the same context and play the same role in the court case (just opposing sides). However, the performance of the JUDGE class is lower, although JUDGE type entities are usually enclosed by prefixes such as "Hon-

| Label | APP | RESP | JUD | AC | RC | CRT | AUTH | WIT | STAT | PREC | DATE | CN | Macro |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Strict F1 | 22.72 | 11.70 | 57.33 | 61.27 | 53.32 | 69.16 | 44.37 | 29.32 | 63.45 | 36.99 | 81.52 | 51.76 | 48.58 |
| Ent-type F1 | 34.14 | 18.01 | 71.30 | 67.18 | 58.80 | 76.09 | 50.13 | 34.21 | 72.43 | 64.49 | 85.06 | 67.73 | 58.30 |

Table 18: Performance of the best model (InLegalBERT) across all labels of the L-NER dataset. All values are in terms of macro-F1(percentage).

ourable Justice" or suffixes such as "J.". The considerable difference in the strict and ent-type scores for the JUDGE class indicates that the model fails to detect the spans properly rather than the class.

Finally, for labels like STAT, PREC, and CASE NO., the spans can be challenging to identify even for human readers since these entities are usually long, occur in multiple forms, and can have extended suffixes. For example, STAT can either be in the full form, such as "Indian Penal Code, 1860" or its abbreviated version "I.P.C.", while PREC entities can sometimes contain the case number of the particular precedent as a suffix. The considerable differences in strict and ent-type scores of these entities also point to this possibility.

## B.2   Rhetorical Role Prediction (RR)

For the task of RR Prediction, we experiment with different approaches, such as passing each sentence individually to BERT, LegalBERT and In-LegalBERT or applying hierarchical approaches to model the entire document together, such as BiLSTM-CRF with sent2vec (Gupta et al., 2019) or BERT embeddings. Malik et al. (2022) suggest an auxiliary task, Label Shift Prediction (LSP), which aims to predict, for sentence $i$ in a document, whether the label changed from sentence $i - 1$ to $i$. This is based on the intuition that RRs tend to maintain some inertia when going from one sentence to another, and changes in RR labels are not abrupt but smooth. BERT-SC is obtained by fine-tuning BERT for the LSP task *only* over the train set of the RR dataset. Finally, the Multi-task Learning (MTL) approach incorporates both RR (main task) and LSP (auxiliary task) prediction. For more details about LSP and MTL, check Malik et al. (2022).

**Results and Analysis:** Table 19 compares the performances of different models for the RR task. It is evident that RR prediction is a challenging task; standard transformer-based models like BERT, LegalBERT and InLegalBERT applied on individual sentences do not perform well. Posing the task as a sequence labeling problem, the hierarchical models employing BiLSTM-CRF show

| Model | IT | CL | IT+CL |
|---|---|---|---|
| BERT | 0.56 | 0.52 | 0.58 |
| LegalBERT | 0.55 | 0.53 | 0.56 |
| InLegalBERT | 0.64 | 0.52 | 0.58 |
| BiLSTM-CRF (sent2vec) | 0.59 | 0.61 | 0.60 |
| BiLSTM-CRF (BERT emb) | 0.63 | 0.63 | 0.63 |
| LSP (BERT-SC) | 0.65 | 0.68 | 0.67 |
| MTL (BERT-SC) | **0.70** | **0.69** | **0.70** |

Table 19: Performance of different models on the RR dataset. All values are in terms of Macro-F1.

improvements. LSP plays a significant role in improving performance, which is seen in the performance of LSP (BERT-SC) over models that do not employ LSP. Harnessing the power of learning both RR and LSP prediction in an end-to-end setup, the MTL model performs the best. However, this is still quite far from human annotations, pointing towards significant scope for improvement.

## B.3   Court Judgment Prediction with Explanation (CJPE)

We use the `ILDC-multi` split for judgment prediction and `ILDC-expert` for explanations, out of the ILDC dataset developed by (Malik et al., 2021). Different transformer-based models (BERT, RoBERTa and XLNet, InLegalBERT) have been tried for the CJPE task. Since these models cannot accommodate large documents, one approach is to make the prediction based on a chunk of 512 tokens. The last 512 tokens are chosen since these parts of the text are likely to contain more information for guiding the final decision (Malik et al., 2021). In other settings, a hierarchical approach is adopted by chunking the entire document into chunks of 512 tokens, passing these to the transformer, and collecting the [CLS] embeddings to be fed to a high-level encoder, such as BiGRU or BiGRU coupled with attention.

For the explanation part, an occlusion method is used by Malik et al. (2021). The primary idea behind this is to mask a chunk of text and then see the change in prediction probability. The prediction probability change indicates the salience of that particular chunk for making the prediction. The

| Model | Macro Precision (%) | Macro Recall (%) | Macro F1 (%) | Accuracy (%) |
|---|---|---|---|---|
| BERT | 69.33 | 67.31 | 68.31 | 67.24 |
| RoBERTa | 72.25 | 71.31 | 71.77 | 71.26 |
| XLNet | 72.09 | 70.07 | 71.07 | 70.01 |
| InLegalBERT | 74.13 | 73.86 | 73.76 | 72.87 |
| BERT + BiGRU | 70.98 | 70.42 | 70.69 | 70.38 |
| RoBERTa + BiGRU | 75.13 | 74.30 | 74.71 | 74.33 |
| XLNet + BiGRU | 77.80 | 77.78 | 77.79 | 77.78 |
| BERT + BiGRU-attn | 71.31 | 70.98 | 71.14 | 71.26 |
| RoBERTa + BiGRU-attn | 75.89 | 74.88 | 75.38 | 74.91 |
| XLNet + BiGRU-attn | 77.32 | 76.82 | 77.07 | 77.01 |
| LegalBERT + BiLSTM-attn | 77.73 | 77.02 | 76.90 | 77.06 |
| InLegalBERT + BiLSTM-attn | 82.15 | 81.45 | 81.31 | 81.41 |

Table 20: Performance of different models on the ILDC-multi dataset. All values are macro-averaged and in terms of percentage.

more the change in probability, the more salient the chunk.

**Results and Analysis:** From Table 20, it is evident that the hierarchical models perform better than their counterparts that take just the last 512 tokens (and thus suffer from loss of information). While adding the attn. The layer to the BiGRU module seems to help BERT and RoBERTa slightly, but the same is not true for XLNet. However, BERT-based models developed for the Indian legal domain, such as InLegalBERT (Paul et al., 2023), outperform the open-domain encoders and achieve state-of-the-art performance in terms of macro F1.

The occlusion approach for extracting explanations can give positive or negative scores to each chunk; we choose the chunks that obtain positive scores. The text from these chunks is concatenated and compared with the expert-annotated chunks (5 different annotations for 5 experts). We only consider all sentences ranked 1 to 10 by the experts as gold-standard explanations (note that many sentences are not ranked). As described by Malik et al. (2021), the occlusion scores are calculated at the chunk level from the hierarchical model and at the sentence level (for a particular chunk) from the flat model. Thereon, chunks with positive score are chosen, and among them, the top 50% sentences in terms of occlusion score are chosen for evaluation w.r.t. gold-standard. The best model, InLegalBERT + BiLSTM-Attn, gives 0.561 Rouge-L score and 0.325 BLEU score averaged across all experts. This demonstrates that explainability is still a big challenge, and the model's understanding of important sentences is quite far off from that of the experts.

| Model | Accuracy | F1 |
|---|---|---|
| IndicBert-First 512 | 0.73 | 0.71 |
| IndicBert-Last 512 | 0.78 | 0.76 |
| TF-IDF+IndicBert | **0.82** | **0.81** |
| TextRank+IndicBert | 0.82 | 0.81 |
| Salience Pred.+IndicBert | 0.80 | 0.78 |
| Multi-Task | 0.80 | 0.78 |

Table 21: Performance of different models over the HLDC dataset. F1 values are macro-averaged, and all values are in terms of percentage.

## B.4 Bail Prediction (BAIL)

We use the HLDC-all-districts (Kapoor et al., 2022) split for all our experiments. For BAIL, we used the multi-lingual IndicBERT (Kakwani et al., 2020) to encode the facts and predict. Since the facts can be long, some unsupervised summarization-based approaches (such as TF-IDF ranking and TextRank) have been tried to shorten the inputs and remove noise. We also experiment with the salience prediction approach demonstrated by Kapoor et al. (2022) that aims to predict the important sentences via supervised learning of salience scores (the gold standard scores are decided by comparing each fact sentence with the final case summary written by the judge). Finally, we also an MTL approach by combining BAIL and salience prediction tasks is also carried out.

