Bridging Law and Data: Augmenting Reasoning via a Semi-Structured Dataset with IRAC methodology

Anonymous EMNLP submission

Abstract

The effectiveness of Large Language Models (LLMs) in legal reasoning is often limited due to the unique legal terminologies and the necessity for highly specialized knowledge. These limitations highlight the need for high-quality data tailored for complex legal reasoning tasks. This paper introduces LEGALSEMI, a benchmark specifically curated for legal scenario analysis. LEGALSEMI comprises 54 legal scenarios, each rigorously annotated by legal ex-011 perts, based on the comprehensive IRAC (Issue, Rule, Application, Conclusion) framework. 012 In addition, LEGALSEMI is accompanied by a structured knowledge graph (SKG). A series of 014 experiments were conducted to assess the usefulness of LEGALSEMI for IRAC analysis. The experimental results demonstrate the effectiveness of incorporating the SKG for issue identification, rule retrieval, application and conclusion generation using four different LLMs. 021 LEGALSEMI will be publicly available upon acceptance of this paper.

1 Introduction

027

034

035

Access to justice is a universal social challenge. Two-thirds of people in the United States experienced at least one legal issue in the past four years, with less than half of those problems completely resolved ¹. In India, more than 10,490 legal cases in the Supreme Court of India have been pending for more than a decade (Madhana and Subhashree, 2022). These backlogs are often caused by the complexity in legal practice, as well as the scarcity of legal professionals. IRAC framework (Metzler, 2002), stands for issue, rule, application, and conclusion, is the problem solving framework widely used by legal professionals to determine the underlying legal issues, followed by extracting and transforming facts in a legal scenario for legal reasoning, which eventually leads to a legal conclusion.

041

042

044

045

047

048

051

056

057

059

060

061

062

063

064

065

066

067

068

069

071

072

073

074

075

076

077

078

AI models, in particular Large Language Models (LLMs), demonstrate great potentials to improve access to justice (Krasadakis et al., 2024). However, it remains a challenge for LLMs to perform IRAC analysis on legal scenarios accurately. A recent study (Kang et al., 2023) identifies two key problems in analysing legal scenarios. First, Chat-GPT draws wrong conclusions on approximately 50% of the legal scenarios on average. Even if the conclusions are correct, there are mistakes in the intermediate reasoning steps. Secondly, Chat-GPT is not able to cite correct legal rules when analysing majority of the legal scenarios. In real world, it is crucial for legal professionals to understand every single reasoning step that leads to the final conclusion. In addition, our empirical study finds that LLMs struggle to cope with the language gaps between legalese and everyday language. We conjecture that LLMs still cannot fully comprehend the underlying legal knowledge and perform complex legal reasoning accurately.

Recent advances show that it is possible to mitigate the hallucination problem of LLMs by leveraging structured knowledge graphs (SKGs) (Pan et al., 2024). SKGs can enhance LLMs in terms of interpretability and faithfulness by providing external knowledge (Kim et al., 2024). If legal knowledge is stored in SKGs, it is also easy to keep it up-to-date, in accordance with the revisions of legislation. Unfortunately, existing IRAC datasets do not contain any SKGs for legal knowledge.

To address the problems above, we curate LEGALSEMI, a dataset comprising legal scenarios pertaining to Malaysian Contract Law, accompanied by rich structured IRAC analysis carried out by top law students, as illustrated in Fig. 1. Compared to (Kang et al., 2023), we do not only extend their dataset by doubling the legal scenarios in Malaysian Contract Law but also introduce new

¹https://iaals.du.edu/publications/justice-needs-andsatisfaction-united-states-america



Figure 1: An example of a legal scenario pertinent to Malaysian Contract Law with annotations for IRAC analysis. The new types of annotations are legal concepts, court cases, and links to SKGs.

annotation types to all 54 scenarios, including legal concepts and court cases. The inter-annotator agreement across all scenarios exceeds 0.8.

To support reasoning with structured legal knowledge, we extract semantic information from a law textbook and a legislation *automatically* to build the SKG. In the SKG, a node represents either a legal concept, a court case, a legal rule, the interpretation of a legal rule or a concept in lay language, or relevant meta information, while an edge between two nodes denotes their relation. The rigorous layout in the textbook and the legislation facilitates rule-based extraction of semantic relations between legal concepts as well as their relations to legal rules and interpretations. Our extensive experiments reveal the following key findings:

- Following (Kang et al., 2023), we evaluate the capability of LLMs on decomposing a legal question into a set of issues. The key difference to the prior work is the incorporation of legal concepts from the SKG, which improves the quality of issue generation by over 21.4% across all evaluated LLMs.
- By enhancing an LLM with the structured legal knowledge in the SKG, we achieve a 60% increase in recall and a 12% improvement in the F1 score at top-5 results of rule retrieval. We find out that legal concepts are significant in bridging both the language gaps and semantic gaps between facts in scenarios and rules in the legislation. The use of interpretations in lay language further reduces language gaps.
- Our findings indicate that while LLMs are adept at identifying high-level legal concepts, there is still a strong need for improving the quality of recognizing the low-level concepts for supporting neuro-symbolic approaches.

2 Dataset

IRAC provides a comprehensive problem-solving framework for legal professionals. It takes four stages to transform facts acquired from a legal scenario into legal conclusions: (i) identifying legal issues, (ii) determining the legal rules and precedents pertinent to the issues, (iii) performing analysis by applying the law to the facts and the issues, which requires strong legal reasoning skills, and (iv) drawing conclusions based on the analysis. Legal reasoning is *defeasible* such that there are often more than one reasoning traces leading to the same conclusion or different plausible reasoning traces lead to different reasonable conclusions (Billi et al., 2021). Given a legal scenario, legal professionals' concern is not just about the final conclusion, but also why the conclusion is drawn. Therefore, it is essential to build automatic IRAC analysis tools that produce outcomes for each stage and help them identify any missing reasoning steps, and suggest alternative analysis, when necessary.

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

While LLMs demonstrate a great potential to automate IRAC analysis in the absence of supervised training data, they suffer from three key limitations: i) wrong references to statutes and precedents, ii) weak legal reasoning capability, and iii) difficulties in filling the gaps between legalese and everyday language. In contrast, prior studies demonstrate the effectiveness of utilizing Retrieval-Augmented Generation (RAG) and neuro-symbolic approaches with knowledge bases to enhance the factuality and reasoning capability of LLMs (Gao et al., 2023). Therefore, LEGALSEMI builds the first SKG as a legal knowledge base to facilitate research on neuro-symbolic approaches for legal reasoning and provides an annotated corpus to evaluate system outcomes for each stage of IRAC.

111

112

113

114

115

116

2.1 Structured Knowledge Graphs

154

155

157

158

160 161

162

163

165

166

167

168

169

171

173

174

175

176

177

178

179

181

182

183

185

186

190

191

193

195

196

197

198

199

201

Neuro-symbolic systems have garnered increasing interest due to their ability in enhancing the reasoning capabilities of deep neural networks by incorporating symbolic reasoning, such as logic. Recent advances indicate that it is possible to mitigate the hallucination problem of LLMs and enhance the factual accuracy of their responses by incorporating knowledge graphs (KGs) (?). These approaches are considered neuro-symbolic because KGs essentially implement the principles of description logic (Baader et al., 2017).

We consider Malaysian Contract Law as the target area of law due to the importance of contracts in everyday life. The corresponding SKG is automatically constructed from the textbook "Law for Business" (Trakic et al., 2022), the Contracts Act 1950 (the primary legislation governing contracts in Malaysia), and 76 court cases pertinent to contracts downloaded from Malaysia e-judgement ². It is easy to implement rules to extract legal knowledge from legal documents because The layout of a legal document often resembles the structure of legal knowledge, as evident by the screenshots of the textbook in Appendix E.

Legal concepts serve as the building blocks of legal doctrine, often act as bridges that connect related knowledge from diverse sources. For example, under the Contract Act 1950, Section 2(a) states: "when one person signifies to another his willingness to do or to abstain from doing anything, with a view to obtaining the assent of that other to the act or abstinence, he is said to make a proposal;". This section is linked to paragraph P4-014 in the text book via the legal concept "offer".

