TathyaNyaya and FactLegalLlama: Advancing Factual Judgment Prediction and Explanation in the Indian Legal Context

Anonymous ACL submission

Abstract

In the landscape of Fact-based Judgment Prediction and Explanation (FJPE), reliance on factual data is essential for developing robust and realistic AI-driven decision-making 005 tools. This paper introduces TathyaNyaya, the largest annotated dataset for FJPE tailored to the Indian legal context, encompassing judgments from the Supreme Court of India and various High Courts. Derived from the Hindi terms "Tathya" (fact) and "Nyaya" (justice), the TathyaNyaya dataset is uniquely 011 designed to focus on factual statements rather than complete legal texts, reflecting real-world judicial processes where factual data drives outcomes. Complementing this dataset, we present FactLegalLlama, an instruction-tuned 017 variant of the LLaMa-3-8B Large Language Model (LLM), optimized for generating high-019 quality explanations in FJPE tasks. Finetuned on the factual data in TathyaNyaya, FactLegalLlama integrates predictive accu-021 racy with coherent, contextually relevant explanations, addressing the critical need for transparency and interpretability in AI-assisted legal systems. Our methodology combines transformers for binary judgment prediction with FactLegalLlama for explanation generation, creating a robust framework for advancing FJPE in the Indian legal domain. TathyaNyaya not only surpasses existing datasets in scale and diversity but also establishes a benchmark for building explainable AI systems in legal analysis. The findings underscore the importance of factual precision and domain-specific tuning in enhancing predictive performance and interpretability, positioning TathyaNyaya and FactLegalLlama as foundational resources for AI-assisted legal decision-making.

1 Introduction

040

043

The integration of AI technologies into the legal domain holds immense potential for improving the efficiency, accessibility, and transparency of judicial processes, especially in jurisdictions like India, where case backlogs place a significant burden on the courts. Among the emerging solutions, Factbased Judgment Prediction and Explanation (FJPE) aims to predict judicial outcomes and provide explainable rationales by focusing exclusively on the factual elements of a case. By highlighting the core facts rather than relying on entire legal documents, FJPE seeks to deliver more realistic early-phase predictions, akin to conditions when a judge must form an opinion before hearing all arguments or awaiting additional documents. 044

045

046

047

051

055

058

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

075

076

078

081

Recent studies have explored FJPE by extracting factual elements from full case texts or summarizing multiple components like statutes, ratios of decisions, and arguments (Nigam et al., 2024b; Nigam and Deroy, 2024). While these efforts represent a step toward fact-centric approaches, they often rely on limited or automatically extracted data that may not always be reliable. This distinction is crucial: unlike Court Judgment Prediction with Explanation (CJPE), which may utilize entire judgments including judicial reasoning and statutory references, FJPE restricts itself to facts presented at the time of decision-making. Such a focus mimics the constraints of real-world legal reasoning, where judges form preliminary opinions based on presented facts before considering lengthy arguments, precedents, or statutes.

In this paper, we introduce TathyaNyaya, the first large-scale, expertly annotated dataset dedicated to FJPE in the Indian legal context. Derived from the Hindi terms "Tathya" (fact) and "Nyaya" (justice), TathyaNyaya places factual segments at the heart of predictive modeling, thereby providing a robust foundation that eschews the noise and complexity of full legal documents. This annotated dataset addresses a critical gap: while prior work has attempted fact-based prediction using incomplete or non-annotated sources, TathyaNyaya represents the first systematic effort to assemble, annotate, and align factual information with judi-

168

169

170

171

172

173

174

175

176

177

178

179

180

181

182

183

134

135

136

cial outcomes and explanations. By simplifying the input space to factual data, our dataset enables models to more accurately reflect real-world conditions, offer transparent and early-stage predictive insights, and reduce reliance on domain-level heuristics or post-hoc explanations.

086

087

090

094

100

103

104

105

106

107

108

110

111

112

113

114

115

116

117

118

119

121

122

123

124

125

127

128

129

130

131

132

133

TathyaNyaya incorporates judgments from the Supreme Court of India (SCI) and various High Courts, ensuring broad coverage and diversity. Comprising four components—NyayaFacts, NyayaScrape, NyayaSimplify, and NyayaFilter, the dataset caters to different facets of factual extraction, simplification, and retrieval, thereby facilitating comprehensive FJPE research.

To complement this dataset, we present FactLegalLlama, an instruction-tuned variant of LLaMa-3-8B fine-tuned specifically for FJPE tasks. FactLegalLlama is explicitly adapted to produce faithful, fact-grounded explanations.

Our key contributions are summarized as:

- TathyaNyaya *Dataset:* We introduce TathyaNyaya, the first extensively annotated, purely fact-centric dataset for judgment prediction and explanation in the Indian legal context, providing a stable benchmark for FJPE research. The dataset is structured into four components, each addressing distinct aspects of factual legal analysis.
- Factual Basis for Prediction: Unlike previous datasets that use full legal texts, TathyaNyaya focuses on factual information, allowing for more realistic and targeted legal judgment predictions.
- FactLegalLlama for Explanation: We introduce FactLegalLlama, a LLaMa-3-8B model finetuned to generate high-quality, factual explanations for judicial outcomes, enhancing interpretability and trust in AI-assisted legal systems. To ensure reproducibility and encourage further

research, the dataset and model code are made publicly available¹.

2 Related Work

The domain of Legal Judgment Prediction (LJP) has seen substantial progress, driven by a need for efficient and transparent decision-making aids in the legal system. Early foundational works (Aletras et al., 2016; Chalkidis et al., 2019; Feng et al., 2021) focused on predicting case outcomes while emphasizing interpretability, laying the groundwork for methodologies and benchmark datasets like CAIL2018 (Xiao et al., 2018) and ECHR-CASES (Chalkidis et al., 2019). These resources fostered the development of advanced models including TopJudge and MLCP-NLN, yet a consistent performance gap remains between automated predictions and actual judicial decisions.

Within the Indian legal ecosystem, researchers have introduced datasets such as ILDC (Malik et al., 2021) and PredEx (Nigam et al., 2024a), alongside related efforts (Nigam et al., 2022; Malik et al., 2022; Nigam et al., 2023), underscoring the significance of domain-specific datasets and the need for explanations that can be integrated into real-world legal workflows. Studies leveraging Large Language Models (LLMs) (Vats et al., 2023; Nigam et al., 2024a) have shown that models such as GPT-3.5 Turbo and LLaMa-2 can adapt to Indian legal texts, further diversifying approaches to LJP and interoperability. Fact-based judgment prediction has gained prominence as a realistic approach to LJP, focusing on predictions derived from case facts rather than full case judgments. Nigam et al. (2024b) and Nigam and Deroy (2024) explore LJP based on facts, arguing that this approach better simulates real-world scenarios.

Cross-jurisdictional research, including that by (Zhao et al., 2018), has extended LJP methodologies across different legal frameworks. Multilingual efforts like (Niklaus et al., 2021; Kapoor et al., 2022) and approaches integrating event extraction and multi-stage reasoning (Feng et al., 2022) continue to broaden the capabilities of LJP systems.

