
Building LLMs for Legal Reasoning

Elven Shum

CPSC490 Project

Advisors: Ruzica Piskac, Scott Shaprio, Timos Antonopoulos

elven.shum@yale.edu

1 Introduction

1.1 Overarching Motivation:

Large Language Models (LLMs) have emerged as powerful tools across various domains, including the legal sector, where they promise to enhance and streamline legal processes. However, despite their utility in basic legal tasks, these models often struggle with more intricate legal reasoning that involves specialized terminologies and complex legal principles. Trained primarily on general language corpora, conventional LLMs tend to underperform when faced with texts dense in legalese, revealing a gap in their capacity to interpret and apply nuanced legal constructs. [5]

Our project seeks to address this gap by not merely fine-tuning an LLM for improved legal reasoning, but by strategically utilizing prompt engineering and chain-of-thought reasoning to elicit more accurate and contextually grounded outputs. Rather than solely focusing on fine-tuning, we emphasize techniques that refine the model's interpretative capabilities—pushing the boundaries of how LLMs process legal predicates and reach conclusions. Collaborating with Yale Law School professor Scott Shapiro and Yale Computer Science professor Ruzica Piskac, we concentrate on proximate cause cases, as detailed in Section 3.1, to develop a robust framework that enhances the model's ability to handle sophisticated legal reasoning tasks. This approach aims not just for better performance, but for more meaningful decision-making, where the model's reasoning mirrors the critical thinking demanded in real-world legal contexts.

1.2 Content of Report, Distinction Between 488 and 490:

Please note this project is separated into two parts, divvied up between CPSC488's final project and CPSC490. The CPSC488 component of this project consists of the following components: Scraped Raw Legal Case Data, Data Optimization Scripts, Optimized Proximate Cause Dataset. The CPSC490 component consists of (1) the Fine-Tuned Model, (2) Fine-Tuning Scripts + Chain of Thought Prompts, (4) Automated Evaluation Script (3) Real World Case Verdict Prediction Evaluation, (4) Proximate Cause Reasoning Evaluation. As will be discussed in later sections, these evaluations are done by both an LLM and by the author, who is not legally trained. Proper holistic evaluation requires professional evaluation from a legal scholar. Please note that section 3.1-3.3 are intended to be in the CPSC488 report, but are included in this CPSC490 report for better context and continuity; to not offend university dual submission requirements, sections 3.1-3.3 are not considered part of the CPSC490 component, and are merely for reference.

2 Related Work

Recent advances in integrating LLMs with legal reasoning have led to the development of several benchmarks and models. LegalBench remains the most comprehensive benchmark, testing LLMs on tasks ranging from issue-spotting to rule application, offering a thorough evaluation of models' performance in practical legal scenarios [5]. Complementing this, LawLLM was developed specifically for

36 U.S. legal tasks, improving the prediction of verdicts, case retrieval, and precedent recommendation
37 through advanced fine-tuning techniques [4].

38 ADAPT, a prompting framework, enhances legal judgment prediction by focusing on multi-task
39 instruction tuning. This framework tackles limitations in traditional LLMs by refining models through
40 specific legal datasets, thus increasing reasoning accuracy across complex legal charges [3]. Similarly,
41 Chain of Logic has demonstrated improved performance over standard methods by applying a rule-
42 based approach to tasks like jurisdiction and contract interpretation, further showcasing the role of
43 specialized prompting techniques in legal applications [1].

44 Beyond the U.S., DISC-LawLLM focuses on legal reasoning for the Chinese legal system, utilizing
45 fine-tuned models to navigate local laws and legal documents. This highlights the growing importance
46 of adapting LLMs for jurisdiction-specific contexts [2]. Additionally, LegalGPT introduces a multi-
47 agent framework for legal reasoning, leveraging chain-of-thought reasoning to enhance interpretative
48 capabilities across diverse legal tasks [6].

49 As the field evolves, domain-specific models such as SaulLM are showing superior performance com-
50 pared to general-purpose models like GPT-4. These legal-specific LLMs provide deeper contextual
51 understanding and higher accuracy in tasks like document review and verdict prediction, underscoring
52 the potential of tailored models in the legal sector [7].