We derive the skeleton of the SKG from the textbook and enrich the skeleton with statutes from the primary legislation. The index of the book organizes the key concepts of contract law hierarchically, as illustrated in Appendix E. We extract those concepts from the index and annotate them as nodes at the corresponding levels, such as *main_concept* and *subconcept*. The children nodes are linked to their parent nodes using the relation *subconcept_of*. Additionally, we represent each chapter title as a node, indicating specific aspects of the Contract Act 1950, such as *Void Agreements*. Furthermore, we extract the titles, section titles

²e-Judgement: https://cms2.kehakiman.gov. my/CommonWeb/ejudgment/SearchPage.aspx? JurisdictionType=ALL etc. from the Contracts Act 1950 and represent each as a node. Then we introduce several relations to associate the nodes derived from the book with the relevant ones in the legislation. For example, each chapter is associated with the relevant sections of the legislation. Figure 2 shows a snippet of the SKG. 202

203

204

205

206

207

208

209

210

211

212

213

214

215

216

217

218

219

220

221

222

223

224

225

226

227

228

229

230

231

232

233

234

235

236

237

238



Figure 2: A snippet of the SKG.

To bridge the language gap between legalese and plain English, we extract interpretations from the book that provides layman explanations of the corresponding statutes in the legislation. Each interpretation is represented as a node in the SKG, and the *mentions* relation is used to link an interpretation to the relevant statute. Overall, the SKG comprises 3,114 nodes and 1,811 edges, stored in Neo4j for easy data exchange. Further details about the SKG can be found in Table 5 in Appendix C.

2.2 Corpus Construction

From a pool of applicants from Malaysia, we carefully selected six data annotators to assist with scenario selection and annotating IRAC analysis. This annotator team comprises four second-year law students from three distinct Malaysian universities and two junior lawyers. The annotated corpus comprises 54 legal scenarios covering five chapters with 55 subtopics in text book of Malaysian Contract Law. Each scenario reflects real-world legal problems. While the rigorous annotation task takes around three hours per scenario for the IRAC analysis, our easy-to-use annotation tool (Appendix **??**) significantly boosts annotators' productivity.

2.2.1 Scenario Selection

To ensure diversity of scenarios and coverage of legal concepts pertinent to *formation of contracts*, we gather scenarios based on the law textbook "Law for Business" (Trakic et al., 2022) used by law students when studying contract law. In particular, we choose five main topics: *offer and acceptance*,

consideration, certainty, capacity, and intention to 240 create legal relations. The corresponding chapters 241 in the text book are Chapter 4 "Formation of Con-242 tract: Proposal and Acceptance", Chapter 5 "Consideration", Chapter 6 "Promissory Estoppel", and 244 Chapter 7 "Intention to Create Legal Relationships 245 and Capacity". The section headings of these chap-246 ters represent the corresponding subtopics, such as 247 proposal, acceptance, and minors etc.. There are 55 unique subtopics extracted from the subhead-249 ings of the text book.

> First, we asked two annotators to create 24 scenarios which were modified from tutorial questions, books, and past exam questions. Next, for the remaining subtopics, we utilized GPT-3.5 TURBO to suggest candidate scenarios with the prompt : " *You are a legal professional, based on the example scenarios, main topic, and subtopics, create a new scenario around avg_length*". To ensure the quality of the scenarios, another two of the six law students evaluated the quality of the candidate scenarios using the following questions, as shown in 6 in appendix A. Based on annotators' feedback, our experts revised all 54 scenarios to ensure their quality. Their average length is 800 words.

2.2.2 IRAC Annotations

254

258

260

261

265

267

269

271

272

273

274

275

Legal concepts act as a bridge, connecting the facts in a scenario to the professional legal knowledge, including statutes and precedents. We adopted the IRAC analysis based on the annotation guideline in Appendix A

Legal Concepts. Using the legal concepts in the SKG, annotators were asked to identify and high-light relevant legal concepts within a given scenario. If a concept, such as "*offeror*", is not absent from the textbook, they are allowed to add new concepts into the SKG.

277**Issues.** A legal issue is a point of dispute involv-278ing the interpretation, application, or violation of279laws. Six annotators identified issues in scenarios,280focusing on whether a valid contract exists between281parties. The main problem is broken down into spe-282cific questions, such as "Whether there was an283acceptance by Vanessa?" in the example scenario284(Fig. 1).

Rules. A rule specifies the laws applicable to
the issues. The annotators are asked to locate the
appropriate cases and/or statute sections from the
Contract Act 1950 pertinent to issues. For statutory

law, the annotation tool offers a drop-down menu to select relevant sections from the 280 sections available in the SKG. For case law, the tool includes text fields for related court cases along with the corresponding page numbers in the cases. For instance, *Eckhardt Marine GMBH v Sheriff, High Court of Malaya, Seremban & Ors [2001] 4 MLJ 4 (CA) [3/4].* To enable reuse and reference of those rules, the provided cases are displayed as buttons in the user interface so that annotators can refer to those cases by clicking on the buttons (see Fig. 9 in the Appendix D). 289

290

291

292

293

294

295

297

298

299

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

324

325

326

327

328

329

330

332

333

334

335

336

337

Application. In the Application section, annotators applied the rules identified in the rules section to the specific facts of the issues in a given scenario step by step. They are encouraged to use the conditional statements in form of "IF...THEN..." to articulate each reasoning step. Figure 7 in Appendix B illustrates the application section of our example. As legal reasoning is defeasible, annotators can make assumptions due to incomplete information. Different assumptions may lead to different conclusions, thus it is essential to discuss and justify these assumptions in the corresponding reasoning step. The application is the most important and challenging section of an IRAC because it develops the answer to the issue at hand.

Conclusion. The conclusion section directly answers the questions in the issue section, without introducing any new rules and analysis. Following the common practice in law, annotators were asked to write the full sentence of a conclusion, such as "*There is no contract between Vanessa and Niko.*"

2.2.3 Data Quality Assurance

Given a scenario, there are many plausible IRAC analysis, because different assumptions and different interpretations of rules may lead to different conclusions. As it takes roughly three hours to perform a single IRAC analysis and it is infeasible to annotate all possible IRAC, we verified the quality of an IRAC analysis by asking another annotator to act as an evaluator. Specifically, an evaluator can either agree, disagree, or partially agree with an IRAC analysis. The ratio of overall agreements across all 54 scenarios exceeds 0.8, indicating a high level of annotation quality.

Annotator Quality. IRAC analysis is challenging. Hence, as mentioned above, we selected only annotators who have a strong legal background.

		No.	Full	Legal	Dulas	Annotated	SKG
		Scenario	IRAC	Concept	Kules	Application	
	SIRAC	40	yes	0	58	Yes	No
Ì	LEGALSEMI	54	yes	297	90	Yes	Yes
Ì	SARA_entailment	277	no	38	9	No	No
	SARA_numeric	100	no	38	9	No	No
	LEGAL BENCH	59	no	0	18	No	No

Table 1: Comparison between LEGALSEMI and the most relevant datasets.

The law students were required to have achieved at least a B grade in related law subjects. All annotators must pass a specialized pre-test before being recruited. Financial compensation is MYR30 per hour, above the minimum wage MYR7.21 in Malaysia, reflecting the complexity and rigour of the annotation tasks.

2.2.4 Summary of the Corpus

338

339

341

342

344

345

346

347

348

354

363

367

372

376

377

Our corpus comprises 54 scenarios, 243 issues as decomposed questions, 197 mentions of legal concepts (70 of them are unique), 268 sections of the Contracts Act 1950 (44 of them are unique), 76 court cases, and 607 reasoning paths. On average, each application involves 11.25 reasoning steps to draw a conclusion for the main questions. The most common legal concepts encountered include "offeror", "offeree", and "proposal", reflecting the frequent focus on contract formation. Similarly, the law sections most often cited are s2(a), s2(d), s7(a), s2(e), and s2(b).

Dataset Comparison. We compare our corpus with the publicly available corpora in Table 1. Among them, SIRAC (Kang et al., 2023) is the only one that includes annotations of full IRAC analysis for legal scenarios. LEGALSEMI improves upon their work by i) adding a SKG to support neuro-symbolic approaches, ii) introducing annotations of legal concepts to facilitate evaluation of neuro-symbolic approaches, iii) decomposing main legal questions into scenario-based issues, instead of using fixed issues across scenarios as in SIRAC, and iv) include a test set with longer scenarios, closer to the real world scenarios, that require more complex reasoning.

SARA (Holzenberger and Van Durme, 2021) focuses on legal question answering in Taxation Law. It annotates structured reasoning paths involving merely nine rules in total for the QAs but does not include any annotations of IRAC. LEGAL BENCH (Guha et al., 2022) covers diverse legal AI tasks but does not include full IRAC analysis, particularly the detailed analysis in Application.