3 Task Description

Our work centers on predicting and explaining legal judgments from the Supreme Court of India (SCI) and various High Court cases using a newly introduced annotated dataset, TathyaNyaya. This dataset is the largest of its kind for factual judgment prediction and explanation in the Indian legal domain. Unlike prior approaches relying on full case texts, TathyaNyaya emphasizes factual information alone, reflecting more realistic conditions for automated legal decision-making.

The Fact-based Judgment Prediction and Explanation (FJPE) task consists of two subtasks:

Task A: Judgment Prediction: This is a binary classification problem. Given the factual information of the legal judgment as input, the objective is to predict whether the decision favors or goes

¹Anonymous GitHub Link



Figure 1: A high-level illustration of the TathyaNyaya dataset creation pipeline, showcasing the development process and interconnections of its four components.

against the appellant. The prediction is represented by binary labels: "1" indicates that the appeal is accepted (i.e., if any part of the appeal is accepted, the decision is considered in favor of the appellant), while "0" indicates that the appeal is rejected.

Task B: Rationale Explanation: This subtask involves generating an explanation or rationale that justifies the predicted decision based on the provided factual information of the case. The goal is to provide a clear understanding of the reasoning behind the predicted outcome, grounded in the facts presented.

Figure 2 in the Appendix illustrates the overall process of fact-based judgment prediction and explanation employed in our study, outlining the sequential steps from prediction to explanation based on the factual data provided.

4 Dataset

185

186

189

190

192

193

194

195

196

197

198

199

In this research, we introduce TathyaNyaya, a comprehensive dataset explicitly designed for Fact-based Judgment Prediction and Explana-204 tion (FJPE) in the Indian legal domain. This dataset consists of four distinct components: (1) NyayaFacts-expert-annotated data that serves as the gold standard for prediction and explanation tasks, (2) NyayaScrape—automated fact-extracted 209 data obtained through machine-driven processes, 210 (3) NyayaSimplify—a user-friendly dataset created by paraphrasing complex legal language, and 212

| Metric | Train (Multi) Train (Single) | | Validation | Test | | | |
|----------------|------------------------------|-------|------------|-------|--|--|--|
| NyayaFacts | | | | | | | |
| # Documents | 13,629 | 8,216 | 1,197 | 2,389 | | | |
| Avg # Words | 855 | 853 | 828 | 865 | | | |
| Acceptance (%) | 55.20 | 47.66 | 47.45 | 47.72 | | | |
| NyayaScrape | | | | | | | |
| # Documents | 8,993 | 3,828 | 548 | 1,095 | | | |
| Avg # Words | 405 | 404 | 412 | 405 | | | |
| Acceptance (%) | 65.77 | 61.44 | 59.85 | 60.55 | | | |

Table 1: Statistics for NyayaFacts and NyayaScrape datasets from the TathyaNyaya corpus.

213

214

215

216

217

218

219

221

222

223

224

225

226

227

229

230

231

233

234

235

237

238

239

240

241

242

243

244

245

247

248

249

250

(4) NyayaFilter—a binary fact vs. non-fact classification dataset designed to streamline the retrieval of relevant factual information. Together, these components form the largest and most diverse factual dataset in the Indian judiciary, enabling the development and evaluation of advanced AI models for transparent and interpretable judgment prediction and explanation. By focusing exclusively on factual data, TathyaNyaya addresses a critical gap in the field, paving the way for more robust and realistic AI-driven solutions tailored to the Indian legal context.

Figure 1 illustrates the TathyaNyaya dataset creation pipeline. It provides a high-level overview of how each component which is derived, from expert-curated facts and machine-driven extraction, to fact segmentation and paraphrasing. This end-toend pipeline ensures that the final dataset captures both breadth and depth in factual legal information, supporting the FJPE task.

4.1 Dataset Compilation and Statistics

The compilation process involved collecting approximately 16,000 judgments from the Supreme Court of India (SCI) and various High Courts through IndianKanoon², a widely used legal search engine known for its comprehensive repository of Indian legal documents. These judgments were then categorized into the following components:

4.1.1 NyayaFacts

NyayaFacts comprises a subset of Supreme Court of India (SCI) and High Court judgments carefully annotated by legal experts. These annotations highlight key factual segments that significantly influence judicial outcomes, serving as high-quality ground truth for both judgment prediction and rationale explanation. After refining and preprocessing, this subset serves as the gold standard for evaluating prediction and explanation tasks.

²https://indiankanoon.org/

In particular, the validation and test data were derived from the NyayaFacts Single subset to maintain consistency during evaluation, while the training data include both single and multi-case judgments, offering a broad learning landscape. Table 1 provides comprehensive statistics. NyayaFacts thus provides a high-quality benchmark for both judgment prediction and explanation tasks.

4.1.2 NyayaScrape

251

257

261

262

263

265

269

270

271

272

277

278

290

291

297

300

NyayaScrape comprises judgments sourced from the Indiankanoon website, where cases are automatically segmented into various categories such as facts, issues, conclusions, and assessments of how the courts have treated certain elements (e.g., "Negatively Viewed by Court," "Relied by Party,"
"Accepted by Court"). Although these segments aim to provide structured insights, the labels are not entirely reliable. They are generated by automated tools rather than human legal experts, resulting in potential inconsistencies and may introduce noise. Moreover, not all judgments contain every type of label, further complicating the data's uniformity.

Despite these limitations, NyayaScrape offers valuable machine-derived factual extractions that enable us to compare expert-driven annotations with automated processes. This comparison helps assess the reliability, quality, and shortcomings of model-based fact identification and segmentation. Document-level statistics and comparisons against NyayaFacts are provided in Table 1.

4.1.3 NyayaSimplify

NyayaSimplify focuses on making complex legal texts more accessible by paraphrasing the NyayaFacts test data into clearer, more concise language. Using LLaMA-3-70B-Instruct, intricate legal jargon is transformed into user-friendly text without altering the factual content or legal reasoning. Since NyayaSimplify is derived directly from NyayaFacts, the majority of dataset statistics remain comparable, with the primary difference being a reduced average word count. This simplification aims to improve both the accuracy and interpretability of models in FJPE tasks. Detailed prompt used for paraphrasing is provided in Appendix Table 7.