**Results and Analysis:** The results are reported in Table 21. As we observe, summarization of the input facts is a better approach than just taking the first or last 512 tokens for passing to IndicBERT. Surprisingly, TF-IDF shows the best performance with 81% macro-F1, even outperforming supervised salience prediction and MTL approaches. This could possibly be because of the large vari-

ation in the nature and dialect of text across the entire dataset.

## B.5 Legal Statute Identification (LSI)

We chose some models from the BERT family – LegalBERT (Chalkidis et al., 2020) and InLegal-BERT (Paul et al., 2023) as baselines for this task. Since fact descriptions (input for LSI) can be long, they may not fit within the maximum 512-token limit for BERT encoders, necessitating a hierarchical model. Examples from the ILSI dataset are pre-segmented into sentences. We pass each sentence individually through the BERT encoder and gather the [CLS] embeddings for each document. The sequence of [CLS]-embeddings are passed through an upper Bi-LSTM layer coupled with attention, yielding a single representation for the entire fact portion. It then passes through a fully connected layer with sigmoid activation to obtain label probabilities. Labels with a probability score $> 0.5$ are considered relevant. Apart from these two models, we also experiment with LeSICiN (Paul et al., 2022), a graph-based deep neural model that also utilizes sent2vec (Gupta et al., 2019) embeddings pre-trained on Indian legal data.

| Encoder Module | mP | mR | mF1 |
|---|---|---|---|
| LegalBERT + LSTM-Attn | 53.79 | 15.72 | 21.74 |
| InLegalBERT + LSTM-Attn | **58.75** | 19.29 | 26.23 |
| LeSICiN | 24.34 | **36.58** | **28.08** |

Table 22: Performance over the ILSI dataset for LSI. All reported values are macro-averaged and in terms of percentage.

**Results:** The results are reported in Table 22. All models perform poorly, indicating the challenging nature of the ILSI dataset. Among the BERT-based methods, InLegalBERT outperforms LegalBERT since the former has been trained on Indian legal documents and is likely to have more inherent domain knowledge. While the BERT-based methods utilize strong contextual representations to identify patterns in the fact text that highly correlate with certain labels (high precision), the low recall suggests that the model is not able to pick up more latent patterns. On the other hand, LeSICiN shows a comparatively better recall since it compares the fact text with the text of the statutes via a graph neural network but has poor precision. Overall, LeSICiN still manages to outperform the BERT-based methods.

## B.6 Prior Case Retrieval

For the PCR task, a classical IR baseline BM-25, apart from some transformer-based approaches, is chosen. We follow the baselines proposed in (Joshi et al., 2023) and perform all the experiments, including the ones where a document is converted to a set of events.

**Results and Analysis:** The results are shown in Table 23. BM-25 seems to be a strong baseline, and BERT-based models fail to outperform this. In fact, the scores of transformer-based approaches are surprisingly low (less than 10% F1). Instead, the event-filtered doc approach works the best. Comparing the two event-based approaches, working directly with the atomic events works better for BM25 approaches with unigrams and bigrams, but for trigram onwards, the event-filtered doc approach outperforms this.

We have observed that the event-based models perform the best but still have a micro F1 score of 39.15, which is relatively low. Given the low scores, there is massive scope for developing better models for PCR.

## B.7 Summarization (SUMM)

Although the `IN-Abs` dataset is meant for abstractive summarization, we can apply both extractive and abstractive methods (Shukla et al., 2022).
(i) **Extractive methods:** We try out approaches like CaseSummarizer (Polsley et al., 2016) (legal-specific, unsupervised), DSDR (He et al., 2012) (open domain, unsupervised), Gist (Liu and Chen, 2019) (legal-specific, supervised) and SummaRuN-Ner (Nallapati et al., 2017) (open domain, supervised). To adapt the abstractive gold-standard summaries for these extractive methods, we use the technique suggested by Narayan et al. (2018).
(ii) **Fine-tuned Abstractive methods:** We try out text generation models both from the open-domain like BART (Lewis et al., 2019), and legal domain like Legal-Pegasus (HuggingFace, b) and Legal-LED (HuggingFace, a). While Legal-LED can accommodate a large number of documents (16,384 token limit), the same is not true for the other models. To overcome this problem, we chunk the document into equal-sized chunks (each chunk size is lesser than the model length limit) and pass each chunk through the model. The summaries for each chunk are concatenated to form the final summary. To convert the overall document summary (gold standard) into chunk-wise summaries, we follow

| Model | | K | Precision | Recall | F1 |
|---|---|---|---|---|---|
| Word Level | BM25 | 5 | 17.11 | 11.64 | 13.85 |
| | BM25 (Bigram) | 7 | 29.30 | 27.91 | 28.59 |
| Segmented-Doc Transformer (full document) | BERT | 6 | 10.28 | 8.40 | 9.24 |
| | BERT (finetuned) | 6 | 8.79 | 7.18 | 7.90 |
| | DistilBERT | 7 | 17.02 | 16.21 | 16.61 |
| | DistilBERT (finetuned) | 5 | 9.70 | 6.60 | 7.86 |
| | LegalBERT | 6 | 7.87 | 9.65 | 8.67 |
| | InCaseLawBERT | 11 | 3.02 | 4.52 | 3.62 |
| | InLegalBERT | 12 | 6.10 | 9.96 | 7.56 |
| Atomic Events | Jaccard Similarity | 7 | 35.12 | 33.28 | 34.17 |
| | BM25 | 7 | 37.69 | 35.90 | 36.77 |
| | BM25 (Bigram) | 6 | 35.39 | 28.89 | 31.81 |
| | BM25 (Trigram) | 6 | 30.71 | 25.07 | 27.61 |
| Events Filtered Docs | BM25 | 5 | 24.26 | 16.50 | 19.64 |
| | BM25 (Bigram) | 6 | 33.69 | 27.50 | 30.28 |
| | BM25 (Trigram) | 6 | **41.35** | 33.76 | 37.17 |
| | BM25 (Quad-gram) | 7 | 40.12 | **38.22** | **39.15** |
| | BM25 (Penta-gram) | 7 | 39.57 | 37.70 | 38.61 |

Table 23: Performance of different models on the IL-PCR dataset. The table shows the K values, and Precision, Recall, and F1 scores (in terms of percentage) for each model.

the approach given by Gidiotis and Tsoumakas (2020). All the models were fine-tuned on the train part of the IN-Abs dataset.

**Model Result and Analysis** The results of all approaches are reported in Table 24. SummaRuN-Ner performs the best among the extractive approaches across three of the four metrics considered (Rouge-1 & 2, and BERTScore). The abstractive approaches show a general improvement over the extractive ones, possibly due to the gold-standard summaries also being abstractive. Despite being open-domain and requiring chunking, the BART model still comes close to or outperforms Legal-LED across different legal domain-specific metrics and can accommodate very long documents. Legal Pegasus beats BART in terms of R-2 and R-L but falls short in terms of R-1. Legal-LED outperforms every other model in terms of BERTScore.

## B.8 Legal Machine Translation (L-MT)

For this task, we employed a host of systems, including Commercial systems such as Google Cloud Translation - Advanced Edition (v3) system[3] (**GOOG**) and the Translation API offered by Microsoft Azure Cognitive Services (v3)[4] (**MSFT**). We also used open-source models such as Indic-Trans, which is a transformer-4x based multilingual

| Algorithm | ROUGE Scores | | | BERTScore |
|---|---|---|---|---|
| | R-1 | R-2 | R-L | |
| *Extractive Methods (U: Unsupervised, S: Supervised)* | | | | |
| DSDR (U) | 0.485 | 0.222 | 0.270 | 0.848 |
| CaseSummarizer (U) | 0.454 | 0.229 | 0.279 | 0.843 |
| SummaRunner (S) | 0.493 | **0.255** | 0.274 | 0.849 |
| Gist (S) | 0.471 | 0.238 | 0.308 | 0.842 |
| *Abstractive Methods* | | | | |
| BART | **0.495** | 0.249 | 0.330 | 0.851 |
| Legal-Pegasus | 0.488 | 0.252 | **0.341** | 0.851 |
| Legal-LED | 0.471 | 0.235 | 0.332 | **0.856** |

Table 24: Document-wide ROUGE-L and BERTScores (Fscore) on the IN-Abs dataset, averaged over the 100 test documents.