3 Experiments and Results

We empirically demonstrate the usefulness of LEGALSEMI for IRAC analysis and highlight the open challenges for future research.

381

382

384

385

386

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

3.1 Legal Concept Identification

Legal concepts play a key role in neuro-symbolic approaches for legal reasoning by linking facts in scenarios with legal knowledge. We investigate how accurate the state-of-the-art (SOTA) LLMs can identify legal concepts in scenarios, because LLMs demonstrate remarkable performance on zero-shot and few-shot learning (Brown et al., 2020).

We adopt four LLMs for legal concepts identification: GPT-3.5 TURBO, LLAMA 2, MISTRAL, and GEMINI. The configurations of those models are detailed in Appendix G. Our prompt for those models is shown in Fig. 13. It begins with instructing the LLM to select relevant concepts from a comprehensive list of concept candidates, followed by providing a scenario and main legal questions. At the end of the prompt, it requires the output format to be a Python list for easy post-processing.

We extract a list of legal concepts from each model output, and compare them with the ground truth concepts in terms of precision, recall and F1 score. As the concepts are organized into a hierarchy in the textbook, we report the results for top-level and lower-level concepts, respectively, in order to highlight open challenges.

Models	High-Level Concepts	Low-Level Concepts
GPT 3.5	35%	8%
LLAMA 2	34%	12%
MISTRAL	32%	10%
GEMINI	34%	11%

Table 2: F1 Score for predicting both high-level and lower-level concepts.

Results. As illustrated in table 2, all four LLMs perform significantly better at predicting top-level concepts compared to the lower-level ones. For the top-level concepts, GPT-3.5 TURBO achieves the highest precision (35%), while GEMINI obtains the highest recall (93%). We conjecture that compared to top-level concepts, e.g. "invitation to treat", the lower-level concepts associate with specific details of contracts, such as "audition" and "advertisement". Hence, they appear less frequently in the pre-training data of LLMs. This sheds light on the importance of constructing dedicated supervised training data for future research.

3.2 Issue Identification

422

423

494

425

426

427

428

429

430

431

432

433

434

435

436

437

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

Prior works (Kang et al., 2023; Guha et al., 2022) employ a set of fixed issues to decompose main legal questions into simpler issues. Since issues can vary significantly from case to case in practice, we investigate the extent to which LLMs can generate scenario-based issues and identify the helpfulness of legal concepts at this stage.

We adopt the same four LLMs as legal concept identification. Their prompt is detailed in Appendix H. In the prompt, we instruct LLMs to break a main legal question into a list of issues by leveraging relevant ground-truth legal concepts, and ask LLM to self-evaluate its outputs by ensuring they are *reasonable*, inspired by (Hao et al., 2023).

Evaluation Details. As issue generation is a language generation task, following (Kang et al., 2023), we apply GPT-3.5 TURBO to compare predicted issues with annotated reference issues with a list of criteria detailed in the prompt (see Appendix H). An LLM is expected to select one option from: strongly agree, neutral, or disagree, which is further mapped to a score of 1, 0, and -1, respectively.

To investigate the quality of this automatic metric, we compare the results of the automatic evaluation with those of human evaluation. In the human evaluation, we assess the quality of an IRAC analysis using a rubric that is widely used in Malaysian contract Law courses (Gerhardt, 2008; Carter, 2006). The issues of an IRAC analysis receive a *Pass* when they satisfy the corresponding criteria detailed in Appendix F.1, otherwise, they are marked as *Fail*.

We ask two annotators to independently assess the issues generated by two best performing models (GEMINI and GPT-3.5 TURBO) in three different configurations (e.g. with or without legal concepts) on 10 randomly selected scenarios. If there is a disagreement between their assessments, we ask the most experienced annotator with a strong legal background to resolve it. Each generated output marked as *Pass* receives a score of 2, otherwise, a score of 1. We then rank all LLM configurations according to their average scores and compute the Spearmann rank correlation with the counterpart using the automatic metric. A strong correlation of 0.8 suggests the effectiveness of the automatic metric. Further details are in Appendix F.

470 **Results.** Figure 3 depicts the average scores of
471 the automatic metric across all 54 scenarios. Le-



Figure 3: The results of issue identification.

Decompose Issues

1. Was there a valid contract between Vanessa and Niko?

4. Did Vanessa breach the contract with Niko by selling the vinyl to Ken?

5. What is the legal concept of 'Notice of revocation'?"

Application

(1) According to {Malaysia Contract Law, Section 9(1)}, an advertisement is
generally considered an invitation to treat, not an offer.

(2) IF the advertisement placed by Vanessa was an invitation to treat THEN
Niko's call to reserve that vinyl is an offer THEN Vanessa's reply to
reserve the vinyl for him until Wednesday 8pm is an acceptance. {Section 4.1} {
.....

Figure 4: The LLM outputs of the running example.

gal concepts are beneficial for all LLMs especially for GPT-3.5 TURBO, which increased by 21.4%. The self-evaluation instruction further enhances the performance of all LLMs, with the best performance observed when both legal concepts and self-evaluation are combined. 472

473

474

475

476

477

478

479

480

481

482

483

484

485

486

487

488

489

490

491

492

493

494

495

496

497

498

499

Based on our manual inspection, the incorrect issues generated by LLMs are often either irrelevant to the main issue and the facts of a given scenario, or logically not plausible based on relevant legal knowledge. For example, the generated issue "*Did Vanessa breach the contract with Niko by selling the vinyl to Kenn?*" in Fig. 4 is not pertinent to the main issue regarding whether there is a valid contract between the two parties. Another common error is the generation of issues like "What is the legal concept of 'Notice of revocation'?", which cannot be considered as subissues.

3.3 Rule Retrieval

Given a scenario annotated with legal concepts, we investigate what information in the SKG is beneficial for retrieving relevant legal rules from the Contract Act 1950. A key challenge herein is the gaps between the lay language used in scenarios and the legalese used to express legal rules. Given a scenario as the query, we apply a TF-IDF based search engine (Pedregosa et al., 2011) to retrieve rules indexed by three types of representations, detailed

No initial ratriaval	index: legalese		index: interpret (text book)			index: GPT_interpret			
no initial fettleval	@ top5	@ top10	@ top50	@ top5	@ top10	@ top50	@ top5	@ top10	@ top50
Precision	2.60%	1.70%	1.40%	4.30%	4.90%	7.80%	3.30%	4.40%	3.20%
Recall	2.90%	3.30%	12.50%	0.90%	1.85%	15.70%	2.30%	9.00%	29.40%
F1 score	2.50%	2.00%	2.50%	1.50%	2.54%	9.50%	2.60%	5.50%	5.60%
Initial retrieval		index: law		index:	interpret (te	xt book)	inde	x: GPT_inte	erpret
with legal concepts	@ top5	@ top10	@ top50	@ top5	@ top10	@ top50	@ top5	@ top10	@ top50
Precision	9.70%	7.50%	3.10%	11.80%	13.30%	11.80%	10.30%	9.00%	4.40%
Recall	32.20%	32.60%	37.20%	35.30%	31.20%	35.30%	33.20%	36.50%	48.50%
F1 score	13.90%	11.50%	5.60%	16.30%	17.20%	16.30%	14.60%	13.50%	7.90%

Table 3: Results for rule retrieval. GPT_interpret denotes using the interpretations generated by GPT-3.5 TURBO.

below. Those experimental results are compared with the setting that only the legal rules associated 501 with the same legal concepts pertinent to the sce-502 nario are considered for retrieval. This is achieved by performing retrieval in two stages: i) sending 504 the legal concepts of a scenario as the structured 505 query to the SKG to identify the set of legal rules associated with those concepts, and ii) using the TF-IDF based search engine to rank the legal rules 508 in the results of the initial retrieval. 509

Indexing. We consider three representation types 510 of a legal rule for indexing: i) original legalese, ii) 511 interpretation extracted from the textbook, and iii) 512 combination of the interpretations from the text-513 book and the additional interpretation generated 514 by GPT-3.5 TURBO (detailed in Appendix H.3), 515 because the textbook contains the interpretations 516 of only 18.5% of the legal rules. 517

Evaluation Metrics. We consider precision, re-518 call, and F1 scores at top-k retrieved results, where 519 k = 5, 10, and 50, respectively.