4.1.4 NyayaFilter

NyayaFilter addresses the challenges of manual annotation by employing a BiLSTM-CRF model to classify sentences as either factual (1) or nonfactual (0). This binary classification replaces

| Metric | Train | Validation | Test | | |
|------------------------------|-----------|------------|---------|--|--|
| Facts | | | | | |
| # Documents | 14,134 | 1,197 | 2,389 | | |
| # Sentences | 15,83,858 | 49,671 | 56,240 | | |
| Avg # Words | 29.03 | 29.00 | 34.00 | | |
| Avg # Facts/Document (%) | 25.1 | 20.61 | 22.7 | | |
| Overall Facts (%) | 20.34 | 12.86 | 18.46 | | |
| Non-Facts | | | | | |
| # Documents | 14,134 | 1,197 | 2,389 | | |
| # Sentences | 4,04,349 | 3,36,478 | 248,433 | | |
| Avg # Words | 28.06 | 27.00 | 30.00 | | |
| Avg # Non-Facts/Document (%) | 74.9 | 79.39 | 77.3 | | |
| Overall Non-Facts (%) | 79.66 | 87.14 | 81.54 | | |

Table 2: Comparison of factual vs. non-factual statistics used during BiLSTM-CRF classifier training for the NyayaFilter dataset.

the traditional multi-label approach, simplifying the task while maintaining a focus on essential factual information. The model was trained on NyayaFacts Single data, with validation and testing on the corresponding splits. This approach achieved approximately 90% accuracy in separating factual statements, as shown in Table 2. This dataset streamlines the retrieval process for FJPE tasks and enables scalable fact extraction. 301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

331

4.2 Annotation Methodology and Quality Assurance

4.2.1 Expert Participation

The annotation process for NyayaFacts was carried out by a team of 10 legal experts, comprising advanced third- and fourth-year law students from premier Indian law colleges. These individuals were chosen based on their academic standing, legal reasoning skills, and familiarity with judicial processes, ensuring that the annotations reflected high-quality and domain-relevant insights.

4.2.2 Timeline and Workload Distribution

The annotation process was conducted over an extended period (April 1, 2022, to October 30, 2023), reflecting the complexity and precision required to analyze diverse legal texts. Each annotator was assigned approximately 30 judgment documents per week, a volume that balanced efficiency with attention to detail. This measured pace allowed the annotators to thoroughly examine the factual segments without compromising quality.

4.2.3 Annotation Protocol

The annotators were tasked with identifying and ex-
tracting specific judgment segments that contained332factual information, without personal interpreta-
tion or summarization. This approach preserved333

430

431

432

433

434

385

386

387

389

the authenticity of the annotations, ensuring that
they faithfully represented the judicial reasoning
within each document.

4.2.4 Quality Control Framework

339

341

342

343

345

351

366

371

372

373

374

376

377

381

384

To maintain annotation consistency and reliability, a multi-layered quality control mechanism was implemented:

- Initial Review: Each case was initially annotated by a single expert. This ensured efficiency while maintaining focus on factual segments. Subsequently, the annotations underwent multiple validation layers.
- Senior Expert Validation: Discrepancies or ambiguous annotations were escalated to a review panel comprising senior legal practitioners, who provided final judgments on contentious segments, enhancing the reliability of the final annotations.
- Training and Alignment Meetings: Regular training sessions and coordination meetings were conducted to align all annotators on annotation protocols, legal conventions, and factual identification criteria. These interactive forums helped minimize subjectivity, solidify common standards, and maintain uniform annotation quality throughout the project's duration.

This robust framework ensured the annotations' accuracy and credibility, making NyayaFacts a reliable benchmark for evaluating prediction and explanation models.

5 Methodology

In this section, we present our overall methodology for extracting factual segments from legal judgments, training our custom model FactLegalLlama for Fact-based Judgment Prediction and Explanation (FJPE), and finally addressing both the prediction-only and prediction-withexplanation tasks. We also detail the prompts we used and instruction-tuning strategies employed to refine our model's outputs.

5.1 Fact Extraction from Full Legal Judgments

To prepare the dataset for Fact-based Judgment Prediction and Explanation (FJPE), we first extracted the factual statements from full-text legal judgments. We adopted a streamlined binary classification approach by fine-tuning a BiLSTM-CRF model (Ghosh and Wyner, 2019), a previous stateof-the-art (SoTA) model for semantic segmentation of legal documents. Instead of using the original multi-class rhetorical role framework, which distinguishes between roles such as issue, statute, precedent, and argument, we simplified the task by treating all non-factual segments as a single class labeled "non-facts."

This transformation into a binary classification problem enabled the model to focus solely on identifying factual segments critical to judgment prediction. Training was conducted using the NyayaFacts multi, which provided expertannotated labels for factual and non-factual segments. By isolating the facts, we laid the groundwork for developing AI models capable of making decisions and generating explanations based solely on factual data. This preprocessing ensured that the subsequent models trained on the dataset remained focused on the most relevant and actionable information in legal cases.

5.2 Training FactLegalLlama

The FactLegalLlama model, based on the LLaMa-3-8B architecture, was fine-tuned specifically for the FJPE task using NyayaFacts. The training process involved instruction-tuning with a diverse set of 16 templates designed to guide the model in judgment prediction and explanation tasks. We utilized low-rank adaptation (LoRA) to optimize model training on limited computational resources. Training parameters, such as quantization to 4-bit precision and gradient accumulation, ensured efficient usage of resources while maintaining model performance.

To further enhance its capabilities, FactLegalLlama was fine-tuned with both prediction-only and prediction-with-explanation tasks, enabling it to handle a wide range of factual judgment scenarios. The fine-tuning process emphasized the use of simplified prompts to ensure clarity and relevance in the generated outputs.

5.3 Fact-Based Judgment Prediction

5.3.1 Language Model-Based Approach

For baseline comparisons, we utilized transformerbased models like InLegalBERT (Paul et al., 2023), and XLNet Large (Yang et al., 2019) for binary classification. Due to the token length constraints of these models, we adopted a chunking strategy by dividing documents into 512-token segments with a 100-token overlap to preserve context. Chunklevel predictions were aggregated to generate final case-level predictions.

435

- 438
- 439 440
- 441 442
- 443 444

445

446

447

448

449

450

451

452

453

454

455

456

457

458

459

460

461

462

463

464

465

466

467

468

469

470

471

472

473

474

475

476 477

478

479

480

481

5.3.2 Large Language Model-based Approach

We utilized FactLegalLlama, our instructiontuned LLaMa-3-8B model (Dubey et al., 2024), for judgment prediction-only instructions, where the model predicts judicial outcomes solely based on the factual inputs. The training data from TathyaNyaya was used to train the factual prediction context, emphasizing precision.

5.4 Fact-Based Judgment Prediction with Explanation (FJPE)

For the combined task of prediction and explanation, we employed FactLegalLlama with modified instruction prompts. Instructions guided the model to first predict the outcome and then generate a rationale grounded in the provided factual data.

5.5 Prompts Used

Prompts for both prediction and explanation tasks were carefully designed and adapted from previous studies (Vats et al., 2023; Nigam et al., 2024c,d). For prediction-only tasks, the prompts instructed the model to output a binary decision. For prediction-with-explanation tasks, the prompts included directives to explain the reasoning behind the prediction. These templates are detailed in Table 6 in the Appendix.

5.6 Instruction Sets

The fine-tuning process for FactLegalLlama involved using a diverse set of 16 instruction templates for judgment prediction and explanation. These templates ensured the model could generalize effectively across a wide range of cases and factual scenarios. The complete list of instruction sets used for tuning is in Table 8 in the Appendix.

6 Evaluation Metrics

To rigorously assess the performance of our models on judgment prediction and factual explanations in the TathyaNyaya test dataset, we employed a suite of evaluation metrics. For judgment prediction, we report Macro Precision, Macro Recall, Macro F1, and Accuracy. For evaluating the quality of explanations, both quantitative and qualitative methods were applied.