---

[3] https://cloud.google.com/translate/docs/samples/translate-v3-translate-text

[4] https://azure.microsoft.com/en-us/products/cognitive-services/translator

| EN → IN | Model | MILPaC-IP | | | MILPaC-Acts | | | MILPaC-CCI-FAQ | | |
|---|---|---|---|---|---|---|---|---|---|---|
| | | BLEU | GLEU | chrF++ | BLEU | GLEU | chrF++ | BLEU | GLEU | chrF++ |
| EN → BN | GOOG | 27.7 | 30.7 | 56.8 | 12.0 | 17.0 | 40.7 | **52.0** | **53.6** | **74.8** |
| | MSFT | **31.0** | **33.8** | **59.4** | 18.4 | **23.1** | **45.6** | 36.5 | 40.4 | 66.2 |
| | IndicTrans | 24.7 | 27.3 | 51.7 | **18.6** | 21.8 | 45.5 | 20.9 | 25.6 | 50.2 |
| EN → HI | GOOG | 36.6 | 35.3 | 53.8 | 21.2 | 26.7 | 47.1 | 46.0 | 48.4 | 67.3 |
| | MSFT | **38.5** | **37.0** | **54.9** | **46.4** | **48.9** | **67.3** | 45.5 | 48.2 | **67.5** |
| | IndicTrans | 27.0 | 28.1 | 45.1 | 45.7 | 48.2 | 66.6 | **49.1** | **49.8** | 67.1 |
| EN → TA | GOOG | **39.3** | **41.8** | **69.4** | 8.1 | 13.7 | 37.0 | **41.4** | **44.0** | **70.7** |
| | MSFT | 35.3 | 38.7 | 68.8 | **12.1** | **17.6** | **46.3** | 29.5 | 33.7 | 64.9 |
| | IndicTrans | 21.4 | 25.5 | 51.9 | 11.1 | 16.7 | 43.7 | 22.9 | 26.8 | 56.1 |
| EN → MR | GOOG | **23.0** | **25.6** | **51.6** | 8.6 | 14.6 | 37.5 | **51.3** | **53.0** | **74.8** |
| | MSFT | 19.4 | 22.8 | 49.6 | **13.9** | **19.6** | **45.0** | 34.1 | 38.3 | 65.8 |
| | IndicTrans | 16.0 | 19.6 | 44.0 | 12.9 | 18.5 | 42.1 | 28.2 | 32.0 | 56.7 |
| EN → TE | GOOG | **22.4** | **23.2** | **48.9** | 6.6 | 11.4 | 28.8 | - | - | - |
| | MSFT | 15.8 | 18.3 | 44.8 | **12.0** | **16.9** | 39.4 | - | - | - |
| | IndicTrans | 15.5 | 17.6 | 40.6 | 11.9 | 16.8 | **40.4** | - | - | - |
| EN → ML | GOOG | 22.3 | 27.7 | 57.5 | 7.3 | 12.4 | 32.2 | - | - | - |
| | MSFT | **34.2** | **37.7** | **66.5** | 10.8 | 17.0 | 46.2 | - | - | - |
| | IndicTrans | 19.8 | 24.5 | 48.9 | **16.6** | **21.2** | **50.3** | - | - | - |
| EN → PA | GOOG | 17.8 | 20.8 | 41.3 | 8.9 | 14.1 | 28.6 | - | - | - |
| | MSFT | **30.2** | **30.5** | **51.3** | **40.1** | **42.4** | **62.5** | - | - | - |
| | IndicTrans | 28.1 | 28.8 | 47.6 | 24.0 | 28.8 | 48.8 | - | - | - |
| EN → OR | GOOG | 2.4 | 6.5 | 29.0 | 4.1 | 8.2 | 26.3 | - | - | - |
| | MSFT | 5.5 | 9.0 | 33.7 | 7.6 | 13.3 | 37.3 | - | - | - |
| | IndicTrans | 4.9 | 8.6 | 30.5 | 8.9 | 15.0 | **40.4** | - | - | - |
| EN → GU | GOOG | 43.6 | 46.0 | 67.8 | 14.3 | 19.5 | 42.1 | - | - | - |
| | MSFT | **47.3** | **49.2** | **70.6** | 21.7 | 26.1 | **51.9** | - | - | - |
| | IndicTrans | 31.3 | 34.9 | 56.3 | **22.9** | **27.0** | 50.9 | - | - | - |
| Average | GOOG | 26.1 | 28.6 | 47.6 | 10.1 | 15.3 | 35.6 | **47.7** | **49.8** | **71.9** |
| | MSFT | **28.6** | **30.8** | **55.5** | 20.3 | 25.0 | 49.1 | 36.4 | 40.2 | 66.1 |
| | IndicTrans | 24.4 | 27.8 | 52.5 | **21.7** | **26.8** | **53.3** | 30.3 | 33.6 | 57.5 |

Table 25: Corpus-level BLEU, GLEU, and chrF++ scores for all MT systems, over three datasets. All values are averaged over all text pairs in a particular dataset. For each dataset and each English-Indian language pair, the best value of each metric is boldfaced.

| Model | BLEU | GLEU | chrF++ |
|---|---|---|---|
| GOOG | 28.0 | 31.2 | 51.7 |
| MSFT | **28.4** | **32** | **56.9** |
| IndicTrans | 25.5 | 29.4 | 54.4 |

Table 26: Corpus-level BLEU, GLEU, and chrF++ scores for all MT systems. All values are averaged over all text pairs, across all languages, and across 3 datasets.

NMT model[5] trained over the *Samanantar* dataset for translation among Indian languages (Ramesh et al., 2022).

**Model Result and Analysis** The performances of all the MT systems across the 3 datasets are presented in Table 25. We find that no single model performs the best in all scenarios. MSFT, GOOG, and IndicTrans are the 3 best models that generally perform the best in most scenarios. The scores for **MILPaC-Acts** are consistently lower than those for other datasets. This is expected since **MILPaC-Acts** has very formal legal language, which is challenging for all MT systems. Interestingly, though MSFT and GOOG perform the best over most datasets, IndicTrans performs better over **MILPaC-Acts** for several Indian languages (e.g., Malayalam & Gujarati). The superior performance of Indic-Trans over **MILPaC-Acts** may stem partly from the fact that it was trained on some legal documents from Indian government websites (such as State Assembly discussions) according to Ramesh et al. (2022). However, it is *not* known publicly over what data commercial systems such as GOOG and MSFT are trained. By looking at the average

[5] https://github.com/AI4Bharat/indicTrans

scores across all 3 datasets and language pairs (see Table 26), we can establish that MSFT performs the best across all metrics.

## C  Additional Experiments with LLMs

The wide generalization capability of large language models has shown tremendous performance across various Natural Language Understanding (NLU) tasks. To validate if the available LLMs generalize enough to domain-specific legal language, we perform a detailed set of experiments by prompting LLMs over the set of proposed tasks in **IL-TUR**. We design prompts based on the available task, the context length, and prior knowledge required for the task, like label definition, which is specific to the legal domain. In recent years, In-Context Learning (ICL) (Brown et al., 2020) has significantly improved LLMs performance on various tasks. Considering the performance boost due to the ICL prompt template, it becomes crucial to consider few-shot prompts. For our experiments with LLMs, we design a prompt template that is compatible with ICL, i.e., the same prompt template can be used to provide a few shot examples as a prompt to the language models. Primarily, we validate the performance of large proprietary LLMs as well as smaller non-commercial LLMs. As some of the tasks require the entire document to be a part of the model's input, evaluating the entire test sets becomes more challenging and time-consuming for tasks with large test sets. Since the primary design of the benchmark is not LLM specific, we perform the LLM validation to obtain a general proxy of LLM performance.

### C.1  Experiments with Proprietary LLMs

For experiments with proprietary LLMs, we consider the widely used models released by OpenAI: GPT-3.5 (`gpt-3.5-turbo-16k`) and GPT-4 (`gpt-4-turbo`). As explained in §4, PCR requires as input the texts of the source document as well as a set of candidate documents. Due to the size of legal documents, such a setup exceeds the token length limit for GPT-3.5 and also for GPT-4 if all candidates are considered. Hence, we could not experiment with LLMs for the task of PCR. We discuss the task-specific prompt design and evaluation strategies and the obtained findings in the subsections below. Table 27 shows the results for various tasks.

### C.1.1  Legal Named Entity Recognition (L-NER)

**Prompt Design:** Although the NER task is known to GPT, LNER involves clearly understanding the meaning of the legal entities. Thus, we provide descriptions of the entities as part of our prompt (Table 28).

**Data Selection:** As discussed in App. B.1, we divided our entire data into 3 folds for testing the other models. In this experiment, we only choose the documents of one particular fold (Fold 1) for passing to GPT. For in-context learning, we randomly sample documents from Fold 2. In some cases, especially for 2-shot prompting, the input did not fit within 16k tokens (for GPT-3.5) even after choosing the shortest in-context (IC) examples. In these cases, we split the document into chunks, passed each chunk to the model along with IC examples, and collated the outputs from each chunk to produce the final output. No such adjustments were needed for GPT-4 due to its larger context length (128k tokens).