53 **3 Methodology**

54 **3.1 Dataset Selection:**

55 Our project revolves around "Proximate Cause" cases in law. These cases are emblematic of
56 systematic legal reasoning, use explicit categorizations of necessary types of legal evidence, and are
57 minimally semantically interpretative. Other legal domains, for example constitutional law cases, would
58 be ill suited to applying LLMs, and they're almost entirely devoted to interpreting particular words
59 in a modern context. The systematic nature of proximate cause cases, coupled with their relatively
60 concise length, renders them apt for the operational scope of LLMs. The structured methodology
61 employed in teaching these legal concepts to law students is rather identical; typically, students are
62 tossed a plethora of cases from a large book and expected to determine the abstract reasoning. This
63 further underscores the feasibility of instilling a refined understanding of these concepts in LLMs
64 through fine-tuning. In the USA, different States, while internally consistent, each possess slightly
65 varying laws regarding their handling of Proximate Cause cases. Therefore, we select our cases from
66 only the state of Texas. Furthermore, in order to get a wide variety of Proximate Cause cases, we
67 select from several intervening causes: Lightning, Disease, Criminal acts, Rescue, Suicide, Medical
68 treatment. These keywords were suggested by Professor of Jurisprudence, Scott Shapiro.

69 **3.2 Dataset Acquisition and Cleaning:**

70 For exhaustive dataset, we leveraged the resources of the Yale Law School library, courtesy of
71 Professor Scott Shapiro. Our data acquisition was facilitated through YLS via LexisNexis, a legal
72 database provider. Legal case repositories are not typical pre-filtered or categorized. Therefore 3rd
73 party Legal Databases, where the categorization has already been completed, is an indispensable tool
74 for legal research, like ours. Cases retrieved from LexisNexis are only in Rich Text Format (RTF),
75 laden with unnecessary metadata and formatting elements.

76 To address this, we developed a data cleaning script, primarily focused on two tasks: converting
77 RTF files to a more manageable plain text (txt) format and excising irrelevant metadata, d formatting,
78 headers, and footers. The resultant text files were then formatted to align with the requirements of our
79 Data Optimizer (discussed in the subsequent section). This streamlined dataset, devoid of extraneous
80 information, ensured that the data fed into the model was of consistent quality.

81 **3.3 Dataset Optimization:**

82 To enhance the proficiency of a Language Model (LM) in not only comprehending legal reasoning
83 but also in deducing verdicts from case facts, our research focused on fine-tuning the model using
84 Proximate Cause legal reasoning cases. It became evident that the utilization of raw case data
85 was suboptimal for several reasons. Primarily, these cases often exceed the limits of many context

86 windows due to their length. Additionally, they contain extraneous elements such as judge names and
87 case references. Most critically, the ingestion of raw cases in their entirety might not effectively train
88 the LM in legal reasoning; it could instead result in a superficial replication of verdict-like responses.
89 The objective was to enable the LM to internalize the underlying logic and reasoning of the cases,
90 necessitating an understanding of the Legal Reasoning Rule and the ability to distill relevant facts
91 from a case for rule application.

92 To foster this advanced learning, significant Dataset Optimization was imperative. We employed
93 a preprocessing LM, ChatGPT-4, to transform the text of raw cases into structured (Background,
94 Verdict) pairs. Recognizing that legal cases inherently comprise both background facts and judicial
95 rulings, our training approach differed; for the fine-tuned LM, we provided only the background,
96 tasking it with verdict prediction. This necessitated the creation of (Background, Verdict) pairs
97 specifically for training, which were then used to prompt the LM to generate verdicts based solely on
98 case backgrounds. These pairs were designed with key attributes in mind: brevity compared to the
99 full-length cases, inclusion of only factual elements in the Background, and ensuring that the Verdict
100 not only elucidated the reasoning but also exclusively referenced facts presented in the Background.
101 Through iterative prompt engineering, we identified the phrase *"Please summarize the following case
102 in the form { Background: [background of the case], Verdict: [verdict of the case] }. In Background,
103 do not include the cases's outcome, only facts. Be extremely detailed in the Background, as any
104 facts could be relevant. In Verdict, explain reasoning for the verdict in future tense, as if you were
105 predicting. Do not cite facts in the Verdict that were not in Background. Be very detailed in Verdict."*
106 that yielded best results.