 Lexical-Based Evaluation: We used traditional lexical similarity metrics, including ROUGE (ROUGE-1, ROUGE-2, and ROUGE-L) (Lin, 2004), BLEU (Papineni et al., 2002), and ME-TEOR (Banerjee and Lavie, 2005). These metrics measure word overlap and sequence alignment between generated explanations and reference texts, providing a quantitative measure of the accuracy of lexical content.

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

500

501

502

503

504

505

507

508

509

510

511

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

- 2. Semantic Similarity Evaluation: To assess the semantic alignment of the generated explanations, we applied BERTScore (Zhang et al., 2020), which evaluates semantic similarity between the generated text and reference explanations. Additionally, BLANC (Vasilyev et al., 2020) was utilized to estimate the contextual relevance and coherence of the generated text in the absence of a gold-standard reference.
- 3. Inter-Annotator Agreement: In the construction of NyayaFacts, each document was annotated by a single annotator due to the complexity and scale of the dataset. While this approach ensured a manageable workflow and timely completion, it precludes the direct calculation of inter-annotator agreement. As a result, we do not report inter-annotator metrics, and future work may consider sampling subsets of the data for multiple annotations to facilitate such evaluations.

7 Results and Analysis

In this section, we present and interpret the performance of our models across various datasets and experimental settings. We focus first on raw judgment prediction results using NyayaFacts and NyayaScrape data, then on the performance improvements or trade-offs observed in the NyayaFilter and NyayaSimplify settings. Finally, we analyze the explanation quality generated by FactLegalLlama using both lexical and semantic metrics.

7.1 Performance on NyayaFacts and NyayaScrape

We begin by examining model performances on the NyayaFacts and NyayaScrape test sets, as reported in Table 3. Each model (InLegalBERT, XL-Net_Large, and FactLegalLlama) was evaluated under different training configurations, including Single and Multi.

Language Model-Based Baselines: Across both NyayaFacts and NyayaScrape test sets, XL-Net_Large consistently outperforms InLegalBERT on macro Precision, Recall, F1, and Accuracy metrics. For instance, when trained on NyayaFacts Single, XLNet_Large achieves a macro F1 of

| Model | Macro Precision | Macro Recall | Macro F1 | Accuracy | Training Data | | |
|----------------------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------------------|-----------------------|--|--|
| Results on NyayaFacts Test Data | | | | | | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.5934 0.6064 0.5416 | 0.5936 0.6040 0.5312 | 0.5935 0.6052 0.5036 | 0.5932 0.6061 0.5386 | NyayaFacts Single | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.6001 0.6145 0.5390 | 0.5836 0.5965 0.5368 | 0.5917 0.6054 0.5318 | 0.5740 0.5908 0.5401 | NyayaFacts Multi | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.5480 0.5807 0.5139 | 0.5192 0.5781 0.5122 | 0.5332 0.5794 0.4922 | 0.5082 0.5756 0.5042 | NyayaScrape Single | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.5735 0.5935 0.4951 | 0.5269 0.5878 0.4966 | 0.5492 0.5906 0.4516 | 0.5157 0.5842 0.4884 | NyayaScrape Multi | | |
| | Results on NyayaScrape Test Data | | | | | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.6718 0.6754 0.5574 | 0.5748 0.6394 0.5372 | 0.6195 0.6569 0.5191 | 0.6521 0.6849 0.6045 | NyayaScrape Single | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.7976 0.8098 0.5439 | 0.7268 0.7781 0.5317 | 0.7606 0.7936 0.5177 | 0.7717 0.8055 0.5877 | NyayaScrape Multi | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.6237 0.5433 0.5832 | 0.5243 0.5282 0.5868 | 0.5697 0.5357 0.5792 | 0.6183 0.5918 0.5840 | NyayaFacts Single | | |
| InLegalBert XLNet_Large FactLegalLlama | 0.6784 0.6124 0.6541 | 0.5027 0.5129 0.6583 | 0.5775 0.5583 0.6552 | 0.6073 0.6119 0.6651 | NyayaFacts Multi | | |

Table 3: Performance metrics of models evaluated on NyayaFacts and NyayaScrape test data. Each block shows results obtained by training on either NyayaFacts or NyayaScrape data (single or multi variants), then testing on corresponding subsets. The best scores in each section are highlighted in bold.

0.6052 and Accuracy of 0.6061, surpassing InLegalBERT's macro F1 of 0.5935 and Accuracy of 0.5932. This trend persists in most training and testing configurations, highlighting XLNet_Large's robust capability for factual judgment prediction in the given domain.

531

532

535

536

FactLegalLlama's Prediction-Only Perfor-537 mance: FactLegalLlama, while instructiontuned for outcome prediction, lags behind the 539 transformer-based baselines in raw prediction 540 performance. For example, when trained on 541 NyayaFacts Single and tested on NyayaFacts, 542 it obtains a macro F1 of 0.5036 compared to XLNet_Large's 0.6052. A similar gap is observed 544 across other splits. Although FactLegalLlama 545 underperforms in direct classification metrics, 546 its strength lies in generating explanations, as discussed later. 548

549Single vs. Multi Cases: Both baselines and550FactLegalLlama exhibit more stable performance551on the Single subsets compared to the Multi sub-552sets. The complexity introduced by multiple pe-553titions with varying outcomes in the Multi cases

| Model | Macro Precision | Macro Recall | Macro F1 | Accuracy | Training Data | | |
|------------------------------------|----------------------------------|-------------------------|----------------------|-------------------------|---------------------|--|--|
| | Results on NyayaFilter Test Data | | | | | | |
| InLegalBert | 0.5870 | 0.5857 | 0.5864 | 0.5885 | NyayaFacts | | |
| XLNet_Large | 0.5805 | 0.5775 | 0.5790 | 0.5818 | Single | | |
| InLegalBert | 0.5886 | 0.5560 | 0.5719 | 0.5421 | NyayaFacts | | |
| XLNet_Large | 0.5977 | 0.5874 | 0.5925 | 0.5797 | Multi | | |
| InLegalBert | 0.5342 | 0.5180 | 0.5260 | 0.5023 | NyayaScrape | | |
| XLNet_Large | 0.5577 | 0.5509 | 0.5543 | 0.5429 | Single | | |
| InLegalBert | 0.5789 | 0.5409 | 0.5592 | 0.5249 | NyayaScrape | | |
| XLNet_Large | 0.5581 | 0.5364 | 0.5470 | 0.5224 | Multi | | |
| Results on NyayaSimplify Test Data | | | | | | | |
| InLegalBert | 0.6199 | 0.6197 | 0.6198 | 0.6167 | NyayaFacts | | |
| XLNet_Large | 0.6179 | 0.6169 | 0.6174 | 0.6200 | Single | | |
| InLegalBert XLNet_Large | 0.6222 0.6160 | 0.5986 0.6002 | 0.6102 0.6080 | 0.5839 0.5878 | NyayaFacts Multi | | |
| InLegalBert | 0.5760 | 0.5311 | 0.5526 | 0.5061 | NyayaScrape | | |
| XLNet_Large | 0.5864 | 0.5845 | 0.5854 | 0.5789 | Single | | |
| InLegalBert | 0.5659 | 0.5215 | 0.5428 | 0.4950 | NyayaScrape | | |
| XLNet_Large | 0.5978 | 0.5891 | 0.5934 | 0.5789 | Multi | | |

Table 4: Model performance on NyayaFilter and NyayaSimplify test datasets. For NyayaFilter, results illustrate how automatically retrieved factual data affects performance when models are trained on NyayaFacts or NyayaScrape datasets. For NyayaSimplify, results show the impact of paraphrasing complex legal texts into simpler language. Bolded scores indicate the best performance in each section.

reduces overall accuracy and F1 scores, emphasizing the challenge of fact-based judgment prediction in more intricate legal scenarios.