**Verbalization:** We expect the model's output to be precisely compatible with JSON. For GPT-3.5, the generated JSON format was sometimes incomplete, and we used string processing to complete these strings for JSON compatibility. For GPT-4, all results were perfectly JSON compatible.

**Results:** GPT returns a list of entities for each class. We mapped all character spans in the document corresponding to each entity and used these character span mappings to generate the BIO sequence that is further used for evaluation. The results for the GPT are mentioned in Table 29. Firstly, we observe that GPT-4 performs very poorly as compared to GPT-3.5 (discussed below). In terms of the strict scores, GPT-3.5 performs much poorly compared to the SOTA models, demonstrating that it cannot understand the legal roles clearly without any fine-tuning. Observing the 1 and 2-shot results, it is clear that providing a single ICL example can mislead the model, and adding 2 examples provides a slight improvement over 0-shot. Finally, as observed for the BERT-based models, there is a significant difference between strict and ent-type scores.

The massive drop in performance for GPT-4 requires further investigation. We experimented by lowering the temperature, but this led to even worse performance. Similarly, we tried to modify the prompt to make the model focus on covering all entities, variations, etc. But none of these techniques

| Task | GPT-3.5 | | | GPT-4 | | | SOTA | Metric |
|------|---------|---|---|-------|---|---|------|--------|
| | 0-Shot | 1-Shot | 2-Shot | 0-Shot | 1-Shot | 2-Shot | | |
| L-NER | 30.59% | 23.68% | 32.84%* | 13.65% | 10.51% | 24.03% | **48.58%** | strict mF1 |
| RR | 30.95% | 30.05% | 30.31% | 37.37% | 37.43% | 38.18% | **69.01%** | mF1 |
| CJPE | 54.17% | 51.46% | 56.74% | 68.29% | 47.26% | 60.44% | **81.31%** | mF1 |
| | 0.30 | 0.29 | 0.30 | 0.40 | 0.39 | 0.43 | **0.56** | ROUGE-L |
| | 0.08 | 0.15 | 0.113 | 0.14 | 0.16 | 0.18 | **0.32** | BLEU |
| BAIL | 51.04% | 46.35% | 61.0% | 51.46% | 56.90% | 66.67% | **81%** | mF1 |
| LSI | 21.55% | 22.61% | 21.40% | 23.99 | 22.26 | 20.53 | **28.08%** | mF1 |
| SUMM | 0.21 | 0.20 | 0.22 | 0.23 | 0.16 | 0.17 | **0.33** | ROUGE-L |
| | 0.85 | 0.84 | 0.84 | 0.85 | 0.81 | 0.81 | **0.86** | BERTScore |
| L-MT | 0.23 | 0.25 | 0.26 | 0.33 | 0.35 | **0.36** | 0.28 | BLEU |
| | 0.28 | 0.28 | 0.29 | 0.36 | 0.38 | **0.39** | 0.32 | GLEU |
| | 0.42 | 0.43 | 0.43 | 0.50 | 0.52 | 0.53 | **0.57** | chrF++ |

Table 27: Performance of Open-AI GPT-3.5 (gpt-3.5-turbo-16k) and GPT-4 (gpt-4-turbo) model on various tasks for zero-shot, one-shot and two-shot settings. The SOTA corresponds to the best performing model as given in Table 3. The best result for each task is marked in **boldface**. The best GPT-based result for each task is underlined.

```
SYSTEM_PROMPT: You are a smart and intelligent Named Entity Recognition
(NER) system. I will provide you with the definition of the entities you
need to extract and the output format. I will also provide you some examples
of the task and the document from where you should extract the entities.
USER_PROMPT: Are you clear about your role?
ASSISTANT_PROMPT: Sure, I'm ready to help you with your NER task. Please
provide me with the necessary information to get started.
INPUT_PROMPT:
Entity Definition:
1. APPELLANT: Name or abbreviation of the person(s) or organization(s)
filing an appeal/petition to a court of law.
2. RESPONDENT: Name or abbreviation of a person(s) or organization(s)
responding/defending to an appeal/petition filed against them in a court
of law.
3. JUDGE: Name of the judge/justice presiding over the case in a court
of law.
4. APPELLANT COUNSEL: Name of the lawyer representing the
appellant/petitioner in a court of law.
5. RESPONDENT COUNSEL: Name of the lawyer representing the respondent in
a court of law.
6. COURT: Name of the court of law
7. AUTHORITY: Name or abbreviation of any organization apart from a
Court, which has administrative, legal or financial authority. This also
includes regulatory and investigative agencies.
8. WITNESS: Name of a person appearing as witness or testifying to a case
in a court of law.
9. STATUTE: Name or abbreviation of a statutory law or legal article.
10. PRECEDENT: Title of a prior court case.
11. DATE: Any format of date, even in natural language.
12. CASE NUMBER: Any format of prior case number or order numbers.
Important Instructions:
1. Salutations or prefixes/suffixes like Mr., Mrs., Smt., Justice, J.,
Dr., P.W., are not part of the named entity.
Output Format:
{"APPELLANT": [list of entities present], "RESPONDENT": [list of entities
present], "JUDGE": [list of entities present], "APPELLANT COUNSEL": [list
of entities present], "RESPONDENT COUNSEL": [list of entities present],
"COURT": [list of entities present], "AUTHORITY": [list of entities
present], "WITNESS": [list of entities present], "STATUTE": [list of
entities present], "PRECEDENT": [list of entities present], "DATE": [list
of entities present], "CASE NUMBER": [list of entities present]}
DO NOT REPEAT THE SAME ENTITY NAME MULTIPLE TIMES.
If no entities are presented in any category, keep an empty list for that
category.
The above format should be a pure JSON format.
Examples:
Document 1: <In-context Document 1 goes here>
Output 1: <Gold-standard Labels for Document 1 goes here>
...
Document n+1: <Test Document goes here>
Output n+1:
```

Table 28: Prompt template for L-NER for both GPT-3.5 and GPT-4 ($n$ in-context examples)

| Method | Strict | | | Ent type | | |
|--------|--------|---|---|----------|---|---|
| | mP | mR | mF1 | mP | mR | mF1 |
| GPT-3.5 0-shot | 48.57 | 24.58 | 30.59 | 65.23 | **34.46** | **42.04** |
| GPT-3.5 1-shot | 39.08 | 18.56 | 23.68 | 56.05 | 26.73 | 34.34 |
| GPT-3.5 2-shot | **51.29** | 26.16 | 32.84 | 65.63 | 32.80 | 41.54 |
| GPT-4 0-shot | 21.44 | 10.10 | 13.65 | 22.62 | 10.64 | 14.37 |
| GPT-4 1-shot | 20.32 | 7.49 | 10.51 | 22.69 | 8.34 | 11.72 |
| GPT-4 2-shot | 47.30 | 18.26 | 24.03 | 53.50 | 20.82 | 27.30 |

Table 29: Performance of GPT-3.5 and GPT-4 over the L-NER dataset. All values are macro-averaged and in terms of percentage.

improved the performance. We manually verified the outputs of GPT-3.5 and GPT-4 in comparison with the gold standard. We observed that GPT-4 was hallucinating to a great extent, returning many new entities that are not present in the input document. It also failed to capture many true entities in the process, which explains the poor precision and recall values. To further highlight the above issue, we show the predictions of GPT-3.5 and GPT-4 in comparison with the gold-standard entities for a particular case "Babulal Badriprasad Varma vs Surat Municipal Corpn. & Ors", Supreme Court of India (2008), in Table 30. Entities not present in the case document are marked in red. We can clearly see that GPT-4 produces hallucinated outputs across all labels. While it can correctly pick up true entities as well, a large number of these hallucinated outputs lead to very poor metrics. On further inspection, we observed that all these hallucinated entities came from the same case, "Sonali Hatua Giri vs Union Of India And Others on 7 April, 2021", Calcutta High Court (Appellate Side) (2021).