521

527

532

534

Results. As shown in Table 3, the naive approach, which sends scenarios as queries to retrieve rules indexed by legalese, achieves top 5 precision, recall and F1 score less than 3%. Such a low performance 524 is mainly caused by the language gap between lay 525 526 language and legalese. For our running example in Fig. 1, there are few words overlapping between the scenario and the relevant rule Section 2a, which states: "when one person signifies to another his willingness to do or to abstain from doing anything, 530 with a view to obtaining the assent of that other to the act or abstinence, he is said to make a proposal." The word "signifies" can indicate a wide range of 533 actions, such as verbal communication, emails or letters etc.. In contrast, people rarely use "signify" 535 in lay language for the same purpose.

Alternatively, when using the legal concepts as-



Figure 5: Results of application generation.

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

559

561

562

563

564

565

sociated with a scenario as queries to retrieve the rules annotated with the same concepts can significantly boost both precision and recall of rule retrieval. However, different rules may still associate with the same concepts. Such an ambiguity can be further mitigated by interpreting rules in lay language in order to improve similarities between relevant interpretations and scenarios via reranking. If the interpretations in the textbook are used, we achieve the highest recall and F1 score at top-5, despite that only 18.5% of the rules have interpretations. To further improve retrieval quality, we still need to tackle two open challenges in future work. First, although rule interpretations can reduce language gaps, but they are still abstract so that reasoning needs to be performed to associate rules with facts. Second, GPT-3.5 TURBO falls short of generating high quality interpretations so that we need to seek new approaches.

3.4 Application Generation

We investigate the effectiveness of utilizing issues and rules for application generation. We reuse the same four LLMs in the previous stages and apply the prompt (Figure 16 in the Appendix H) to generate legal analysis.

Evaluation Details. Similar to issue generation, we apply GPT-3.5 TURBO to compare the generated application sections with the annotated refer-

ence for each scenario. The possible outcomes of 566 the evaluation include strongly agree (1), neutral 567 (0), or disagree (-1). We also conduct a similar 568 human evaluation to assess the quality of this automatic evaluation metric, and obtain a correlation of 0.86. The details are covered in Appendix F. 571

572

573

574

576

577

578

581

582

586

588

593

594

598

601

611

Results. Figure 5 shows the results when the prompts were incrementally added with issues and rules. Notably, all LLMs greatly benefit from the identified issues except for GEMINI. GPT-3.5 TURBO shows the most significant increase in performance, with an improvement of 18.9% from the self-evaluation prompt to the self-evaluation prompt with issues and rules. Generated applications are effective in improving the quality of conclusions, detailed in Appendix H.4.

Even though we provide a set of correct rules to LLMs, they still produce substantial errors in applying those rules. One of the common mistakes is to associate a wrong rule with facts in a reasoning step. For example, a reasoning step generated by GPT-3.5 TURBO states "According to Malaysia Contract Law, Section 9(1), an advertisement is generally considered an invitation to treat, not an offer." Although the statement "an advertisement is generally considered an invitation to treat" is correct, but it has nothing to do with Section 9(1), even though the content of the rules are provided as a part of model inputs. the other common errors are caused by misunderstanding of rules, reasoning errors, or hallucinations so that logically implausible reasoning steps are generated. For example, in the reasoning step "(3) IF Niko's call to reserve that vinyl is an offer THEN Vanessa's reply to reserve the vinyl for him until Wednesday 8pm is an acceptance, according to Section 4.1. ", both the rule and the derived conclusion is wrong. The correct reasoning should be "According to section 7a, an acceptance must be absolute and unqualified. SINCE Vanessa's response is not absolute and unqualified THEN there is no acceptance.".

Related Work 4

Legal Reasoning Savelka et al. (2023) analyzed how effectively GPT-4 produces definitions for legal terms found in legislation. Huang et al. (2023) addressed the challenge of improving Large Language Models (LLMs), such as LLAMA 2, for 612 domain-specific tasks in the legal field. LEGAL 613 BENCH (Guha et al., 2022) is created through an 614 interdisciplinary procedure for legal scenario anal-615

ysis using the IRAC methodology. However, their work did not utilize the same legal scenarios for the completed IRAC tasks. Large Language Models (LLMs) have demonstrated significant reasoning abilities, especially when chain-of-thought (CoT) prompting (Wei et al., 2022) is employed. Hu et al. (2023) applied LLM to generate a reasoning chain along with the final answer given the legal question. Hao et al. (2023) proposed Reasoning via Planning (RAP). RAP enhances the LLM with a world model and employs principled planning, namely Monte Carlo Tree Search (MCTS), to generate high-reward reasoning traces following effective exploration, demonstrating its superiority over several contemporary CoT-based reasoning approaches. However, these approaches, including RAP, have yet to be applied in the legal domain as artificial intelligence (AI) for legal tasks requires highly domain-specific legal knowledge rather than just common sense knowledge.

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

652

653

654

655

656

657

658

659

660

661

662

663

665

Structured Knowledge Graph SKILL (Moiseev et al., 2022) demonstrated that the results show improvements with pre-trained models on the Wikidata KG, beating the T5 baselines on FreebaseQA, WikiHop, and the Wikidata-answerable subset of TriviaQA and NaturalQuestions. Knowledge graphs with external knowledge can help the model improve accuracy and reduce confusion. Leveraging the power of structured knowledge graphs is able to enhance the performance of the LLMs. For legal reasoning, we need highly specialized legal knowledge to ensure accurate answers. Unfortunately, the current approach mainly focuses on common sense knowledge.

Conclusion 5

We introduce LEGALSEMI, which consists of 54 scenarios annotated with IRAC analysis in Malaysian Contract Law, and a SKG for legal knowledge extracted from a law textbook and legislation. The SKG covers legal concepts, legal rules, interpretations in lay language and their relations. Legal concepts from the SKG are particularly useful for improving the quality of issue generation, which in turn significantly enhance legal analysis in Application across all four evaluated LLMs. Besides, LLMs fall short of identifying relevant legal rules accurately by having the mean precision at top-5 below 3%. By leveraging the SKG, we achieve a significant improvement in rule retrieval, with an increase of 17.2% in F1 score at top-5.

Ethical Statement

666

691

707

712

713

714

Our research practices align with the principles of the ACL Code of Ethics. Our investigation complies with these ethical standards. LEGALSEMI was created and evaluated with a keen awareness of ethical considerations, especially regarding the 671 involvement of human annotators. We recognize 672 that the necessity for human-annotated data to train 673 conditional independence classifiers in our method demands significant effort. We have taken care-675 ful measures to ensure that this process is ethi-676 cally sound, honoring the annotators' contributions 677 by respecting their time and providing equitable compensation. Moreover, the central objective of LEGALSEMI is to create an IRAC methodologybased benchmark. It is designed without generating 681 any information that could be deemed harmful or violate privacy. We will not release the content of the textbook; however, we will provide scripts and detailed procedures to enable those who have purchased the book or hold the copyright to reproduce our experiments.

Limitation

In this study, our primary emphasis revolves around examining scenarios that pertain specifically to the 'Formation of Contract' as delineated within Malaysian Contract Law. While our dataset may exhibit limitations in terms of the breadth of legal scenarios available for analysis, it remains robust in its coverage of all essential topics related to contract formation. Despite potential constraints, such as data availability or accessibility, our dataset is meticulously curated to encompass a comprehensive spectrum of scenarios relevant to the legal domain, ensuring a thorough investigation into the intricacies of contract formation under Malaysian law.

Furthermore, an additional limitation inherent in our study lies in the selection of LLMs employed for our experiments. Our study opts for a more focused approach by utilizing a limited subset of these models. While this decision may result in a narrower scope of analysis compared to studies incorporating a broader array of LLMs, it ensures consistency and reliability in our experimental methodology. Despite this limitation, our choice of employing the most widely used and recognized LLM ensures that our findings are grounded in established practices within the field of natural language processing and legal analysis. 715

References

Franz Baader, Ian Horrocks, Carsten Lutz, and Uli Sattler. 2017. Introduction to description logic. Cambridge University Press.