554

555

556

557

558

559

560

561

562

563

564

565

566

567

568

569

570

571

572

573

574

575

576

577

7.2 Impact of Fact Retrieval (NyayaFilter) and Text Simplification (NyayaSimplify)

Table 4 reports model performances on the NyayaFilter and NyayaSimplify test datasets. These results highlight how the preprocessing choices affect model accuracy on automatic fact retrieval and paraphrasing complex legal texts.

NyayaFilter **Results:** When comparing NyayaFilter results to the original NyayaFacts and NyayaScrape sets, we see that while performance can fluctuate, some models benefit from training on data where fact and non-fact segments are clearly distinguished. For example, on the NyayaFilter test set derived from NyayaFacts Single, InLegalBERT attains a macro F1 of 0.5864, maintaining competitive performance. XLNet_Large, although not always the top performer here, still sustains a strong baseline. These findings suggest that automatically retrieved factual subsets can be used without severely degrading model performance.

| Training Testing | | Lexical Based Evaluation | | | | | Semantic Evaluation | |
|--------------------|---------------|--------------------------|---------|---------|--------|--------|---------------------|--------|
| Data | Data | Rouge-1 | Rouge-2 | Rouge-L | BLEU | METEOR | BERT Score | BLANC |
| No Training | NyayaFacts | 0.2757 | 0.0998 | 0.1448 | 0.0395 | 0.1650 | 0.5250 | 0.0814 |
| No Training | NyayaScrape | 0.1877 | 0.0750 | 0.1297 | 0.0371 | 0.1848 | 0.4819 | 0.0927 |
| NyayaFacts Single | NyayaFacts | 0.3216 | 0.1085 | 0.1897 | 0.0419 | 0.1798 | 0.5785 | 0.0958 |
| NyayaFacts Multi | NyayaFacts | 0.3383 | 0.1133 | 0.1950 | 0.0483 | 0.2120 | 0.5843 | 0.1031 |
| NyayaScrape Single | NyayaScrape | 0.1172 | 0.0480 | 0.0864 | 0.0211 | 0.0967 | 0.3927 | 0.0609 |
| NyayaScrape Multi | NyayaScrape | 0.1720 | 0.0762 | 0.1269 | 0.0314 | 0.1287 | 0.4516 | 0.0782 |
| NyayaSimplify | NyayaSimplify | 0.2818 | 0.0764 | 0.1816 | 0.0242 | 0.1689 | 0.5559 | 0.0731 |

Table 5: Performance of FactLegalLlama on the FJPE task. The base model is LLaMa-3-8B. "No Training" indicates results from the unmodified (vanilla) model. Other rows show improvements after fine-tuning with different subsets of the TathyaNyaya data. Bolded values represent the best performance within a given evaluation scenario.

NyayaSimplify Results: Paraphrasing complex legal language into simpler text (the NyayaSimplify scenario) generally helps models retain or slightly improve performance. For instance, with NyayaFacts Single, InLegalBERT reaches a macro F1 of 0.6198 and XLNet_Large hits an Accuracy of 0.6200 on the simplified data, both representing small yet noteworthy improvements compared to their performance on the original complex texts. This trend indicates that reducing linguistic complexity can aid models in understanding and classifying factual statements more accurately.

578

579

583

584

585

587

588

589

590

592

594

595

596

7.3 Quality of Explanations from FactLegalLlama

Table 5 presents the evaluation of FactLegalLlama on the explanation generation task, measured through both lexical (Rouge, BLEU, METEOR) and semantic (BERTScore, BLANC) metrics. We compare a "No Training" scenario (using the LLaMa-3-8B model) with fine-tuned versions on different subsets of TathyaNyaya data.

601Fine-tuning Benefits:Fine-tuning LLaMa-3-8B602(FactLegalLlama) on factual data substantially im-603proves its explanation quality. For NyayaFacts,604training on the Multi subset yields the strongest605results, with Rouge-1 at 0.3383 and a BERTScore606of 0.5843, outperforming both the "No Training"607scenario and the Single subset training. This sug-608gests that exposure to more complex, multi-petition609cases helps the model generate richer, more contex-610tually sensitive explanations.

611**Domain-Specific Fine-tuning:** The contrast be-612tween "No Training" and the various training con-613figurations highlights the necessity of domain-614specific adaptation. Without fine-tuning, the

model's explanations remain weak and less aligned with factual inputs, as indicated by lower Rouge and BLEU scores. After training with NyayaFacts Multi, the model better captures the underlying legal rationale, producing explanations that align more closely with reference annotations. 615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

650

651

8 Conclusions and Future Work

We introduced TathyaNyaya, a fact-focused dataset for judgment prediction and explanation within the Indian legal domain, and FactLegalLlama, an instruction-tuned model delivering fact-grounded rationales. By emphasizing factual content rather than full judgments, TathyaNyaya aligns more closely with actual legal decision-making scenarios, while FactLegalLlama highlights the value of coupling predictive accuracy with transparent explanations. Preprocessing steps such as fact filtering and paraphrasing further enhance model clarity and performance, and domain-specific fine-tuning proves essential for capturing legal subtleties. Future work may extend these findings to other jurisdictions, refine fact extraction techniques, integrate, and interpretability frameworks. These efforts collectively advance transparent, accessible, and reliable AIassisted judicial processes.

Limitations

This study faced several limitations that influenced both the scope and outcomes of our research. A key constraint was the reliance on a 4-bit quantized model due to resource limitations, which restricted our ability to experiment with larger parametric models, such as 70B or 40B parameter LLMs. Additionally, the high computational costs and token limitations associated with cloud-based services further hindered our capacity to perform extensive inference and fine-tuning. This restricted explo-

728

702

ration may have limited the depth of insights and performance metrics achievable with FactLegalLlama.

653

654

661

667

670

672

673

674

675

677

681

682

684

688

693

697

700

701

Another significant limitation was the lack of extensive expert evaluation for the generated explanations. While we used high-quality annotations for the TathyaNyaya dataset, resource-intensive processes required for legal expert reviews made it impractical to evaluate the entire dataset. Instead, evaluations were conducted on a smaller subset, which, while insightful, may not fully represent the model's performance across diverse legal scenarios.