| Class | Gold Standard | GPT-3.5 | GPT-4 |
|---|---|---|---|
| APP | Babulal Badriprasad Varma | Babulal Badriprasad Varma | Sonali Hatua Giri, Babulal Badriprasad Varma |
| RESP | Surat Municipal Corporation | Surat Municipal Corporation | Union of India, Surat Municipal Corporation & Ors. |
| JUDGE | V.S. Sirpurkar, S.B. Sinha, Johan Romilly M. R. | S.B. Sinha | Sabyasachi Bhattacharyya, S.B. Sinha |
| A.COUNSEL | Lalit, U.U. Lalit | Mr. U.U. Lalit | Sankar Nath Mukherjee, Niraj Gupta, Priyabrata Shah, Aditya Biswas, U.U. Lalit |
| R.COUNSEL | Prashant G. Desai | Mr. Prashant G. Desai | Y. J. Dastoor, Rudraman Bhattacharya, Anuran Samanta, Prashant G. Desai |
| WITNESS | | | |
| COURT | Supreme Court of India, Criminal Court, Apex Court, High Court of Gujarat, High Court | Supreme Court of India | Punjab and Haryana High Court, Supreme Court, High Court of Gujarat |
| AUTHORITY | Government of Gujarat, Parliament | Government of Gujarat, High Court of Gujarat, Municipal Corporation | Ministry of Home Affairs, Central/State Government, Government of Gujarat |
| STATUTE | Gujarat Town Planning and Urban Development Rules, 1979; Transfer of Property Act; Bombay Town Planning Rules, 1955; Real Property Act; Bangalore Development Authority Act, 1976; Indian Penal Code; Bombay Provincial Municipalities Act; Gujarat Town Planning and Urban Development Act | Gujarat Town Planning and Urban Development Act, 1976; Gujarat Town Planning and Urban Development Rules, 1979; Bombay Provincial Municipalities Act; Indian Penal Code | Central Samman Pension Scheme; Constitution of India; Hindu Succession Act, 1956; Indian Succession Act, 1925; Hindu Marriage Act; Indian Divorce Act, 1969; Muslim Women (Protection of Rights on Divorce) Act, 1986; Parsi Marriage and Divorce Act, 1936; Gujarat Town Planning and Urban Development Act, 1976; Gujarat Town Planning and Urban Development Rules, 1979 |
| PRECEDENT | Maneklal Chhotalal & Ors. v. M.G. Makwana & Ors. [(1967) 3 SCR 65]; Maneklal Chhotalal (supra); Bhikhubhai Vithlabhai Patel & Ors. v. State of Gujarat & Anr. 2008 (4) SCALE 278; Manak Lal v. Dr. Prem Chand [AIR 1957 SC 425]; Director of Inspection of Income Tax (Investigation), New Delhi and Another v. Pooran Mal & Sons and Another [(1975) 4 SCC 568]; Bank of India v. O.P. Swarnakar (2003) 2 SCC 721, Mansukhlal (supra); Jaswantsingh Mathurasingh and Another v. Ahmedabad Municipal Corporation and Others [1992 Supp (1) SCC 5]; State of Gujarat v. Shantilal Mangaldas & Ors. [1969 (3) SCR 341]; Sureshchandra C. Mehta v. State of Karnataka and Others 1994 Supp (2) SCC 511; Mansukhlal Jadavji Darji (supra); West Bengal Housing Board etc. v. Brijendra Prasad Gupta and Others, etc. [AIR 1997 SC 2745]; Phillips v. Martin; Mansukhlal Jadavji Darji and Others v. Ahmedabad Municipal Corporation and Others [(1992) 1 SCC 384]; Vyvyan v. Vyvyan [(1861) 30 Beav. 65, 74; 54 E.R. 813, 817]; Jaswantsingh Mathurasingh (supra); Ramdev Food Products Pvt. Ltd. v. Arvindbhai Rambhai Patel and Ors. [2006 (8) SCALE 631]; Krishna Bahadur v. Purna Theatre [(2004) 8 SCC 229]; Wilson v. McIntosh | Mansukhlal Jadavji Darji v. Ahmedabad Municipal Corporation; Jaswantsingh Mathurasingh v. Ahmedabad Municipal Corporation; State of Gujarat v. Shantilal Mangaldas & Ors. | Municipal committee Patiala Vs. Model town Residents Association; Khajani Devi Vs. Union of India and others; Tulsi Devi Vs. Union of India and another; Mansukhlal Jadavji Darji and Others v. Ahmedabad Municipal Corporation and Others; Jaswantsingh Mathurasingh and Another v. Ahmedabad Municipal Corporation and Others |
| DATE | 27.12.2006; 1.06.1999; 23.11.2006; May 02, 2008; 31.03.2000; 15.01.2000; 1.7.1999; 15.1.2000 | 27.12.2006; 23.11.2006; 1.06.1999; 15.01.2000; 31.03.2000 | 07.04.2021; December 4, 2012; March 19, 1999; February 18, 2019; July 29, 2016; September 27, 2019; July 18, 2019; May 28, 2020; 27.12.2006, 23.11.2006, 1.06.1999, 15.01.2000, 31.03.2000 |
| CASE NO. | CIVIL APPEAL NO. 3203 OF 2008; SCA No. 7092 of 2001; Letters Patent Appeal No. 1611 of 2006; SLP (Civil) No. 568 of 2007 | CIVIL APPEAL NO. 3203 OF 2008; SLP (Civil) No. 568 of 2007; Letters Patent Appeal No. 1611 of 2006; SCA No. 7092 of 2001 | WPA 13806 of 2019; CWP No.1504 of 2019; No.17706 of 2017; CIVIL APPEAL NO. 3203 OF 2008, SLP (Civil) No. 568 of 2007, Letters Patent Appeal No. 1611 of 2006, SCA No. 7092 of 2001 |

Table 30: Predictions of GPT-3.5 and GPT-4 for the LNER task, compared with the gold standard. Entities not present in the original case document are marked in red.

## C.1.2 Rhetorical Role Labeling (RR)

**Prompt Design:** The RR task can be considered a semantic role labeling task over the sentences. Such a variant of the task and the definition of the rhetorical roles themselves are probably not clearly known to the GPT models; hence, we give explicit guidelines on how to carry out the labeling task. We tried out some initial prompts considering document-level inputs, i.e., passing the entire document (list of sentences) to GPT-3.5 and asking it to generate a list of labels corresponding to each sentence. This approach had several challenges, such as the output not having the same number of labels as input sentences, random token generation, etc. This problem became more pronounced in the ICL setting. Further, input text and sample output for IC examples were becoming too long. Thus, for both GPT-3.5 and GPT-4, we frame the task as a simple sentence classification task, asking the models to predict the label of an individual sentence. The final prompt is shown in Table 31. We run both GPT-3.5 and GPT-4 over all sentences in a document to get all corresponding label predictions.

```
SYSTEM_PROMPT: You are a smart and intelligent legal semantic role
labeling system. In Indian Court judgment documents, each document
sentence can be assigned a legal semantic role. Your task is, given
a sentence from an Indian Court case document, to identify the given
sentence's semantic role. I will provide you with the descriptions
of the legal semantic roles. I will also provide you with some
examples.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Absolutely, I understand my role. You would like
me to identify a sentence's legal semantic role label in an Indian
court case document. Please provide me with the descriptions of the
legal semantic roles to help guide me in accurately assigning the
role to the given sentence.

INPUT_PROMPT:
Legal Semantic Role Descriptions:
1. Fact: The actual facts and events that led to the case.
2. Argument: Legal arguments which have been put forward by either
lawyer.
3. RulingByLowerCourt: Decisions of the lower courts, if any.
4. Statute: References or citations to statutory laws and articles
referred in the case.
5. Precedent: Sentences containing References or citations to
precedents (prior cases).
6. RatioOfTheDecision: The reasoning which has been established by
the judge in the current judgment.
7. RulingByPresentCourt: The final decision of the current court.
ANSWER ONLY WITH ONE OF THE ABOVE CHOICES, DO NOT PROVIDE ANY EXTRA
OUTPUT.
Examples:
Sentence 1: <In-context Sentence 1 goes here>
Output 1: <Gold-standard Label for Sentence 1 goes here>
...
Sentence n+1: <Test Sentence goes here>
Output n+1:
```

Table 31: Prompt template for RR (for $n$ in-context examples) for both GPT-3.5 and GPT-4.

**Data Selection:** We used all sentences from all documents in the CL and IT test sets (5 documents each). For in-context samples, we randomly choose sentences from all these documents except the document from which the test sentence (sentence for which we expect GPT prediction) is sampled.

| Model | CL | IT | CL + IT |
|---|---|---|---|
| GPT-3.5 0-shot | 0.25 | 0.37 | 0.31 |
| GPT-3.5 1-shot | 0.24 | 0.36 | 0.30 |
| GPT-3.5 2-shot | 0.23 | 0.38 | 0.30 |
| GPT-4 0-shot | **0.29** | 0.46 | 0.37 |
| GPT-4 1-shot | **0.29** | 0.45 | 0.37 |
| GPT-4 2-shot | 0.28 | **0.49** | **0.38** |

Table 32: Macro-F1 scores for RR datasets

**Verbalizer:** In most cases, GPT-3.5 answers with the exact label name. In some cases, it can be accompanied by extra erroneous words. In case the prediction is a sequence of words, we iterate over the words and choose the first word that corresponds to an RR. If no such word is found, GPT-3.5 prediction has failed, and we randomly choose a label to substitute its decision. We did not observe such anomalies for GPT-4.