716

717

718

719

720

721

722

723

724

725

726

727

729

730

731

732

733

734

735

736

737

738

739

740

741

742

743

744

745

746

747

748

749

750

751

752

753

754

755

756

757

758

759

760

761

762

763

764

765

766

767

769

- Marco Billi, Roberta Calegari, Giuseppe Contissa, Francesca Lagioia, Giuseppe Pisano, Galileo Sartor, and Giovanni Sartor. 2021. Argumentation and defeasible reasoning in the law. J, 4(4):897–914.
- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, et al. 2020. Language models are few-shot learners. Advances in neural information processing systems, 33:1877-1901.
- John W Carter. 2006. Carter's guide to australian contract law. (No Title).
- Yunfan Gao, Yun Xiong, Xinyu Gao, Kangxiang Jia, Jinliu Pan, Yuxi Bi, Yi Dai, Jiawei Sun, and Haofen Wang. 2023. Retrieval-augmented generation for large language models: A survey. arXiv preprint arXiv:2312.10997.
- Michael J Gerhardt. 2008. How a judge thinks. Minn. L. Rev., 93:2185.
- Neel Guha, Daniel E Ho, Julian Nyarko, and Christopher Ré. 2022. Legalbench: Prototyping a collaborative benchmark for legal reasoning. arXiv preprint arXiv:2209.06120.
- Shibo Hao, Yi Gu, Haodi Ma, Joshua Jiahua Hong, Zhen Wang, Daisy Zhe Wang, and Zhiting Hu. 2023. Reasoning with language model is planning with world model. arXiv preprint arXiv:2305.14992.
- Nils Holzenberger and Benjamin Van Durme. 2021. Factoring statutory reasoning as language understanding challenges. arXiv preprint arXiv:2105.07903.
- Tongxin Hu, Zhuang Li, Xin Jin, Lizhen Qu, and Xin Zhang. 2023. Tmid: A comprehensive real-world dataset for trademark infringement detection in ecommerce. In Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing: Industry Track, pages 176-184.
- Jie Huang, Xinyun Chen, Swaroop Mishra, Huaixiu Steven Zheng, Adams Wei Yu, Xinying Song, and Denny Zhou. 2023. Large language models cannot self-correct reasoning yet. arXiv preprint arXiv:2310.01798.
- Xiaoxi Kang, Lizhen Qu, Lay-Ki Soon, Adnan Trakic, Terry Yue Zhuo, Patrick Charles Emerton, and Genevieve Grant. 2023. Can chatgpt perform reasoning using the irac method in analyzing legal scenarios like a lawyer? arXiv preprint arXiv:2310.14880.
- Kyungha Kim, Sangyun Lee, Kung-Hsiang Huang, Hou Pong Chan, Manling Li, and Heng Ji. 2024. Can llms produce faithful explanations for fact-checking? towards faithful explainable fact-checking via multiagent debate. arXiv preprint arXiv:2402.07401.

- 771 775 778 783 784 785 786 787 789 796 798 799 803 805
- 810
- 811 812
- 813
- 814

815

816

817

818

819

Annotation Guidelines A

Project Overview Develop a machine learning system for in-depth analysis of legal scenarios, specifically focusing on Contract Law utilising the IRAC (Issue, Rule, Analysis, and Conclusion) methodology.

Panteleimon Krasadakis, Evangelos Sakkopoulos, and

Vassilios S Verykios. 2024. A survey on challenges and advances in natural language processing with

a focus on legal informatics and low-resource lan-

B Madhana and S Subhashree. 2022. A study on back-

Jeffrey Metzler. 2002. The importance of irac and legal

Fedor Moiseev, Zhe Dong, Enrique Alfonseca, and

Shirui Pan, Linhao Luo, Yufei Wang, Chen Chen, Ji-

apu Wang, and Xindong Wu. 2024. Unifying large language models and knowledge graphs: A roadmap.

IEEE Transactions on Knowledge and Data Engi-

F. Pedregosa, G. Varoquaux, A. Gramfort, V. Michel,

B. Thirion, O. Grisel, M. Blondel, P. Prettenhofer,

R. Weiss, V. Dubourg, J. Vanderplas, A. Passos,

D. Cournapeau, M. Brucher, M. Perrot, and E. Duchesnay. 2011. Scikit-learn: Machine learning in

Python. Journal of Machine Learning Research,

Haziqa Sajid. 2023. Leveraging the power of knowl-

Jaromir Savelka, Arav Agarwal, Christopher Bogart,

Laerd Statistics. 2013. Spearman's rank-order correla-

Adnan Trakic, Nagiah Ramasamy, Cheah You Sum,

Paul Linus Andrews, Sri Bala Murugan, P Vijayganesh, and Kanchana Chandran. 2022. Law for

Business, Third Edition. Sweet & Maxwell Malaysia.

Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou,

et al. 2022. Chain-of-thought prompting elicits rea-

soning in large language models. Advances in Neural

Information Processing Systems, 35:24824–24837.

Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten

and Majd Sakr. 2023. Large language models (gpt)

struggle to answer multiple-choice questions about

structured knowledge. Wisecube AI Blog.

code. arXiv preprint arXiv:2303.08033.

tion. Laerd Statistics, page 33.

edge graphs: Enhancing large language models with

Martin Jaggi. 2022. Skill: structured knowledge infusion for large language models. arXiv preprint

writing. U. Det. Mercy L. Rev., 80:501.

log of cases. Issue 5 Int'l JL Mgmt. & Human., 5:942.

guages. *Electronics*, 13(3):648.

arXiv:2205.08184.

neering.

12:2825-2830.

Methodology: Apply Contract Law principles to annotate data using the IRAC framework.

Project Requirements

• Contract Law Expertise: A comprehensive 823 understanding of Contract Law, particularly 824 in relation to contract formation, is essential. 825 You need to have B+ and above for the related 826 subject. 827

822

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

854

855

856

857

858

- · Responsibility and Time Management: Commitment to assigned tasks and timely completion is crucial.
- · Basic IT Knowledge: Familiarity with computer systems and basic IT concepts is preferred.
- · Communication and Teamwork: Strong communication skills and ability to collaborate effectively within a team are important.
- · Pass the pre-test before starting the real annotation work.

Data Annotation Outcomes

- Publication: The annotated dataset will be used for benchmarking and may be published in a journal or presented at a conference.
- Further Research: The annotated data will serve as a resource for subsequent machine learning research.

Benefits

- Research Assistant Experience: Opportunity to work as a Data Annotator on a research project.
- Flexibility: Remote work with flexible hours.
- Compensation: RM 30 per hour.

Annotation Tasks

- · Evaluation of Legal Scenarios: Analyse and evaluate legal scenarios as per the IRAC framework as shown in figure 6.
- IRAC Analysis for Contract Formation: Apply IRAC methodology to analyse contract formation in provided scenarios.
- · Decomposed Questions and Court Case Ref-859 erences: Generate relevant decomposed ques-860 tions for each IRAC segment and include re-861 lated court cases with page numbers. 862

Vanessa, a vinyl store owner, advertised a limited edition vinyl for \$500. On Monday, Niko reserved it but couldn't pay immediately, so Vanessa agreed to hold it until Wednesday &pm. On Tuesday, Ken offered \$700, which Vanessa accepted. She informed Niko via text, but he was out of town and unreachable. When Niko returned on Wednesday at 4pm, he discovered the vinyl was sold and threatened to sue for breach of contract.					
Main Topic: offer and acceptance	Sub Topic: invitation to treat ; proposal ; acceptance; revocation of acceptance				
If the scenario contains the Main topic? (YES/NO/PARTIAL) If the scenario contains the Subtopic? (YES/NO/PARTIAL) The scenario is coherent? (YES/NO/PARTIAL) Accepted with revision? (Accepted with no revision/Accepted with minor revision/Accepted with minor revision/Accepted with no revision acceptance) Details of revision : Revised the scenario					

Figure 6: A scenario with quality assessment questions.

B Examples of the Application annotation.

863

864

872

873

874

877

878

879

Figure 7 exemplifies our annotation process for legal scenario analysis using the IRAC methodology. It demonstrates the structured approach we take to break down and evaluate each aspect of a legal problem. The figure uses logical steps to progress from identifying an initial legal issue to applying relevant rules and statutes, analyzing the facts, and drawing a conclusion. Each step is clearly annotated with references to legal cases and statutes, ensuring that the reasoning is well-supported and transparent. The annotations also include conditional statements and assumptions, highlighting how various legal principles and factual circumstances are considered to reach a final conclusion.



Figure 7: An example of the application section.

Table 4 lists all the condition types used in the annotation. We have a total of six different condition types. The most commonly applied condition type is the IF...THEN... structure.