The model's performance on scrapped datasets was also not fully evaluated due to configuration constraints, leaving gaps in understanding its generalizability to non-annotated factual data. Furthermore, challenges such as hallucinations in generative outputs and maintaining factual consistency in explanations remain unresolved, which can impact the reliability of the model in real-world legal applications.

Lastly, the dataset used in this study comprises only English-language judgments, which limits its applicability in multilingual contexts, especially in jurisdictions where regional languages dominate legal proceedings. This exclusion highlights the need for more inclusive datasets that reflect the linguistic diversity of legal documents in India and beyond.

These limitations underscore the challenges of applying LLMs to specialized legal tasks such as judgment prediction and explanation. They also point to areas requiring further research, including resource optimization, multilingual dataset development, and enhancing the factual consistency and reasoning capabilities of AI models.

Ethics Statement

This research was conducted with a strong commitment to ethical considerations, particularly given the sensitive nature of legal data and the implications of deploying AI in legal contexts. The TathyaNyaya dataset, central to this study, was compiled from publicly accessible sources, such as Indian legal search engines, ensuring adherence to data privacy and usage regulations. To further safeguard privacy, we removed identifiable metainformation, including judge names, case titles, and case IDs, from the dataset.

The computational resources used for model

training and evaluation were obtained through ethical and legitimate means. These resources were either institutional or subscribed services, ensuring compliance with licensing agreements and financial support for these platforms. By adhering to these practices, we ensured that our research activities aligned with sustainable and lawful resource usage.

Transparency and reproducibility were foundational principles of this study. The TathyaNyaya dataset and the code for FactLegalLlama will be made publicly available, enabling researchers to replicate and extend our findings. This open-access approach is intended to foster collaboration within the research community and drive further advancements in AI-assisted legal decision-making.

We recognize the potential societal impact of AI applications in the legal domain, particularly regarding fairness, accountability, and the risk of misuse. Our models are explicitly designed to assist legal professionals rather than replace human judgment, emphasizing the necessity of human oversight in AI-assisted decision-making processes. As we continue this line of research, we remain vigilant in addressing ethical challenges and aligning our efforts with principles of fairness, transparency, and societal benefit.

References

729

733

736

740

741

743

744

747

748

751

755

763

770

771 772

773

774

775

776

777

778

779

- Nikolaos Aletras, Dimitrios Tsarapatsanis, Daniel Preoțiuc-Pietro, and Vasileios Lampos. 2016. Predicting judicial decisions of the european court of human rights: A natural language processing perspective. *PeerJ computer science*, 2:e93.
- Satanjeev Banerjee and Alon Lavie. 2005. METEOR: An automatic metric for MT evaluation with improved correlation with human judgments. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, pages 65–72, Ann Arbor, Michigan. Association for Computational Linguistics.
- Ilias Chalkidis, Ion Androutsopoulos, and Nikolaos Aletras. 2019. Neural legal judgment prediction in english. Association for Computational Linguistics (ACL).
 - Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. 2024. The llama 3 herd of models. *arXiv preprint arXiv:2407.21783*.
- Yi Feng, Chuanyi Li, Jidong Ge, Bin Luo, and Vincent Ng. 2021. Recommending statutes: A portable method based on neural networks. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 15(2):1–22.
- Yi Feng, Chuanyi Li, and Vincent Ng. 2022. Legal judgment prediction via event extraction with constraints. In Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers), pages 648–664, Dublin, Ireland. Association for Computational Linguistics.
- Saptarshi Ghosh and Adam Wyner. 2019. Identification of rhetorical roles of sentences in indian legal judgments. In *Legal Knowledge and Information Systems: JURIX 2019: The Thirty-second Annual Conference*, volume 322, page 3. IOS Press.
- 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 2022*, pages 3521– 3536, Dublin, Ireland. Association for Computational Linguistics.
- Diederik P Kingma and Jimmy Ba. 2014. Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Chin-Yew Lin. 2004. ROUGE: A package for automatic evaluation of summaries. In *Text Summarization Branches Out*, pages 74–81, Barcelona, Spain. Association for Computational Linguistics.

Vijit Malik, Rishabh Sanjay, Shouvik Kumar Guha, Angshuman Hazarika, Shubham Kumar Nigam, Arnab Bhattacharya, and Ashutosh Modi. 2022. Semantic segmentation of legal documents via rhetorical roles. In *Proceedings of the Natural Legal Language Processing Workshop 2022*, pages 153–171, Abu Dhabi, United Arab Emirates (Hybrid). Association for Computational Linguistics. 783

784

785

786

787

790

791

792

793

794

795

797

800

801

802

803

805

806

807

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

- 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 (Volume 1: Long Papers), pages 4046–4062, Online. Association for Computational Linguistics.
- Shubham Nigam, Anurag Sharma, Danush Khanna, Noel Shallum, Kripabandhu Ghosh, and Arnab Bhattacharya. 2024a. Legal judgment reimagined: PredEx and the rise of intelligent AI interpretation in Indian courts. In *Findings of the Association for Computational Linguistics ACL 2024*, pages 4296–4315, Bangkok, Thailand and virtual meeting. Association for Computational Linguistics.
- Shubham Kumar Nigam and Aniket Deroy. 2024. Factbased court judgment prediction. In *Proceedings* of the 15th Annual Meeting of the Forum for Information Retrieval Evaluation, FIRE '23, page 78–82, New York, NY, USA. Association for Computing Machinery.
- Shubham Kumar Nigam, Aniket Deroy, Subhankar Maity, and Arnab Bhattacharya. 2024b. Rethinking legal judgement prediction in a realistic scenario in the era of large language models. In *Proceedings of the Natural Legal Language Processing Workshop* 2024, pages 61–80, Miami, FL, USA. Association for Computational Linguistics.
- Shubham Kumar Nigam, Navansh Goel, and Arnab Bhattacharya. 2022. nigam@ coliee-22: Legal case retrieval and entailment using cascading of lexical and semantic-based models. In JSAI International Symposium on Artificial Intelligence, pages 96–108. Springer.
- Shubham Kumar Nigam, Shubham Kumar Mishra, Ayush Kumar Mishra, Noel Shallum, and Arnab Bhattacharya. 2023. Legal question-answering in the indian context: Efficacy, challenges, and potential of modern ai models. *arXiv preprint arXiv:2309.14735*.
- Shubham Kumar Nigam, Balaramamahanthi Deepak Patnaik, Shivam Mishra, Noel Shallum, Kripabandhu Ghosh, and Arnab Bhattacharya. 2024c. Nyayaanumana & inlegalllama: The largest indian legal judgment prediction dataset and specialized language model for enhanced decision analysis. *arXiv preprint arXiv:2412.08385*.