107 3.4 Fine-tuning Methodology: OpenAI api

108 To fine-tune the model on proximate cause legal reasoning, we utilized the OpenAI API. The decision
109 to use this API was driven by several factors. First, our goal was to produce the most accurate legal
110 reasoning model possible, specifically targeting proximate cause cases. Using state-of-the-art (SOTA)
111 models was crucial to ensure the highest chance of capturing the nuances in legal reasoning. Open-
112 source SOTA models were considered; however, the computational costs associated with self-hosting
113 these models were prohibitive, given our hardware limitations. Thus, the OpenAI API, particularly
114 the GPT-4o model, was selected for fine-tuning due to its robust performance in non-mathematical
115 reasoning tasks.

116 We fine-tuned the model on a dataset of 440 legal cases related to proximate cause. The dataset was
117 split into 95% for training and 5% for evaluation, resulting in 22 evaluation cases. This split was
118 designed to test the model's ability to generalize legal reasoning, ensuring that it could pick up the
119 subtle complexities of proximate cause reasoning in unseen cases. The proof of concept aimed to
120 determine whether such models could accurately replicate legal reasoning with minimal data.

121 3.4.1 Chain of Thought System Prompt

122 We crafted a Chain of Thought (CoT) system prompt to structure the model's reasoning during
123 fine-tuning. It's intentionally rather general to encourage the model to run through a step-by-step
124 analysis of legal facts and the principles of proximate cause, while not restricting it on-high to a
125 proximate cause framework. Ideally, the model learns from the training cases *how* to reason about
126 proximate cause, as opposed to directly telling it. In this way, it may be able to capture more general
127 principles than what a rigid CoT prompt would've permitted. See section 5 for discussion on potential
128 creation of a more ideal prompt. The following prompt was used:

129 *You are an expert in law with a focus on reasoning about proximate cause. Your*
130 *goal is to carefully analyze the background of legal cases and determine the*
131 *correct verdict using step-by-step legal reasoning.*

132
133 *To reach a conclusion, follow these steps:*

- 134 *1. Identify all relevant facts from the background.*
- 135 *2. Analyze these facts in the context of proximate cause, highlighting key points of*
136 *causation and foreseeability.*
- 137 *3. Consider whether the defendant's actions could have reasonably led to the*
138 *plaintiff's harm.*

- 139 4. Discuss any potential intervening causes, ambiguities, or contributing factors.
140 5. Synthesize these considerations into a coherent argument, and clearly state
141 the final verdict in the last line of your reasoning. The verdict must be presented
142 without additional quotes or formatting.

143 3.4.2 Fine-tuning Setup

144 In our setup, the case backgrounds were presented to the model as the user input, while the verdict was
145 framed as the assistant response. This structuring choice reflects the natural division between the legal
146 facts (input) and the legal decision (output), closely mirroring real-world legal reasoning workflows.
147 By using this approach, the model learns to process background information, perform legal analysis,
148 and predict a verdict, effectively simulating the role of a legal expert. This approach, combined with
149 the CoT prompt, ensured the model was exposed to a clear, step-by-step legal reasoning process
150 during training. We trained for 3 epochs.

151 4 Results and Evaluation

152 4.1 Evaluation Setup

153 Our evaluation was designed to assess two core objectives: (1) the accuracy of the fine-tuned model’s
154 verdict predictions in proximate cause cases, and (2) the quality of its legal reasoning compared to
155 the base model, GPT-4o. We benchmarked the fine-tuned model against GPT-4o using a set of 22
156 evaluation cases, which were randomly selected from our dataset. These cases were chosen to test
157 whether the fine-tuned model could effectively generalize legal reasoning to unseen cases.

158 To evaluate both objectives, we implemented the following approaches:

- 159 1. **Verdict Prediction:** We compared the verdicts predicted by GPT-4o and the fine-tuned
160 model to the ground truth verdicts, which were based on real-life case outcomes.
- 161 2. **Reasoning Quality:** We used GPT-4o to evaluate the reasoning provided by both models,
162 determining which model’s reasoning best aligned with the legal standards of proximate
163 cause.

164 While comprehensive legal evaluation requires input from trained legal professionals, this proof of
165 concept aims to demonstrate that accurate legal reasoning can be learned. With an intentionally
166 limited evaluation set of 22 cases, we focused on whether the model could capture the core principles
167 of proximate cause reasoning, despite the absence of expert reviewers.

