

# LegalViz: Legal Text Visualization by Text To Diagram Generation

Eri Onami<sup>1,2</sup> Taiki Miyanishi<sup>3,6,2</sup> Koki Maeda<sup>4,7</sup> Shuhei Kurita<sup>5,7,2</sup>

<sup>1</sup>Nara Institute of Science and Technology <sup>2</sup>RIKEN AIP <sup>3</sup>The University of Tokyo

<sup>4</sup>Institute of Science Tokyo <sup>5</sup>National Institution of Informatics <sup>6</sup>ATR <sup>7</sup>NII LLMC

onami.eri.ob6@is.naist.jp, taiki.miyanishi@weblab.t.u-tokyo.ac.jp,  
koki.maeda@nlp.c.titech.ac.jp, skurita@nii.ac.jp

## Abstract

Legal documents including judgments and court orders require highly sophisticated legal knowledge for understanding. To disclose expert knowledge for non-experts, we explore the problem of visualizing legal texts with easy-to-understand diagrams and propose a novel dataset of LegalViz with 23 languages and 7,010 cases of legal document and visualization pairs, using the DOT graph description language of Graphviz. LegalViz provides a simple diagram from a complicated legal corpus identifying legal entities, transactions, legal sources, and statements at a glance, that are essential in each judgment. In addition, we provide new evaluation metrics for the legal diagram visualization by considering graph structures, textual similarities, and legal contents. We conducted empirical studies on few-shot and finetuning large language models for generating legal diagrams and evaluated them with these metrics, including legal content-based evaluation within 23 languages. Models trained with LegalViz outperform existing models including GPTs, confirming the effectiveness of our dataset.

## 1 Introduction

Driven by the rapid advancements in large language model (LLM) performance (Brown et al., 2020; OpenAI, 2023), adaptation to specialized domains in Natural Language Processing (NLP) receives increasing attention in many fields (Lu et al., 2022; Kung et al., 2023; Guha et al., 2023). Specifically, the application of LLMs to the legal field holds the potential to automate significant tasks and support roles traditionally occupied by lawyers and judges (Choi et al., 2021; Frankenreiter and Nyarko, 2022). The understanding of legal documents poses unique challenges for NLP applications, as legal reasoning requires not only interpreting the surface utterance but also implicit rules often omitted from the legal documents. It also requires following legal syllogisms, understanding the implications of

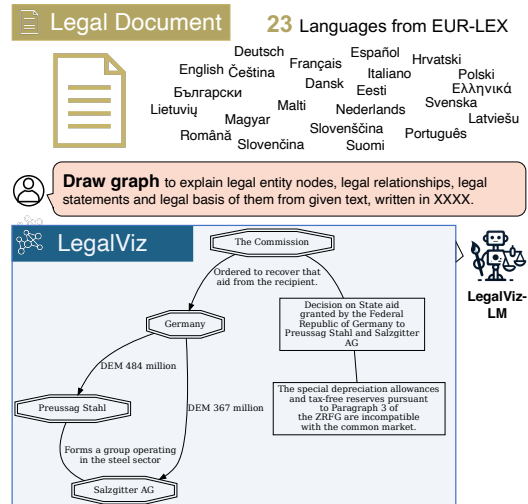


Figure 1: Model input and expected output of legal text visualization drawn by Graphviz.

related regulations, and applying them to specific case facts to deduce the final consequences.

At an early stage of legal NLP, there are several studies applying traditional NLP approaches for legal documents, such as named entity recognition (Angelidis et al., 2018; Luz de Araujo et al., 2018; Pais et al., 2021; de Gibert Bonet et al., 2022), summarization (Elaraby and Litman, 2022; Aumiller et al., 2022), text classification (Chalkidis et al., 2019) and text segmentation (Aumiller et al., 2021). In addition, some notable studies focus on capturing the structural legal meanings inherent in legal documents, such as learning judgment facts and results (Niklaus et al., 2021), assessing the fairness of law (Chalkidis et al., 2022b), and using the facts and attributes to predict charges (Hu et al., 2018). However, there are still numerous gaps between current legal domain status and whole automation of legal tasks by LLMs, especially judgment generation.