**Results:** The SOTA model achieves a macro-F1 of 70% over the combined (IT + CL) test set. In comparison, GPT-3.5 can only achieve a macro-F1 of 31%, showing that it is not straightforward for the LLM to assign semantic labels to sentences. On manual inspection, we observed that the model was prone to assign the FAC label to all sentences with the model temperature set to 0. On increasing the temperature to 0.95 (temperature 1 was not giving stable results), we observe that the model is still prone to assigning labels like FAC, ARG-P, ARG-R, and RPC (frequent labels) to most sentences. GPT-4 consistently outperforms GPT-3.5, with the improvements being more significant for IT documents. Also, it seems that ICL has no positive impact on either GPT-3.5 or GPT-4, with there being minimal or no improvements at all. It could be possible that just the description of the labels is not enough; GPT models might need 1/2 examples from each class to clearly understand the meaning of the RRs. However, this approach is likely to increase the context length significantly.

## C.1.3 Court Judgment Prediction and Explanation (CJPE)

**Prompt Design:** For the prediction aspect of this task, we ask GPT to read the content of the entire document and predict the final "accept"/"reject" decision (Table 33). For the explanation aspect, we modify the prompt, asking GPT first to predict the accept/reject decision and then extract important sentences of the text that led to its decision (Table 35). While the exact same prompt is used

for both GPT-3.5 and GPT-4 for the classification task, slightly different prompts were used for the explainability task for each model.

```
SYSTEM_PROMPT: You are a smart and intelligent system, trained to
act like a judge in the Indian Supreme Court. A court case document
in the Indian Supreme Court can consist of one or more appeals by
a particular party. Your task is, given such a case document, to
predict whether the appeals will be accepted or rejected. For cases
containing multiple appeals, you will predict either 'accept' if at
least one of the appeals can be accepted or 'reject' if none of the
appeals can be accepted. PLEASE ANSWER ONLY WITH EITHER 'ACCEPT'
OR 'REJECT'. I will provide you with some examples of this task and
the case document you need to make the prediction for.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Sure, I'm ready to help you with your court
judgment prediction task. Please provide me with the examples
and the case document I'm supposed to make the prediction for.

INPUT_PROMPT:
Examples:
Case Document 1: <In-context Document 1 goes here>
Output 1: <Gold-standard Label for Document 1 goes here>
...
Case Document n+1: <Test Document goes here>
Output:
```

Table 33: Prompt template for CJPE Prediction for both GPT-3.5 and GPT-4 ($n$ in-context examples)

**Data Selection:** For prediction, we divide the `ILDC-multi` test set into positive and negative examples and randomly sample 50 positive and 50 negative examples. For ICL, we randomly sample examples from the remaining test set documents such that the final prompt is within the GPT token limit. For explanation, we use all 56 documents from `ILDC-expert` for prompting. We sample the IC examples from this set itself. The gold standard outputs, in this case, are the important sentences with rank 1 and 2, as per the ranking given by either expert 3 or expert 4, chosen randomly (since these experts had the highest agreement according to Malik et al. (2021)). In both cases, for 2-shot prompting, we sample one document each from the positive and negative classes.

**Verbalizer:** For classification, the model always answers with either ACCEPT/REJECT. For explanation, in the case of GPT-3.5, we did not issue strict output format instructions since GPT-3.5 was unable to understand them correctly (based on few documents of the validation set). Thus, there were a few variations in the output format, but included a list of the important sentences, marked either with bullet points, numbering, or other delimiters. We used these cues to extract the exact sentences. However, GPT-4 can understand and adhere to complex instructions much better and thus we could specify stricter rules regarding the output format, making the task of verbalization easier for GPT-4.

**Results:** For prediction, both GPT-3.5 and GPT-4 tend to predict "reject" in favor of "accept" in

```
SYSTEM_PROMPT: You are a smart and intelligent system, trained to
act like a judge in the Indian Supreme Court. A court case document
in the Indian Supreme Court can consist of one or more appeals
by a particular party. Your task is, given such a case document,
to predict whether the appeals will be ACCEPTED or REJECTED. You
will also have to explain your prediction by QUOTING VERBATIM the
important sentences of the input text that led to your decision.
I will provide you with examples of the task, and then the input
document.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Sure, I'm ready to help you with your court
judgment prediction task. I will also quote verbatim the important
areas of the input text that led to my prediction. Please provide
me with the case document I'm supposed to make the prediction for.

INPUT_PROMPT:
IMPORTANT: For explaining your prediction, quote important
sentences verbatim from the input text. DO NOT PARAPHRASE OR
SUMMARIZE THESE SENTENCES.
Examples:
Case Document 1: <In-context Document 1 goes here>
Output 1: The appeals will be <Gold-standard ACCEPT/REJECT Label
for Document 1 goes here>. Here are the verbatim sentences that led
to this decision: <Gold-standard important sentences for Document
1 goes here>
...
Case Document n+1: <Test Document goes here>
Output:
```

Table 34: Prompt template for CJPE Explantion for GPT-3.5 ($n$ in-context examples)

```
SYSTEM_PROMPT: You are a smart and intelligent system, trained to
act like a judge in the Indian Supreme Court. A court case document
in the Indian Supreme Court can consist of one or more appeals
by a particular party. Your task is, given such a case document
from the Indian Supreme Court, predict whether the appeals will be
ACCEPTED or REJECTED. You will also have to explain your prediction
by printing the important sentences of the input text that primarily
influenced your decision. Make sure to include any sentence which is
even slightly relevant for your final decision. Do not paraphrase,
summarize or change the sentences in any way while printing, you
must quote the sentences verbatim from the input case document. I
will provide you with the output format and the input case document.
I will also provide some examples of the task to help you learn
from.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Sure, I'm ready to help you with your court
judgment prediction task. I will also quote verbatim the important
sentences of the input text that led to my prediction. Please
provide me with the example and the case document.

INPUT_PROMPT:
IMPORTANT: Strictly adhere to the output format given below. Quote
the important sentences verbatim from the input document.
Output Format:
DECISION: ACCEPTED/REJECTED
Important Sentences:
1. First important sentence
2. Second important sentence
...
Examples:
Case Document 1: <In-context Document 1 goes here>
Output 1:
DECISION: <Gold-standard ACCEPT/REJECT Label for Document 1 goes
here>
Important Sentences:
1. <Gold-standard First Important Sentence of Document 1 goes here>
2. <Gold-standard Second Important Sentence of Document 1 goes
here>
...
Case Document n+1: <Test Document goes here>
Output n+1:
```

Table 35: Prompt template for CJPE Explanation for GPT-4 ($n$ in-context examples)

| Model | ILDC$_{\text{multi}}$ | | | ILDC$_{\text{expert}}$ | |
|---|---|---|---|---|---|
| | **mP** | **mR** | **mF1** | **R-L** | **BLEU** |
| GPT-3.5 0-shot | 57.14 | 56.00 | 54.17 | 0.301 | 0.077 |
| GPT-3.5 1-shot | 57.06 | 55.00 | 51.46 | 0.292 | 0.155 |
| GPT-3.5 2-shot | 65.79 | 60.42 | 56.74 | 0.299 | 0.113 |
| GPT-4 0-shot | 70.88 | **69.00** | **68.29** | 0.398 | 0.137 |
| GPT-4 1-shot | 55.31 | 53.00 | 47.26 | 0.395 | 0.161 |
| GPT-4 2-shot | **71.88** | 64.00 | 60.44 | **0.426** | **0.184** |

Table 36: Performance over the CJPE datatsets. P, R and F1 values are macro-averaged and in terms of percentage.

most cases. Only by tweaking the temperature up to as high as $0.98$, we could observe more "accept" predictions. Despite this, GPT-3.5 significantly underperforms compared to SOTA approaches, barely performing better than random choice (see Table 36). The result turns even worse with 1-shot prompting, possibly making the model biased towards the class of the IC example. 2-shot prompting gives the best result among these settings for GPT-3.5. GPT-4, however, produces quite decent performance in 0-shot setting. It seems that the addition of ICL examples is greatly detrimental to GPT-4, producing the worst performance in 1-shot setting, while also showing a drastic decrease in performance for the 2-shot setting.

The explanation is a more difficult task than the prediction, and GPT-3.5 again underperforms compared to the SOTA approach, especially considering the BLEU score (Table 36). ICL does not change the R-L score but has a positive impact on the BLEU score in the 1-shot setting only. GPT-4 shows significant improvement over GPT-3.5, both in terms of R-L and BLEU scores. ICL improves the R-L score under 2-shot setting only, while BLEU score improves progressively from 0-shot to 2-shot. The superior understanding capability of GPT-4 helps it perform better for both prediction and explanation, as compared to GPT-3.5.