C Structured Knowledge Graphs

Structured Knowledge Graphs (SKGs) significantly enhance Large Language Models (LLMs) by providing organized, interconnected data representations. This methodical arrangement allows LLMs to make coherent and clear interpretations, aligning seamlessly with their ability to recognize data patterns and relationships. This is particularly beneficial in domains that demand precision, such as scientific research, financial analysis, and medical diagnostics (Sajid, 2023). 883

885

886

887

888

889

890

891

892

893

894

895

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

Legal text often resembles structured knowledge. For example, under the Contract Act 1970, Section 2(a) states: "when one person signifies to another his willingness to do or to abstain from doing anything, with a view to obtaining the assent of that other to the act or abstinence, he is said to make a proposal;" .This section is related to the legal concept "offer" and corresponds to paragraph P4-014 in the text book.

Given the nature of legal knowledge and the benefits of SKGs for LLMs, we design an SKG based on the legal knowledge from book paragraphs, legal concepts, laws, and court cases. The details of the SKG is shown in the Table 5. Figure 2 shows partial of the graph. This graph illustrates the structure of a Social Knowledge Graph (SKG) in Neo4j, showcasing the relationships between various sections, chapters, interpretations, and main concepts through nodes and edges.

Figure 8 illustrates a sections of the SKG, demonstrating the hierarchical and relational structure of legal statutes. The central pink node represents a Chapter, which connects to various Section nodes (yellow) through "BELONGS_TO" relationships. Each Section, like Section with a Title (dark brown) via "HAS_TITLE" relationships and has Interpretations (orange) connected by "HAS_INTERPRETATION" links. These Interpretations, such as P4-109, detail specific provisions and connect to main Concepts (green) and sub-concepts (light brown).

D Annotation Tool

To facilitate this intricate annotation process, we developed an online data annotation platform, grounded in the principles of IRAC methodology. It is designed for universal accessibility, requiring only an internet connection. It features a 'Review' function, allowing annotators to refine and adjust their inputs as necessary. Data output is organized

Conditional Type	Use	Example	Count
If Then	Used to state a condition and its consequence	"IF {She placed an advertisement on a social media platform selling a limited edition vinyl for the price of \$500 [advertisements: invitation to treat]} is an invitation to treat, THEN the call from a customer, Niko, to reserve that vinyl is an offer."	54
According to	Used to refer to legal cases, statutes, or authoritative sources	"ACCORDING TO {Eckhardt Marine GmbH v Sheriff (2001) 4 MLJ 49[53-54]}, IF {Lowel[offeror]} puts up {advertisement[Advertisements to treat]}, THEN it is not an offer but an invitation to treat."	39
Since Then	Used to show a reason and its result	"SINCE {Lowel responded by saying that 500 is too cheap, and that the lowest price she is willing to sell is RM 700 .[Counter offers of initial proposal by proposee]} is not absolute and unqualified THEN there is no valid acceptance."	52
However If Then	Used to introduce an exception or a contrasting condition and its consequence	"{HOWEVER} IF the agreement between Penny and Tina is not supported by considerations THEN it will be void even if it is supported by intention to create legal relation."	42
Even if Then	Used to express a consequence that applies despite a condition	"EVEN IF {Nina[Capacity contracting with]} has lied about her age, she would not be estopped from pleading incapacity."	10
Only if Then	Used to indicate that a consequence will occur solely under a specific condition	IF (7) THEN agreement is supported by intention ONLY IF stated condition is fulfilled. {Confetti Records (A Firm) and others v Warner Music UK Ltd and another [2003] EWHC 1274 (Ch)[]}	12

Table 4: Application conditional types



Figure 8: Details example of SKG in Neo4j.

Node name	Details	No of nodes	Example	
Chapter	Each chapter covers a specific aspect	24	Void Agreements	
Chapter	of the Contract Act 1950.	24		
Title	Each title focuses on specific legal	210	Misrepresentation, Acceptance	
THE	points within the Contract Act 1950.	210	must be absolute	
	Sections of the Contract Act 1950,			
Section	providing detailed legal	204	Section 5.2, An acceptance may be revoked at any	
Section	provisions and serving as the main	504	time before the communication	
	reference for the statute.			
	Interpretations of the contract law,		If a statement does not satisfy the elements of proposal	
Terterment of an	automatically extracted from the	1622	If a statement does not satisfy the elements of proposal	
Interpretation	book content. Each is labeled with	1625	as discussed above, the statement would	
	its content and paragraph IDs.		more likely be an invitation to treat or	
	Footnotes or extensions of the			
Extend content	interpretations, sharing the same	307	The defendent education the	
	paragraph ID as the related interpretation.		The defendant advertised the	
	Key legal concepts summarized from the			
Main concept	interpretations, auto extracted from the	189	proposal revocation	
	index of the book content.			
Sub Concept	Detailed extensions of the main concepts.	351	communication	
Sub sub concept	More detailed information on the sub-concepts.	106	thrid party	
	Total	3114		
Edge name	Details	No of eges	Example	
			chapter_title : OF CONTRACTS, VOIDABLE CONTRACTS	
	Connects the Section and Chapter.		AND VOID AGREEMENTS	
Belongs to		304	content: the court regards it as immoral,	
Delongs_to		501	or opposed to public policy.	
			section_id : 24e	
			title; What considerations and objects are lawful, and what not	
has_title	Connects the Section and Title.	304	Same as above	
			The definition of "agent" and "principal" is provided	
mentions	Connects the Interpretation and Section.	193	in section 135 of;the Contracts Act 1950	
			title "Agent" and "principal"	
related_to	Connects the Interpretation and Extend Content.	307	Hughes v Metropolitan Rly Co (1877) 2 App	
concept_of	Connects the Main Concept and Interpretation.	184	Acceptance concept_of Under section 3, an acceptance can	
subconcept_of	Connects the Main Concept and Sub Concept.	364	communication sub concept of proposal revocation	
subSubconcept_of	Connects the Sub Sub Concept and Sub Concept.	155	third paty sub sub concept of communication	

Table 5: The table illustrates the structure of the Social Knowledge Graph (SKG) implemented in Neo4j. It includes detailed information about the nodes representing entities and edges depicting relationships between these entities. The nodes can represent various entities such as people, organizations, and concepts, while the edges capture the interactions and connections among them. Attributes associated with nodes and edges are also detailed, providing a comprehensive view of the SKG.

980

into a structured .json and ./txt format, significantly
enhancing efficiency and streamlining the data processing workflow for subsequent analysis. Figure
936
9 shows an example of the annotation.

E Textbook details

937

939

940

941

943

947

949

951

953

954

955

957

960

961

962

963

964

965

966

967

968

969

970

971

972

973

975

976

977

Figure 10 showcase various sections and subsections that illustrate the organization of legal knowledge within the textbook. The index of the book's structured format, including headings, subheadings, and bullet points, mirrors the hierarchical nature of legal documents, making it conducive for rulebased knowledge extraction. At the end of the index is a link to the paragraph on that legal concept.

F Human Evaluation

Three human evaluators participate in the evaluation session. We select 10 scenarios and two models (GEMINI and GPT-3.5 TURBO) with all the experiment settings for them to evaluate. We select GPT-3.5 TURBO and GEMINI since it has the best performance from the auto evaluation result. They attend a briefing meeting to discuss and clarify their understanding of the marking rubric. After the briefing, they independently evaluate the scenarios. The third evaluator, who is more experienced, serves as the final decision-maker in cases where the first two annotators disagree, ensuring the reliability and accuracy of the final results. This method follows the steps for identifying issues, which involve decomposing questions for this experiment. The remainder of the evaluation focuses on the application of these guidelines

F.1 Human Evaluation Guidelines

These guidelines are based on the marking rubric and evaluation criteria for contract law (Carter, 2006). They determine how and why the conclusion is made. The answer needs to demonstrate the facts, which include the description of the circumstances, the main questions, and the issues (decomposed questions) raised by the question. It should also cover the rules and principles related to the issues, as well as the conclusions drawn. The proper approach is similar to that of the court's judgments. The grade in figure 11 shows is a helpful step-bystep process to evaluate the legal reasoning process.

978Human Evaluation ResultsHuman evaluation979is conducted for issue identification, application,

and conclusion. These aspects require human judgment based on expertise and credentials to ensure accuracy and reliability in the evaluation process.

We compared the human evaluation results with the auto-evaluation results by examining the rankings of the models' outcomes.

Figure 12 displays the human evaluation results for all experiments. We calculate the Spearman's rank correlation coefficient (Statistics, 2013) for each experiment to assess the alignment between human and automated evaluations.