168 4.2 Evaluation Process

169 We automated the evaluation process by comparing both models’ verdicts to the ground truth. First,
170 we assessed whether each model’s predicted verdict agreed with the ground truth. For this, we utilized
171 GPT-4o to determine if the predicted and actual verdicts were in agreement. Next, we compared
172 the reasoning quality between the models. GPT-4o was used to evaluate three verdicts per case: the
173 ground truth, GPT-4o’s verdict, and the fine-tuned model’s verdict. The model then judged which of
174 the two predictions (GPT-4o or fine-tuned) best aligned with the reasoning in the ground truth verdict.

175 This process generated a CSV file containing detailed results, including verdict agreement and
176 reasoning quality. This file serves as the primary deliverable for comparing the models’ performance
177 across all evaluation cases.

178 Subtle adjustments were made to ensure that the base model’s system prompt provided minimal
179 guidance, whereas the fine-tuned model’s prompt encouraged a structured, step-by-step analysis. This
180 allowed us to maintain consistency while still evaluating the effects of fine-tuning on legal reasoning.

181 4.3 Results

182 In our evaluation of 22 cases, we observed a 20% increase in verdict accuracy with the fine-tuned
183 model compared to the base GPT-4o model. Specifically, the fine-tuned model correctly predicted
184 verdicts in 18 cases, while the base model succeeded in 14. Though the dataset is small, the results

185 indicate that fine-tuning increases the likelihood of reaching the correct verdict in real-world legal
186 cases.

187 More importantly, the fine-tuned model demonstrated a significantly improved ability to reason
188 through the complexities of proximate cause. In head-to-head comparisons of legal reasoning, our
189 model aligned with real-life court reasoning in 15 of the cases, while GPT-4o matched in 7. This
190 difference highlights the fine-tuned model’s enhanced capacity to handle complex legal concepts and
191 doctrines. The fine-tuned model can be accessed via the OpenAi API as model:

192 `ft:gpt-4o-2024-08-06:yale:AAynVGq6`

193 4.4 Case Study: Verdict Comparison

194 To illustrate the effectiveness of our fine-tuned model, we compare the verdicts generated by the base
195 GPT-4 model and our fine-tuned model against the ground truth verdict of a legal case. Table ??
196 presents selected quotations from each verdict, highlighting the similarities between the ground truth
197 and the fine-tuned model.

Table 1: Comparison of Verdicts (Eval Case 20)

Key Point	Ground Truth Verdict	Base Model Verdict	Fine-Tuned Model Verdict
Applicability of Chapter 95	Weekley did not conclusively show that the Chapter 95 protections were applicable.	Unless appellants can produce compelling evidence, the initial motion for summary judgment by Weekley Homes might stand.	The appellate court will reverse the trial court’s judgment and remand the case based on the existence of genuine issues of material fact regarding both prongs of Chapter 95.
Error in Granting Summary Judgment	The trial court’s granting of the traditional summary judgment based on Chapter 95 was erroneous.	The trial court likely granted summary judgment assuming lack of substantial evidence.	The decision will be based on the existence of genuine issues of material fact.
Negligence and Premises Liability Claims	Appellants presented more than a scintilla of evidence demonstrating a duty owed and the possibility of a nexus between their injuries and the control of their working conditions by Weekley.	Without clear evidence showing Weekley’s direct control, appellants may struggle to prove control.	Evidence could show that Weekley had control over the timing and conditions of appellants’ work.
Denial of Motion for Reconsideration	Denial of appellants’ motion for reconsideration and new trial was an abuse of discretion by the trial court.	Unless appellants can produce compelling evidence, the trial court’s summary judgment might stand.	The trial court’s denial of the motion for reconsideration was an erroneous decision warranting reversal.
Gross Negligence Claim	Appellants failed to provide substantive evidence regarding Weekley’s actual, subjective awareness of risk.	Appellants would need robust evidence proving Weekley’s actual knowledge.	The court will likely affirm the trial court’s verdict on the gross negligence claim.
Final Judgment	The appellate court affirmed the trial court’s verdict on the gross negligence claim, reversed the verdict on the negligence and premises liability claims, and remanded those parts back to the trial court.	Thus, if we assume appellants could not conclusively show these elements, the summary judgment would be upheld.	The appellate court will reverse the trial court’s judgment and remand the case for further proceedings.