- 840 841
- 842 843
- 845 846 847
- 848 849 850 851

- 853 854 855 856 857 858 859 860 861 862 863
- 865 866 867 868 869 870 871 872 873
- 874 875 876 877
- 878 879 880 881
- 8
- 88 88
- 88

892 893

893 894 895

- Shubham Kumar Nigam, Anurag Sharma, Danush Khanna, Noel Shallum, Kripabandhu Ghosh, and Arnab Bhattacharya. 2024d. Legal judgment reimagined: Predex and the rise of intelligent ai interpretation in indian courts. *Preprint*, arXiv:2406.04136.
- Joel Niklaus, Ilias Chalkidis, and Matthias Stürmer. 2021. Swiss-judgment-prediction: A multilingual legal judgment prediction benchmark. *arXiv preprint arXiv:2110.00806*.
- Kishore Papineni, Salim Roukos, Todd Ward, and Wei-Jing Zhu. 2002. Bleu: a method for automatic evaluation of machine translation. In *Proceedings of the* 40th Annual Meeting of the Association for Computational Linguistics, pages 311–318, Philadelphia, Pennsylvania, USA. Association for Computational Linguistics.
- Shounak Paul, Arpan Mandal, Pawan Goyal, and Saptarshi Ghosh. 2023. Pre-trained language models for the legal domain: A case study on indian law. In *Proceedings of 19th International Conference on Artificial Intelligence and Law - ICAIL 2023*.
- Oleg V. Vasilyev, Vedant Dharnidharka, and John Bohannon. 2020. Fill in the BLANC: human-free quality estimation of document summaries. *CoRR*, abs/2002.09836.
- Shaurya Vats, Atharva Zope, Somsubhra De, Anurag Sharma, Upal Bhattacharya, Shubham Nigam, Shouvik Guha, Koustav Rudra, and Kripabandhu Ghosh.
 2023. LLMs – the good, the bad or the indispensable?: A use case on legal statute prediction and legal judgment prediction on Indian court cases. In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 12451–12474, Singapore. Association for Computational Linguistics.
- Chaojun Xiao, Haoxi Zhong, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei Han, Zhen Hu, Heng Wang, et al. 2018. Cail2018: A large-scale legal dataset for judgment prediction. arXiv preprint arXiv:1807.02478.
- Zhilin Yang, Zihang Dai, Yiming Yang, Jaime Carbonell, Russ R Salakhutdinov, and Quoc V Le. 2019.
 Xlnet: Generalized autoregressive pretraining for language understanding. *Advances in neural information processing systems*, 32.
- Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. Bertscore: Evaluating text generation with BERT. In 8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020. OpenReview.net.
- Jieyu Zhao, Yichao Zhou, Zeyu Li, Wei Wang, and Kai-Wei Chang. 2018. Learning gender-neutral word embeddings. In *Proceedings of the 2018 Conference on Empirical Methods in Natural Language Processing*, pages 4847–4853, Brussels, Belgium. Association for Computational Linguistics.

A Experimental Setup and Hyper-parameters

896

897

898

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

916

917

918

919

920

921

922

923

924

925

926

927

928

929

930

931

932

933

934

935

936

937

938

939

940

941

942

943

In this section, we detail the experimental configurations, training procedures, and hyper-parameters employed to develop and evaluate our models. We first describe the training of transformer-based baseline models for fact-based judgment prediction, then outline the instruction-tuning process used to adapt FactLegalLlama for both prediction-only and prediction-with-explanation tasks.

A.1 Transformers Training Hyper-parameters

To establish competitive baselines, we fine-tuned transformer models such as InLegalBERT and XL-Net_Large on the NyayaFacts dataset. Each model was trained with a batch size of 16 using the AdamW optimizer (Kingma and Ba, 2014) and a learning rate of 2e-6. We ran the training for three epochs, adopting default hyper-parameter settings from the HuggingFace Transformers library. Experiments were carried out on an NVIDIA A100 40GB GPU, ensuring adequate computational resources for handling extensive legal text. This training protocol allowed the models to capture the nuances of fact-based segments and reliably predict judicial outcomes.

A.2 FactLegalLlama Instruction Fine-Tuning

To develop FactLegalLlama, we began with the meta-llama/Meta-Llama-3-8B base model. We applied 4-bit quantization to optimize memory usage and introduced Low-Rank Adaptation (LoRA) with a rank of 16 for parameter-efficient fine-tuning. The maximum input sequence length was set to 2,500 tokens, accommodating the substantial factual inputs characteristic of legal documents.

We employed the paged AdamW optimizer in 32-bit precision with a learning rate of 1e-4 and implemented a cosine decay learning rate scheduler for smoother convergence. Mixed-precision training (fp16) and a gradient accumulation of 4 steps were used to further manage GPU memory. We utilized a per-device batch size of 4 and trained the model for three epochs, a process that required approximately 38 hours on an NVIDIA A100 40GB GPU. Under these conditions, the model achieved a training loss of 1.5060 and a validation loss of 1.6745, indicating effective adaptation to the underlying factual patterns in the data.

948

951 952

957

958

959

960 961

962

963

964

965

966

967

968

969

970

972

973

974

975

976

977

979

981

982

984

985

989

991

A.3 Training Objectives

instruction-based The fine-tuning of FactLegalLlama targeted two primary objectives: fact-driven judgment prediction and fact-driven prediction with explanation. By employing a carefully designed set of instructions and incorporating LoRA-based parameter updates, the model learned to generate outcomes and accompanying rationales rooted in the factual segments. This combination of parameter-efficient fine-tuning and instruction-oriented training yielded a model well-suited for practical applications in legal NLP, balancing computational feasibility with interpretability and domain relevance.

A.4 Training Procedure for Hierarchical BiLSTM-CRF Classifier

The Hierarchical BiLSTM-CRF classifier is designed to classify sentences in legal documents into factual and non-factual categories by leveraging the hierarchical structure of the data. The model architecture comprises a word-level BiL-STM coupled with a CRF layer and a sentencelevel BiLSTM. The word-level BiLSTM encodes contextual dependencies within sentences, while the CRF ensures coherence in predicted tag sequences. The sentence-level BiLSTM aggregates these representations to capture inter-sentence dependencies, enabling the model to account for both local and global patterns in the data.

Training is conducted using the AdamW optimizer with a learning rate of 2e-6, a batch size of 16, and for five epochs. A CRF-based loss function is used to optimize sequence-level tagging accuracy. During training, metrics such as precision, recall, F1-score, and loss are evaluated on a validation set after each epoch to monitor performance and ensure generalization. The model configuration includes a word embedding size of 100 and a sentence embedding size of 200, with training conducted on an NVIDIA A100 40GB GPU.

To enhance generalization, K-fold crossvalidation is employed, where the dataset is split into multiple folds, and the model is trained and validated on different subsets. The average performance across folds provides a robust measure of the model's capability. Checkpoints are saved periodically during training, enabling the model to be restored for inference or further fine-tuning.

Template 1 (prediction only) prompt = f^{*****} ### Instructions: Given the facts of the case,just predict the outcome as '1' for acceptance or '0' for rejection. #### Input: <{case_facts}> #### Response: """ Template 2 (prediction with explanation) prompt = f^{*****} ### Instructions: Given the facts of the case,first predict the outcome as '1' for acceptance or '0' for rejection. Then, provide key sentences from the facts or clear reasoning that support your decision. ### Input: <{case_facts}> ### Response: """

Table 6: Prompts for Factual Judgment Prediction and Explanation used for instruction fine-tuned models. Instructions were selected based on the templates provided in Table 8.