We rank models based on the 'Pass' rate and compare it to the 'Agree' rate from the automated evaluation matrix. For human evaluation, we use a five-level grading scale: 1 (Fail), 2 (Pass), 3 (Credit), 4 (Distinction), and 5 (High Distinction), and compare these rankings with the 'Agree' rate from the automated evaluation. In conclusion, we compare the entailment rate with the 'Agree' rate from the automated evaluation matrix.

From the result, the average Spearman's rank correlation coefficient, indicated by the red dashed line, is approximately 0.89. This average further emphasizes the overall strong positive correlation between human and automated evaluations across different criteria.

G Models Details

We apply four Large Language Models (LLMs): GPT-3.5 TURBO, LLAMA 2, MISTRAL, and GEM-INI.

In the GPT-3.5 TURBO, our settings include a temperature of 0.7, a common and general setting for GPT models, balancing creativity and coherence in responses. We also set the maximum token count to 1000, allowing for extensive and detailed answers.

We choose LLAMA 2 70B, MISTRAL 7B, and GEMINI as comparative models to analyze their performance against GPT-3.5 TURBO Turbo in handling complex legal scenarios. This comparison aims to assess the efficacy of each model in terms of accuracy, coherence, and relevance to the given legal context. LLAMA 2 is selected for its extensive parameter count, which may enhance its ability to understand intricate details. MISTRAL 7B is known for its efficiency and speed, making it an interesting contrast to the larger models. GEMINI is included for its promising performance in previous legal text analyses, providing a benchmark for evaluation. 1028 By comparing these models, we aim to determine 1029

Example of annotat	ion					
FC1		~				
CHANGE MODE: ANNOTATION Scenario Text						
Vanessa is a vinyl records store o Niko, to reserve that vinyl, Howe a few more days. Vanessa replied customer, Ken, called to purchas his phone was unreachable at th When he came back to collect th if any, to Niko and Ken.	wher. She placed an advertisement on a social media platform selling a limited ec er, Niko told Vanesa that he is unable to make payment at the moment as he is signing that "Time," will reserve the winy! for you until Weindesday Bym. If I don't that viny, he offers \$700 to purchase it in which Vanesa accept. <u>Vanesa then</u> time. Vanesa then sent Niko a text message saying that he sold the winy to avoid to Weindesday at 4pm, he was annoyed to find out that the viny! was sold	Sition vinyl for the price of \$500. On Monday, she receives a call from a customer, faring finan-sit difficulties, he then requested Vanessa to keep the vinyl for him for hear from you by then. Then cit the vinyl for someone dee? On Tuesday, another immediately calls Niko to inform him Then Yau, you was sold to someone else but mone des. Turso u.Niko went out of town and the protocogo and service. Niko now threatens to sue for breach of contract. Advise Vanessa as to her liability.	Sceanrio Sceanrio Highlighted legal			
			concepts			
IRAC Analysis						
Issues	Was there a valid contract be	tween Vanessa and Niko?	→ Main Issue			
Decomposition Questions						
Was the advertisement put of Was there an acceptance on Whether Vanessa was bound Was there a communication/	ut by Vanessa an invitation to treat or an offer? the part of Vanessa? It o reserve the vinyl for Niko until Wednesday 8pm? notice of revocation?		→ Decompose Issues			
Select Related sections	<u>6 selecte</u>	d				
Court cases						
Add courtcase + Remo	ve courtcase 🖥 Generate courtcase buttons		Pulas: Sactions			
No.	Related court case	Paragraph/Page number	→ of Statues and			
#1	Eckhardt Marine GMBH v Sheriff, High Court of Malaya, Sere	3/4	with page			
#2	Preston Corp Sdn Bhd v Edward Leong [1982] 2 MLJ 22 (FC)	2/4	number			
#3	The Ka Wah Bank Ltd v Nadinusa Sdn Bhd [1998] 2 MLJ 350	3/14				
#4	Affin Bank Bhd v Mohd Kasim Ibrahim [2012] MLJU 1794	125/17				
Analysis			Application:			
1. The advertisement placed by V [3/4]}{LEGALLIZATION}.	anessa to sell a limited edition vinyl for \$500 is an Invitation to treat (Eckhardt Ma	arine GMBH v Sheriff, High Court of Malaya, Seremban & Ors [2001] 4 MLJ 4 (CA)	Analysis of the			
2. IF Vanesa's advertisement is an invitation to treat, then (she receives a call from a customer, Niko, to reserve that vinyl.[Niko's reply to the invitation to treat]) is an offer (Section 2a)[Preston Corp Sdn Bhd v Edward Leong [1982] 2. MLJ 22 (FC)[2/4]](LEGALLIZATION).						
3. IF (Fine, I will reserve the viny) acceptance, then there is a valid a	for you until Wednesday 8pm. If I don't hear from you by then, I will sell the vinyl acceptance (Section 7a).	for someone else[Vanessa's reply to the offer]) is an absolute and unqualified	steps.			
Conclusion	Vanessa had breached the contract by selling the vinyl to a third party when s	he was still in an agreement with Niko.	→ Conclustion			
Export File	Daport					

Figure 9: The web-based annotation tool developed to enable the legal scenario analysis.

BECONNICO OF CONTRACT: PROPOSAL CHAPTER I FORMATION OF CONTRACT: PROPOSAL Auto ACCEPTANCE Introduction introduction Proposal, invitation to treat and request for information > Communication of proposal Acceptance > Communication of acceptance Terminating proposals and acceptances > Revocation of acceptances Concluding thoughts

be communicated. There are times when the method of communication is prescribed and if the prescribed method of proposal is not used then the proposal is deemed not to have been communicated (at 432).

[4.050] Mere communication alone is not sufficient for a proposal to be deemed to be communicated. As provided by section 4(1) of the Contracts Act 1950, the communication is only effective when it is brought to the knowledge of the person to whom the proposal is made. A person who has no knowledge of the proposal cannot accept the proposal.⁴³ However, once a person has knowledge of the offer, the motive for accepting the offer is irrelevant.

[4.051] In the Australian case of *Williams v Cawardine*,⁴⁴ a notice was published where reward was promised in return for information given which leads to the apprehension of a criminal in question. The plaintiff in this case had knowledge of the offer of reward for information. However, the plaintiff gave information not motivated by the reward to give the information sought but rather by guilt for her own misconduct with the criminal in question. The Court held that when the plaintiff gave the information sought she had satisfied the conditions of the offer. Her motive in accepting the offer was irrelevant.⁴⁵

ACCEPTANCE

(b)

[4.052] The need for the existence of acceptance in response to a proposal to create a promise is provided for in section 2(b) of the Contracts Act 1950, which provides that:

When the person to whom the proposal is made signifies his assent thereto, the proposal is said to be accepted: a proposal, when accepted, becomes a promise.

[4.053] However, it may be observed that section 2(b) of the Contracts Act 1950 on its own does not explain what characteristics a statement should have in order to be an acceptance. In order to determine whether a particular statement made amounts to an acceptance, reference must be made to spection 7 of the Contracts Act 1950, which provides that:

In order to convert the proposal into a promise, the acceptance must-

(a) be absolute and unqualified



be expressed in some usual and reasonable manner, unless the proposal prescribes the manner in which it is to be accepted. If the proposal prescribes a manner in which it is to be accepted, and the acceptance is not made in that manner, the proposer may, within a reasonable time after the acceptance is communicated to him, insist that his proposal shall be accepted in the prescribed manner, and not otherwise; but, if he fails to do so, he

Abortion

Malaysia, [1.045] United States, [1.041] * The index of the legal concept links tot he paragraph of the text book

Acceptance, [4.052]-[4.063] absolute and unqualified, [4.054] case example, [4.055] acceptor's advantages, [4.119] characteristics, [4.053] communication of, [4.069]-[4.086] completion, [4.077] conditions unfulfilled despite, [4.075] deeming, [4.071] effective, determination, [4.077] English law, [4.069], [4.070], [4.073] general rule, [4.069] Malaysia, [4.070] means, [4.072] omission, acceptance by, [4.072], [4.074] silence, by, [4.073], [4.076] timing, [4.078] case, [4.079] case analysis, [4.080], [4.081] unilateral proposals, [4.075] Contracts Act 1950 s 2(b), [4.052]

Figure 10: The screenshot of the structure of the textbook.



Figure 11: Human Evaluation Guidelines.