### C.1.4 Bail Prediction (BAIL)

**Prompt Design:** BAIL is a binary classification task, and in terms of understanding and format, it is very similar to the CJPE task, the only difference being that the HLDC dataset for BAIL contains Hindi text rather than English. We use the same prompt for both GPT-3.5 and GPT-4, asking the models to read the application's content and

provide the final decision, i.e., if the bail will be granted or dismissed (see Table 37).

```
SYSTEM_PROMPT: You are a smart and intelligent system, trained to
act like a judge in a district court of India. Most criminal cases
in district courts involve bail applications written in Hindi. The
application can be 'granted' if the judge believes the applicant
deserves relief or 'dismissed' if the crime is too grave to grant
relief. Your task is, given such a bail application, to predict
if the bail will be 'granted' or 'dismissed'. PLEASE ANSWER ONLY
WITH EITHER 'GRANTED' OR 'DISMISSED'. I will provide you with some
examples of this task and the application document you need to
make the prediction for.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Sure, I'm ready to help you with your bail
application prediction task. Please provide me with the examples
and the bail application I'm supposed to make the prediction for.

INPUT_PROMPT:
Examples:
Bail Application 1: <In-context Application 1 goes here>
Output 1: <Gold-standard Label for Application 1 goes here>
...
Bail Application n+1: <Test Application goes here>
Output:
```

Table 37: Prompt template for BAIL Prediction for both GPT-3.5 and GPT-4 ($n$ in-context examples)

**Data Selection:** We divide the `HLDC-all-districts` test set into positive and negative examples and randomly sample 50 positive and 50 negative examples that can be accommodated in the token length limit. For ICL, we sample examples at random from the rest of the test set. For the 2-shot setting, we always sample one example each from the positive and negative classes.

**Verbalizer:** Both GPT-3.5 and GPT-4 outputs only GRANTED/DISMISSED, so we directly take the model output as the predicted label.

**Results:** Both GPT-3.5 and GPT-4 perform poorly in the 0-shot setting. As with CJPE, we observed a much higher proportion of negative class predictions. However, unlike CJPE, adjusting the temperature did not help too much. 1-shot ICL degrades the performance for GPT-3.5, possibly biasing the model to the class of the ICL example. GPT-4 is able to overcome that bias, showing improvements. Both models perform significantly better with 2

| Model | mP | mR | mF1 |
|---|---|---|---|
| GPT-3.5 0-shot | 52.22 | 52.00 | 50.74 |
| GPT-3.5 1-shot | 46.85 | 47.00 | 46.35 |
| GPT-3.5 2-shot | 63.37 | 62.00 | 61.00 |
| GPT-4 0-shot | 57.06 | 55.00 | 51.46 |
| GPT-4 1-shot | 58.91 | 58.00 | 56.90 |
| GPT-4 2-shot | **71.43** | **68.00** | **66.67** |

Table 38: Performance over the HLDC dataset for BAIL. P, R and F1 values are macro-averaged and in terms of percentage.

ICL examples, with GPT-4 performing the best by a margin.

### C.1.5 Legal Statute Identification (LSI)

**Prompt Design:** The Indian Penal Code (IPC) is already known to GPT models since both GPT-3.5 and GPT-4 can accurately answer when asked about the content of different Sections of IPC. In an initial setting, we asked both GPT-3.5 and GPT-4 to just output the list of relevant Section numbers of IPC for a given input. We observed that the models produced hallucinated or completely non-relevant outputs; in this case, the output for GPT-3.5 often consisted of *non-existent* IPC Section numbers. Now, each IPC Section contains a corresponding title, which is a very short description of the entire statute. In a second setting, we asked both models to output the section numbers and their corresponding titles. For instance, if, for a particular case, Section 302 of the IPC is relevant, the model was expected to output just "Section 302" in the first setting, whereas it was expected to answer "Section 302 — Punishment for murder" in the second setting. We observed that this second setting reduced the hallucination to a great extent and improved performance. We did not specify any strict output format for GPT-3.5 since it was causing wrong predictions. However, we used this strategy for GPT-4, asking the model to print the relevant statutes in a structured format. The prompts are shown in Table 39 (GPT-3.5) and Table 40.

```
SYSTEM_PROMPT: You are an intelligent Legal Crime Classification
system. In the Indian legal system, the Indian Penal Code (IPC) is
an Act in the Indian legislature that contains many legal articles
or 'Sections' that codify different laws. Your task is, given the
facts or evidence of an Indian court case as input, to predict the
relevant or violated 'Sections' of the IPC as output. I will provide
you some examples of this task and the facts of the case to make
predictions for.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Yes, I understand my role as an intelligent Legal
Crime Classification system for the Indian legal system. You can
provide me with the facts of a court case, and I will identify
the relevant or violated sections of the Indian Penal Code (IPC)
based on the provided input and output format. Please go ahead and
provide me with the examples and the necessary information for the
case you'd like me to analyze.

INPUT_PROMPT:
Output Format:
List of relevant Sections and their titles
Examples:
Facts 1: <In-context Facts 1 go here>
Output 1: <Gold-standard Labels for Facts 1 go here>
...
Facts n+1: <Test Facts go here>
Output:
```

Table 39: Prompt template for LSI for GPT-3.5 ($n$ in-context examples)

**Data Selection:** We randomly chose 100 documents (in this case, fact portions) from the ILSI

```
SYSTEM_PROMPT: You are an intelligent Legal Statute Identification
(LSI) system. You will be provided the the facts of an Indian court
document. You need to output the Sections of the Indian Penal Code
(IPC) and their corresponding titles which are possibly violated
given the facts of the document. You should strictly adhere to the
output format. Do not output anything else. It has been found out
that on average, each Indian Court Document has may contain between
1 and 12 relevant statutes per document. Keep this in mind while
finding out relevant statutes.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Yes, I understand my role as an intelligent Legal
Statute Identification (LSI) system. My task is to identify and
output the IPC sections and their corresponding titles that may be
violated provided the facts of the document. I will adhere strictly
to the specified output format. Please provide me with the facts,
and I will do my best to assist you.

INPUT_PROMPT:
Output Format:
<Statute 1>
<Statute 2>
...
<Statute x>
Examples:
Facts 1: <In-context Facts 1 go here>
Output 1: <Gold-standard Labels for Facts 1 go here>
...
Facts n+1: <Test Facts go here>
Output:
```

Table 40: Prompt template for LSI for GPT-4 ($n$ in-context examples)

| Model | mP | mR | mF1 |
|---|---|---|---|
| GPT-3.5 0-shot | 21.60 | **32.55** | 21.55 |
| GPT-3.5 1-shot | **27.06** | 22.07 | **22.61** |
| GPT-3.5 2-shot | 25.35 | 21.53 | 21.40 |
| GPT-4 0-shot | 25.31 | 26.74 | 23.99 |
| GPT-4 1-shot | 27.13 | 23.22 | 22.26 |
| GPT-4 2-shot | 25.16 | 20.89 | 20.53 |

Table 41: Macro-averaged scores for ILSI dataset

test set, all of which satisfied the length constraints of GPT. For ICL, we sample other documents from the test set while satisfying the length constraints. Also, for IC examples, we collate the gold-standard Section numbers and their respective titles in the form *Section x — Title of Section x*, create a numbered list, and pass it to GPT.

**Verbalizer:** Due to the flexibility of the output format for GPT-3.5, it can output a lot of Sections from the IPC and even other acts. We filtered the outputs by considering if either the Section number OR the Section title matched with any of the 100 IPC Section numbers and the corresponding titles of the ILSI candidate statute set. The OR condition was necessary since we observed that even with the second setting, GPT still suffers from the hallucination problem, sometimes providing the correct Section titles with non-existent Section numbers. For instance, consider the GPT output *"Section 1565 of the Indian Penal Code (IPC) - Liability of abettor when one act abetted and different act done"*. This is a hallucinated output since IPC does

not have more than 600 Sections. But, the title actually corresponds to a Section in IPC, namely Section 111. Although the output format for GPT-4 was stricter, we still had to perform these processing steps so that we could match with either the Section number or the title.

**Results:** The ILSI dataset is quite challenging, as seen in the SOTA results. The results for the GPT models are listed in Table 41. In such a comparison, GPT-3.5 does not perform too badly, as compared to other tasks. GPT-4 improves upon this score. ICL does not seem to help too much, with 0,1 and 2-shot settings showing very little difference in results for GPT-3.5. However, the performance of GPT-4 actually decreases significantly with ICL, even performing worse than some GPT-3.5 settings.