Template 1 (Paraphrasing facts)

prompt = f^{*****} **### Instructions**: You are an Indian legal expert with extensive knowledge of legal terms, statutes, and laws. Your task is to explain a legal case to your clients in simple and understandable language. Avoid legal jargon and focus on conveying the meaning of the case in everyday language, making it clear and easy for someone without legal knowledge to understand. While simplifying, ensure that the key points of the case, including the facts, legal claims, and decisions, are clearly communicated without losing any critical information. You should Preserve the key legal terms and references, Clarify complex legal processes, Avoid excessive legal jargon, Be concise but complete, Explain court actions clearly, Provide Only Paraphrased Outcome

Input: Paraphrase the following text:<{case_facts}>
Response: """

Table 7: Prompt for paraphrasing facts to change legal jargons to interpretable terms.



Model Output: Prediction and Rationale Explanation for the Prediction Made

Figure 2: Illustration of the Fact-based Judgment Prediction and Explanation (FJPE) pipeline using the FactLegalLlama model.



Model Output: Prediction and Rationale Explanation for the Prediction Made

Figure 3: Training dynamics of FactLegalLlama for the combined judgment prediction and explanation task. The model learns to produce both the outcome and its underlying rationale directly from factual inputs, guided by instruction-based fine-tuning.



Model Output: Prediction and Rationale Explanation for the Prediction Made

Figure 4: Overview of the simplification and fine-tuning process. First, complex legal facts are paraphrased into simpler language using LLaMA-3-70B, creating the NyayaSimplify dataset, followed by supervised fine-tuning (SFT) using LLaMa-3-7B for the FJPE task.



Figure 5: The Fact vs. Non-Fact segmentation framework employing a BiLSTM-CRF model. This segmentation step separates factual statements from non-factual content in legal judgments, creating the NyayaFilter dataset. The refined dataset is subsequently used for downstream judgment prediction and explanation tasks.

| | Instruction sets for Predicting the Decision |
|----|-------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------------|
| 1 | Analyze the facts presented in the case and predict whether the outcome will be favorable (1) or unfavorable (0). |
| 2 | Based on the facts provided, determine the likely outcome: favorable (1) or unfavorable (0) for the appellant/petitioner |
| 3 | Review the facts of the case and predict the decision: will the court rule in favor (1) or against (0) the appellant/petitioner? |
| 4 | Considering the facts and evidence in the case, predict the verdict: is it more likely to be in favor (1) or against (0) the appellant? |
| 5 | Examine the facts of the case and forecast whether the appeal/petition is likely to be upheld (1) or dismissed (0). |
| 6 | Assess the facts of the case and provide a prediction: is the court likely to rule in favor of (1) or against (0) the appellant/petitioner? |
| 7 | Interpret the facts of the case and speculate on the court's decision: will the appeal be accepted (1) or rejected (0) based on the provided information? |
| 8 | Given the specifics of the case facts, anticipate the court's ruling: will it favor (1) or oppose (0) the appellant's request? |
| 9 | Scrutinize the facts and arguments presented in the case to predict the court's decision: will the appeal be granted (1) or denied (0)? |
| 10 | Analyze the facts presented and estimate the likelihood of the court accepting (1) or rejecting (0) the petition. |
| 11 | From the facts provided in the case, infer whether the court's decision will be favorable (1) or unfavorable (0) for the appellant. |
| 12 | Evaluate the facts and evidence in the case and predict the verdict: is an acceptance (1) or rejection (0) of the appeal more probable? |
| 13 | Delve into the case facts and predict the outcome: is the judgment expected to be in support (1) or in denial (0) of the appeal? |
| 14 | Using the case facts, forecast whether the court is likely to side with (1) or against (0) the appellant /petitioner. |
| 15 | Examine the case facts and anticipate the court's decision: will it result in an approval (1) or disapproval (0) of the appeal? |
| 16 | Based on the facts and evidence in the case, predict the court's stance: favorable (1) or unfavorable (0) to the appellant. |
| | Instruction sets for Integrated Approach for Prediction and Explanation |
| 1 | First, predict whether the appeal in case proceeding will be accepted (1) or not (0), and then explain the by identifying crucial sentences from the document. |
| 2 | Determine the likely decision of the case facts (acceptance (1) or rejection (0)) and follow up with an explanation highlighting key sentences that support this prediction. |
| 3 | Predict the outcome of the case based on the facts provided (acceptance (1) or rejection (0)) and explain your reasoning by extracting key sentences that justify the decision. |
| 4 | Evaluate the case facts to forecast the court's decision (1 for yes, 0 for no), and elucidate the reasoning behind this prediction with important textual evidence from the case. |
| 5 | Ascertain if the court will uphold (1) or dismiss (0) the appeal based on the case facts, and then clarify this prediction by discussing the critical sentences that support the decision. |
| 6 | Judge the probable resolution of the case based on the facts (approval (1) or disapproval (0)), and elaborate on this forecast by extracting and interpreting significant sentences from the case facts. |
| 7 | Forecast the likely verdict of the case (granting (1) or denying (0) the appeal) based on the facts, and rationalize your prediction by pinpointing and explaining pivotal sentences in the case document. |
| 8 | Assess the case to predict the court's ruling (favorably (1) or unfavorably (0)) based on the facts, and expound on this prediction by highlighting and analyzing key textual elements from the case facts. |
| 9 | Assess the case to predict the court's ruling (favorably (1) or unfavorably (0)) based on the facts, and expound on this prediction by highlighting and analyzing key textual elements from the case facts. |
| 10 | Conjecture the end result of the case (acceptance (1) or non-acceptance (0) of the appeal) based on the facts, followed by a detailed explanation using crucial sentences from the case facts. |
| 11 | Predict whether the case will result in an affirmative (1) or negative (0) decision for the appeal based on the facts, and then provide a thorough explanation using key sentences to support your prediction. |
| 12 | Estimate the outcome of the case (positive (1) or negative (0) for the appellant) based on the facts, and then provide a reasoned explanation by examining important sentences within the case documentation. |
| 13 | Project the court's decision (favor (1) or against (0) the appeal) based on the case facts, and subsequently provide an in-depth explanation by analyzing relevant sentences from the document. |
| 14 | Make a prediction on the court's ruling (acceptance (1) or rejection (0) of the petition) based on the case facts, and then dissect the case to provide a detailed explanation using key textual passages. |
| 15 | Speculate on the likely judgment (yes (1) or no (0) to the appeal) based on the case facts, and then delve into the case to elucidate your prediction, focusing on critical sentences. |
| 16 | Hypothesize the court's verdict (affirmation (1) or negation (0) of the appeal) based on the case facts, and then clarify this hypothesis by interpreting significant sentences from the case. |

Table 8: Instruction sets for Prediction and Explanation using factual data from case proceedings.