Figure 12: Human Evaluation Results.

which is best suited for tasks requiring precise legal1030understanding and reasoning.1031

1032

1033

H Experiment Details

H.1 Legal Concepts prediction tasks

Legal concept prediction experimentFigure 131034displays the structure of the legal concepts predic-
tion. The figure shows the different components of
the prompt. For different experimental settings, we
sometimes remove the legal concepts list from the
potential legal concepts to compare the results.103410351036

Legal Concept EvaluationWe use automatic1040evaluation metrics for this task. The outcome of1041the legal concept list is compared with the ground1042truth. The comparison is separated into two differ-1043ent levels: top-level and lower levels. The top level1044refers to more general concepts, such as "invitation1045to treat." The lower level includes more detailed1046aspects of the concept. For example, under "invi-1047



Figure 13: The prompt for legal concept identification.

tation to treat," there are specifics like "audition," "advertisement," etc.

To evaluate the accuracy of our predictions, we use precision, recall, and F1 score. Precision measures the proportion of correctly identified legal concepts out of all identified concepts. Recall measures the proportion of correctly identified legal concepts out of all relevant concepts in the ground truth, indicating the model's completeness. The F1 score provides a mean of precision and recall, offering a single metric that balances both aspects of accuracy.

H.2 Issue identification task

1048

1049

1052

1053

1055

1056

1057

1058

1060

1061

1064

1065

1066

1067

1068

1069

1070

Issue Identification experiment Figure 14 displays the structure of the issue identification. The figure shows different components of the prompt. For different experimental settings, we sometimes remove the legal concepts list from the ground truth to compare the results.



Figure 14: Details of the experiment of the decompose questions

Issue Identification Evaluation Figure 15 displays the structure of the issue identification evaluation prompt. The evaluation is based on the evaluation guidelines.



Figure 15: Prompts used for the automatic evaluation of the decomposed questions.

H.3 Application Generation

Application Generation experiment Figure 16 displays the structure of the application generation prompt. The figure shows different components of the prompt. For different experimental settings, we sometimes remove the issues or self-evaluation prompt to compare the result.



Figure 16: Application Prompt.

Application Generation Evaluation Figure 17 displays the structure of the application evaluation prompt. The evaluation is based on the evaluation guidelines.

1071

1072

1073

1074

1075

You are a legal professional, Evaluate the provided analysis in comparison to the ground truth based on the evaluation criteria. Assign scores for the analysis and alignment with the corresponding ground truth analysis. a Python list that includes: Illustration: A simplified explanation of the given law in Analysis : {1. Was there a valid contract between Vanessa and the law. Ground Truth: {1. Whether the advertisement put out by offeree was an invitation to treat or an offer?......} Example: Given Law: Evaluation guidelines: Do you agree that the analysis are "A "bailment" is the delivery of goods by one person to Thorough examination of the facts and application of the law to the facts.Consideration of counterarguments or alternative interpretations. Transparency regarding any assumptions made during the analysis. 1: Strongly agree - The document's analysis is logically structured and coherent, comprehensive, and makes reasonable assumptions. It addresses counterarguments and alternative interpretations 0: Neutral - The analysis is logically structured but may lack comprehensiveness or may have some minor logical flaws according Example output : [{ to the document. Illustration: "Every person has a right to enter into any type of -1: Disagree - The analysis lacks clear logical structure. comprehensiveness, or contains major logical flaws as indicated in he document of a party to enter into contracts. The \nreason for such Please indicate the evaluation score for the analysis compared to 1950. \nwhich provides: the ground truth. For instance, if the analysis matches perfectly, assign a score of 1. If it doesn't align, assign a score of -1; if it's somewhat related but lacks specific relevance, assign a score of 0. You only need to show the evaluation result based on this criterion ir the Python list.

Figure 17: Application Evaluation Prompt.

Rule Retrieval: Details GPT-3.5 TURBO interpretation We generate the interpretation of the Malaysia Contract Act 1950 using GPT-3.5 TURBO to interpret the law. The interpretation is then compared with the interpretation from textbooks.Figure 18 display the structure of the prompt. Table 6 shows the example of the output of the interpretation.

Conclusion generation H.4

1082

1084

1085

1086

1087

1088

1089

1090

1091

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

Prompt In the conclusion prompt, we include the scenario, main questions, and the ground truth of the analysis, figure 19 in Appendix H shows the details of the structure. We then ask the LLMs to produce the most reasonable model. The constraint for the conclusion is that it must be in one sentence, answering the main question based on the analysis. According to Kang's work (Kang et al., 2023), providing intermediate steps helps improve the conclusion result. Therefore, we add the ground truth of the reasoning step to obtain the conclusion.

Evaluation Details. The evaluation prompt in-1102 cludes the three scores based on the model output 1103 and ground truth. Each outcome is evaluated with 1104 one of three options: strongly agree (1), neutral (0), 1105 or disagree (-1). These evaluations are based on 1106

Please provide an explanation of a given law with illustration and a real-life example, structured as a Python list containing elements for the law, illustration, and example. The output

common English and explanation to further clarify the law Example: A real-life scenario demonstrating the application of

another for some purpose, upon a contract that they shall. when the purpose is accomplished, be returned or otherwise disposed of according to the directions of the person delivering them. The person delivering the goods is called the "bailor The person to whom they are delivered is called the "bailee".

contract they decide\n to be appropriate for them. However there are some circumstances where\n the law limits the ability limitation usually is based on policy. The requirement \nof capacity may be found in section 10(1) of the Contracts Act Example: "A 16-year-old high school student attempts to take out a loan for a large sum of money to start a business. The bank, upon discovering the student's age, declines to process the loan application, citing the student's lack of legal capacity to

enter into such a contract according to the Contracts Act 1950. This act protects minors from being bound by contracts that they may not fully understand or that may not be in their best nterest."}]

Figure 18: Rule interpretation Prompt.

criteria shown in the figure 20.

Results. Figure 21 in presents the automated eval-1108 uation results for a range of AI models under dif-1109 ferent conditions. The data indicates that GPT-3.5 1110 TURBO and GEMINI demonstrate a more effective 1111 adaptation to the diverse complexities presented in 1112 the experiments. Notably, GPT-3.5 TURBO shows 1113 a significant improvement in the Application & 1114 Self-Evaluation prompt condition, indicating its 1115 potential for enhanced performance in more chal-1116 lenging scenarios. The performance gap between 1117 GEMINI and GPT-3.5 TURBO is relatively narrow, 1118 especially when contrasted with the other models, 1119 LLAMA 2 and MISTRAL, which display lower and 1120 less consistent performances. 1121

1107

1122

1123

1124

1125

1126

1127

Conclusion Generation Figure 19 displays the structure of the conclusion generation prompt. The figure shows different components of the prompt. For different experimental settings, we sometimes remove the application or self-evaluation prompt to compare the result.

Conclusion Generation Evaluation Figure 20 1128 displays the structure of the conclusion evaluation 1129

section_id	content	interpreation	real life example
		A bailment is simply the	John wants to sell his car and
		;transfer of possession and	approaches his friend
	when one person signifies to	control of personal property	Mark with an offer.
	another his willingness to do	(goods) from one	John tells Mark that he
a 2a	or to abstain from doing anything,	;person (the bailor) to another person (the bailee)	is willing to sell his car for
8_2a	with a view; to obtaining the assent	for a specific purpose,	\$10,000. In this scenario,
	of that other to the act or abstinence,	;with the understanding that the goods will b	John is making a proposal by
	he is said to make a proposal;	e returned to the bailor or	expressing his willingness
		otherwise disposed of according to their	to sell his car at a certain price.
		instructions once the purpose is fulfilled.	Mark can either accept or reject the proposal.
			"In a real-estate
			transaction, the seller proposes
	the person making the proposal ; is called the "promisor" and the person accepting the proposal is called the "promisee";	In a legally enforceable promise, there are	to sell their property to the buyer at a
		two parties involved. The person who	specific price. The buyer accepts this offer,
s_2c		makes the offer is known as the 'promisor' and	resulting in a legally binding
		the person who accepts the offer is	contract between them. In this scenario,
		referred to \as the 'promisee.'	the seller is the 'promisor' because
			they made the offer, and the buyer is the
			'promisee' because they accepted the offer.

Table 6: Example of the interpretation.



Figure 19: Conclusion Prompt.

1130prompt. The evaluation is based on the evaluation1131guidelines.

1132Conclusion ResultFigure 21 display the result1133of the conclusion. The comparison is between1134adding the application and not adding the appli-1135cation in the prompt.



Figure 20: Conclusion Evaluation Prompt.



Figure 21: Result of Conclusion.