198 The fine-tuned model’s verdict closely mirrors the ground truth in both outcome and legal reasoning,
 199 demonstrating its enhanced capability to understand and apply complex legal principles. Both the
 200 ground truth and the fine-tuned verdict recognize that Weekley didn’t conclusively establish the
 201 applicability of Chapter 95 protections due to genuine issues of material fact concerning control and
 202 knowledge of the dangerous condition. Also, fine-tuned model accurately identifies the trial court’s
 203 error in granting summary judgment on the negligence and premises liability claims, therefore the
 204 appellants presented sufficient evidence to raise genuine issues for trial. It also concurs with the
 205 ground truth in affirming the trial court’s verdict on the gross negligence claim: lack of substantive
 206 evidence regarding Weekley’s actual awareness of risk.

207 4.5 Case Study 2: Verdict Comparison

208 To further demonstrate the effectiveness of our fine-tuned model, we present a second case study
 209 comparing the verdicts generated by the base GPT-4o model and our fine-tuned model against the
 ground truth verdict. Table 2 similarly showcases selected quotations from each verdict.

Table 2: Comparison of Verdicts (Eval Case 14)

Key Point	Ground Truth Verdict	Base Model Verdict	Fine-Tuned Model Verdict
Affirmation of Judgment	The trial court’s judgment was upheld.	The appellate court should consider remanding for a new trial.	The appeals court would likely affirm the trial court’s rulings.
Limitation of TDHS Records	Limitation of TDHS records to those concerning Watson was within the court’s discretion.	The trial court excluded these records, likely to prevent prejudice.	The court will opine that any error in excluding other staff supplements was harmless.
Exclusion of Dr. Cox’s Testimony	Restrictions on the admissibility of Dr. James Cox’s testimony were deemed non-arbitrary as his conclusions were speculative.	The restriction on Dr. James Cox’s testimony could imply the trial court unnecessarily limited crucial evidence.	The court will conclude that the exclusion of certain expert testimonies, including those involving Dr. Cox, did not likely cause harm to appellants.
Confidentiality of TDHS Photographs	The TDHS photographs were considered confidential and not subject to disclosure.	If the photographs were highly relevant and non-prejudicial, their exclusion might constitute an error.	The court will find that appellants did not preserve their complaint concerning TDHS photographs for appeal.
Negligence Per Se Claims	Negligence per se claims were correctly struck as the cited laws did not provide for tort liability.	If the trial court erred in this determination, it may affect the basis upon which the jury considered negligence.	The court will argue that appellants lack standing because the NHRA does not confer an independent private right of action.
Overruling Appellants’ Issues	The appellants’ issues were overruled in their entirety, and the take nothing judgment against them was affirmed.	The appellate court should consider remanding for a new trial.	Predicting the outcome, the court will express that all points of error by appellants are overruled and the trial court’s judgment is affirmed.

210
 211 Once again, the fine-tuned model’s verdict demonstrates a remarkable alignment with the ground
 212 truth, both in the final judgment and the underlying legal reasoning. It accurately reflects the appellate
 213 court’s decision to affirm the trial court’s rulings, where by capturing the nuances of the court’s
 214 discretion and the application of legal doctrines. The fine-tuned model recognizes that any error in
 215 excluding additional TDHS records was deemed harmless, which mirrors the ground truth’s emphasis
 216 on the court’s discretionary power and the lack of substantial prejudice resulting from such exclusions.

Also, the fine-tuned model correctly identifies the appellants lack of standing to bring negligence per se claims, aligning with the ground truth’s reasoning that the cited statutes did not confer an independent private right of action. It also acknowledges procedural intricacies, such as the appellants’ failure to preserve certain complaints for appeal and the insufficiency of objections during the trial, which the ground truth verdict highlights as critical factors in the appellate court’s decision.

4.6 Results Conclusion

The fine-tuned model shows significant progress in handling complex legal reasoning tasks, particularly in predicting verdicts for proximate cause cases. Its ability to produce verdicts that closely align with the ground truth demonstrates a promising understanding of legal principles. For instance, the model’s handling of intervening causes like "acts of God" and foreseeability mirrors the nuanced reasoning observed in actual court decisions. In one case, the model’s verdict on the applicability of Chapter 95 protections accurately captured the key issue of control, while in another, it correctly identified the trial court’s error in granting summary judgment on negligence claims—both examples highlighting the model’s deepened understanding of proximate cause.

Our approach combining a broad Chain of Thought (CoT) prompt with targeted fine-tuning has proven effective in fostering structured, step-by-step legal reasoning. Most compellingly, the fine-tuning process allowed the model to internalize legal primitives such as causation and foreseeability, which are fundamental to proximate cause cases. The success of these techniques suggests that the model is not just replicating outcomes, but genuinely learning to apply core legal concepts. This opens up exciting possibilities for future work in expanding its capabilities across different legal domains.