### C.1.6 Summarization (SUMM)

```
SYSTEM_PROMPT: You are a smart and intelligent summarization system,
trained to read, understand and summarize Indian court case documents.
Your task is, given a court case document and a target summary length,
generate a detailed summary of the case in your own words within
the specified length. The summary should contain ALL the important
legal aspects of the case. I will provide you with the document to
be summarized.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Sure, I'm ready to help you with your court
judgment summarization task. Please provide me with the examples
and the case document I'm supposed to summarize.

INPUT_PROMPT:
Output Format: Generate the summary in a few simple paragraphs. Do
not use any paragraph headers, bullet points, or any other such
formatting.
Examples:
Case Document 1: <In-context Document 1 goes here>
Summary 1 (in <Length of reference summary goes here> words):
<Reference summary for Document 1 goes here>
...
Case Document n+1: <Test Document goes here>
Summary n+1 (in <Length of reference summary goes here> words):
```

Table 42: Prompt template for SUMM for both GPT-3.5 and GPT-4 ($n$ in-context examples)

**Prompt Design:** GPT is known to be more conversant with the abstractive summarization task. Hence, we provide the model with the summary length limit and ask it to generate the summary (see Table 42). A large majority of the judgments (more than 95%) can be passed as a whole to GPT-3.5. For the rest of the (longer) documents, we break the documents into two chunks, summarize each chunk individually, and then append the chunk summaries to get the final summary. For GPT-4, all documents can be passed without chunking.

**Data Selection:** We chose all 100 documents from the test set of In-Abs for passing to ChatGPT. For ICL, we sample from this set of documents itself. We try to sample the smallest samples for the longer input examples to fit the entire prompt within GPT token length limit.

| Model | R-1 | R-2 | R-L | BERTScore |
|---|---|---|---|---|
| GPT-3.5 0-shot | 0.392 | 0.165 | 0.208 | 0.847 |
| GPT-3.5 1-shot | 0.385 | 0.141 | 0.201 | 0.835 |
| GPT-3.5 2-shot | 0.419 | 0.164 | 0.215 | 0.838 |
| GPT-4 0-shot | **0.472** | **0.183** | **0.228** | **0.848** |
| GPT-4 1-shot | 0.304 | 0.059 | 0.161 | 0.807 |
| GPT-4 2-shot | 0.324 | 0.080 | 0.171 | 0.813 |

Table 43: Rouge-1,2,L and BERTScore scores for SUMM

**Verbalizer:** The entire output returned by GPT is considered as the abstractive summary. We specifically instruct GPT to output simple text, without headers and bullet points. We observe that both GPT-3.5 and GPT-4 mostly adhere to this instruction.

**Results:** GPT results are shown in Table 24. Even for this task (which GPT is quite conversant with), the gap in performance compared to the SOTA is significant. GPT-4 performs slightly better than GPT-3.5 in a zero-shot setting. For both models, 1-shot prompting reduces the performance, although the difference is much larger for GPT-4, which performs very poorly in a 1-shot setting. On using 2-shot setting, we observe an improvement again, but it is not too much (for GPT-4, 2-shot still performs significantly poorly compared to 0-shot). It is possible that the IC examples are actually confusing the model rather than helping it.

### C.1.7 Legal Machine Translation (L-MT)

**Prompt Design:** GPT is known to perform translations effectively. Hence, we provide the model with just the input sentence (in English), and we ask the model to translate the sentence to the desired target language (see Table 44).

**Data Selection:** We randomly choose 5 samples from each target language from each MILPaC dataset. This gives us 45 documents each for MILPaC-IP and MILPaC-Acts (9 target languages), and 20 documents for MILPaC-CCI-FAQ (4 target languages), giving us a total of 110 samples. It should be noted that all datasets contain two types of samples – questions and answers. However, the answers from the MILPaC-CCI-FAQ dataset consist of just a single number corresponding to different choices in the MCQ setting. Thus, we do not choose the answer samples from MILPaC-CCI-FAQ. For ICL, we randomly choose samples from the same target language in the same dataset.

**Verbalization:** We directly take the entire GPT

```
SYSTEM_PROMPT: You are a smart and intelligent machine translation
system, trained to read Indian legal texts and translate them to
Indian languages. Your task is, given an English language sentence
from a legal document, translate it to the given target Indian
language. I will provide you with the input/output format, target
language and the sentence to be translated. I will also provide
some examples of the task.

USER_PROMPT: Are you clear about your role?

ASSISTANT_PROMPT: Sure, I'm ready to help you with your legal
translation task. Please provide me with the sentence and the
target language I am supposed to translate to.

INPUT_PROMPT:
Examples:
Sentence 1 in English: <In-context Sentence 1 goes here>
Sentence 1 in <Target language goes here>: <Reference translation
for Sentence 1 goes here>
...
Sentence n+1 in English: <Test Document goes here>
Sentence n+1 in <Target language goes here>:
```

Table 44: Prompt template for L-MT for both GPT-3.5 and GPT-4(for $n$ in-context examples)

output as the translation.

**Results:** GPT-3.5 produces decent results for L-MT as compared to SOTA approaches, possibly due to GPT's prior knowledge on this task. This is further improved upon by GPT-4, which actually outperforms SOTA on average for two of three metrics. For both models, there is a drop in performance for Acts, possibly due to the more complex nature of the text in the Acts dataset (Mahapatra et al., 2023). We see a gradual improvement across all metrics and all datasets for both GPT-3.5 and GPT-4 with an increasing degree of ICL, with 2-shot prompting producing the best results.

## C.2 Experiments with Smaller LLMs

In addition, we also experimented with other large language models with smaller parameter sizes. Specifically, we experimented with GPT-Neo (Black et al., 2021) family of three models (GPT-Neo-125M, GPT-Neo-1.3B, GPT-Neo-2.7B) trained on the Pile dataset (Gao et al., 2020), GPT-J-6B (Wang and Komatsuzaki, 2021), Llama-2-7b-chat-hf (Touvron et al., 2023), and recently released Mistral-7B-v0.1 (Jiang et al., 2023) language models for our experiments. The primary challenge when validating the smaller language model is the prompt design. Following previous works (Brown et al., 2020; Robinson and Wingate, 2023), we pose the prompt in a multiple-choice question-answering format (a prompt sample for various tasks present in the benchmark can be found in the supplementary material) and validate the performance using the obtained log probability of the predicted tokens as highlighted in (Robinson and Wingate, 2023). Moreover, since the tasks are more complicated with larger context lengths, the generative models

sometimes generate some irrelevant tokens. For those cases with random token generation, we consider it to be a failure case and use a random prediction as a proxy of predictions. Overall, we observed that all the language models perform poorly with near-random predictions over the proposed set of legal language understanding tasks.

We speculate two primary reasons for this finding. First, the language models we used are not explicitly designed to capture the question-answering format for a larger context. Since the context length of the task in the proposed benchmark is significantly higher than the other NLU tasks, it becomes more challenging for smaller language models to decode the question-answer format required for performing these tasks. Second, these models lack the instruction tuning strategies followed by larger models like GPT3.5, making it much harder to capture the context. Moreover, our experiments with GPT3.5 also suggest that if the context is large, even the larger models fail to capture the requested instructions present in the query prompt.

| Dataset | # Shots | GPT-3.5 | | | GPT-4 | | |
|---|---|---|---|---|---|---|---|
| | | BLEU | GLEU | chrF++ | BLEU | GLEU | chrF++ |
| MILPac-IP | 0 | 26.2 | 30.3 | 45.3 | 36.3 | 39.2 | 53.5 |
| | 1 | 27.8 | 31.5 | 46.3 | 37.6 | 40.5 | 54.0 |
| | 2 | 27.9 | 31.0 | 45.4 | **38.0** | **40.6** | **54.5** |
| MILPaC-CCI-FAQ | 0 | 24.1 | 28.2 | 43.9 | 35.0 | 36.1 | 50.0 |
| | 1 | 25.9 | 28.7 | 43.8 | 39.0 | 41.4 | 56.6 |
| | 2 | 27.9 | 30.6 | 44.9 | **42.2** | **43.3** | **57.2** |
| MILPaC-Acts | 0 | 18.2 | 23.1 | 36.0 | 29.0 | **32.6** | 45.7 |
| | 1 | 19.5 | 23.6 | 36.6 | 28.8 | 32.3 | 45.3 |
| | 2 | 21.2 | 24.8 | 38.2 | **29.1** | 32.4 | **46.2** |
| Average | 0 | 22.8 | 28.2 | 43.9 | 33.4 | 36.1 | 50.0 |
| | 1 | 24.4 | 27.9 | 42.3 | 35.1 | 38.0 | 52.0 |
| | 2 | 25.6 | 28.8 | 42.8 | **36.4** | **38.7** | **52.6** |

Table 45: Corpus-level BLEU, GLEU, and chrF++ scores for GPT-3.5 and GPT-4 prompting with 0, 1 and 2 shot settings