5 Discussion

Our results suggest that the fine-tuned model has taken significant steps toward successful legal reasoning in proximate cause cases. The fine-tuning process, coupled with a broad Chain of Thought (CoT) system prompt, allowed the model to learn key reasoning primitives, such as causation and foreseeability. This facilitated the model’s ability to predict verdicts with reasoning that aligns closely with the ground truth in many instances. However, the model’s verdict success rate, while improved over the baseline GPT-4o, remains modest, and its performance could be bolstered with further refinements and more data.

One of the key limitations of this work is the relatively small dataset of 450 proximate cause cases. Although this was sufficient for the model to internalize fundamental legal reasoning structures, it likely missed more nuanced or esoteric legal principles that emerge in particularly intricate cases. Additionally, while the fine-tuned model outperformed the baseline in reasoning tasks, its overall verdict accuracy was only slightly better, demonstrating that there is still room for substantial improvement. Furthermore, GPT-4o is already a very strong reasoner, so we are starting from a high baseline. Future efforts should aim to push the model to a level where it can more consistently predict correct outcomes across a wider variety of cases.

In terms of future work, while our broad CoT approach successfully encouraged structured reasoning without explicitly defining proximate cause rules, there is a case to be made for including more domain-specific legal reasoning rules directly in the CoT system prompt. A well-structured prompt based on actual legal reasoning frameworks might yield short-term improvements in verdict prediction, particularly if the model had more cases to train on. We intentionally avoided this due to the limited scope of our evaluation cases and the desire to let the model learn these rules autonomously. However, with a larger dataset, incorporating a more structured legal framework might be an avenue worth exploring.

Additionally, one of the largest areas for improvement lies in the "Case Optimization" process. Our data was optimized using base GPT-4, rather than more advanced models like GPT-4o or newer reasoning-optimized models from OpenAI. With stronger models, it’s likely that the case summaries and reasoning processes could be much more robust, leading to better downstream performance in legal reasoning tasks. We recommend exploring new model architectures to generate more detailed reasoning tokens, which could then be summarized by another model. This iterative reasoning and summarization approach might unlock more sophisticated legal analysis. Unfortunately, we were

268 constrained by the release timeline of these newer models, but refining this data optimization step
269 should be the focus of future improvements.

270 From a broader perspective, this work demonstrates the potential for AI systems to engage in serious
271 legal reasoning. While current state-of-the-art (SOTA) models have not been explicitly trained for
272 this purpose, our fine-tuning approach shows that it is possible to tailor models for legal tasks like
273 proximate cause reasoning. In the future, such systems could significantly enhance the efficiency of
274 legal professionals by generating structured analyses as a first pass, freeing up lawyers to focus on
275 higher-order decision-making rather than tedious case review.

276 Finally, I would like to extend my sincere thanks to my advisors and mentors: Professors Ruzica
277 Piskac, Scott Shapiro, and Timos Antonopoulos. Their unrelenting patience and guidance throughout
278 this process were invaluable, especially as my health issues caused significant delays in the completion
279 of this project. I am grateful for the opportunity to be involved in their research, and for their support
280 in shaping this idea. Additionally, I would like to thank my Yale College Dean, Ployd, for his
281 encouragement and patience throughout this journey.

282 **References**

- 283 [1] Chain of Logic: Rule-Based Reasoning with Large Language Models. Accessed: 2024-09-23.
- 284 [2] DISC-LawLLM: Fine-Tuning Large Language Models for Intelligent Legal Services. Accessed: 2024-09-23.
- 285 [3] Enabling Discriminative Reasoning in LLMs for Legal Judgment Prediction. Accessed: 2024-09-23.
- 286 [4] LawLLM: Law Large Language Model for the US Legal System. Accessed: 2024-09-23.
- 287 [5] LegalBench: A Collaboratively Built Benchmark for Measuring Legal Reasoning in Large Language Models.
288 Accessed: 2023-10-16.
- 289 [6] LegalGPT: Legal Chain of Thought for the Legal Large Language Model Multi-agent Framework. Accessed:
290 2024-09-23.
- 291 [7] SaulLM-7B: A Pioneering Large Language Model for Law. Accessed: 2024-09-23.