The main challenges of LLM for legal applications are as follows: (i) LLMs need to understand

which legal entities are involved, the relationships between them, and the relevant legal rules. They must also interpret the meaning of legal actions in judgments and court decisions. If models overlook legal entities, their rights, their obligations, or key facts in interpreting the law, they fail to fulfill legal requirements, or deduce inappropriate conclusions. (ii) LLMs must articulate why quoted rules are interpreted in their view, explaining the requirements and effects of the rules should become as they assert. This process should adhere to the procedure of legal syllogism, requiring the recognition of potential legal entities, relationships of them, and related rules. To address these challenges, extensive legal document resources are crucial for effectively tuning LLMs to perform well in legal domains. However, despite the wide accessibility of plain texts of laws and court judgments on the internet, there remains a significant lack of legal datasets with professional annotations that can enhance the capability for legal syllogism. Moreover, LLM technology should enhance legal adaptation capabilities, supporting not only professionals but also non-experts, as everyone has legal rights and should have the opportunity to benefit from the law.

To meet these demands, we explore the novel dataset LegalViz, an automatic visualization task, generating legal diagrams that describe the legal entities, their legal relationship, related rules, and summary of the key legal facts for legal interpretation from input legal plain texts. We introduced this visualization task with an existing diagram visualization tool of GraphViz because diagrams can succinctly elucidate complex legal relationships, allowing viewers to grasp the essentials at a glance without consulting the original article. In fact, the visualization of legal concepts is employed in various contexts, such as textbooks for judicial examinations, university classrooms, and TV news segments. This approach provides non-experts with easy-to-interpret visual and conceptual representations of legal materials, enhancing accessibility and understanding. By training with our dataset, models can accurately recognize legal rules concerned in the case, identify legal entities capable of exercising rights, and understand legal transactions, and statements from professional legal documents. LegalViz consists of 7,010 pair professional legal documents and diagrams of DOT language code of Graphviz, with 23 languages of EUR-LEX. Figure 1 from the LegalViz dataset illustrates a legal diagram that explains a case where the commission

required Germany to recover aid based on the common market principles, and Germany subsequently issued recovery requests. We assume LegalViz is the first work to utilize LLMs for the visualization of legal documents.<sup>1</sup>

## 2 Related Work

NLP applications in the legal domain are several core areas (Katz et al., 2023) such as information extraction, classification, summarization, judgment prediction, and resources and benchmarks.

**Judgment prediction.** In this task, models predict the outcomes of legal cases based on given facts. Previous studies provide judgment data from various courts of diverse countries, including decisions from the Supreme Court of the United States (Katz et al., 2017) and the European Court of Human Rights (Medvedeva et al., 2020; Kaur and Bozic, 2019). Additionally, judgment prediction research has covered Switzerland (Niklaus et al., 2021), Chinna (Ye et al., 2018), criminal law (Chen et al., 2019; Xiao et al., 2018), and asylum decisions (Chen and Egel, 2017; Dunn et al., 2017).

**Legal resources and benchmarks.** Datasets and benchmarks, covering a broad range of legal domains and languages, have been proposed. These include English Tax Law (Holzenberger et al., 2020), European Legislation and the European Court of Human Rights (Chalkidis et al., 2019), Corporate and Contract Law (Hendrycks et al., 2021; Tuggener et al., 2020), Supreme Court cases and US court cases (Zheng et al., 2021). The scope extends to German legal cases (Urchs. et al., 2021), a mixture of Korean legal text summarization, prediction and text classification (Hwang et al., 2024), and refugee cases (Barale et al., 2023). Additionally, multilingual and multi-legal domain datasets have been developed, such as a multilingual corpus of English, German, Italian, Polish (Drawzeski et al., 2021), and LEXGLUE (Chalkidis et al., 2022a) which covers six predictive tasks over five datasets made of English from the US, EU, and Council of Europe. Furthermore, Lexfiles (Chalkidis et al., 2023) offers a comprehensive dataset of comprised of US, UK, Canada, India, European Court of Human Rights, and Lextreme (Niklaus et al., 2023), which covers wide-range of tasks and countries among EU nations. However, none of these datasets are

<sup>1</sup>Our dataset is available at <https://github.com/mizuumi/LegalViz>.

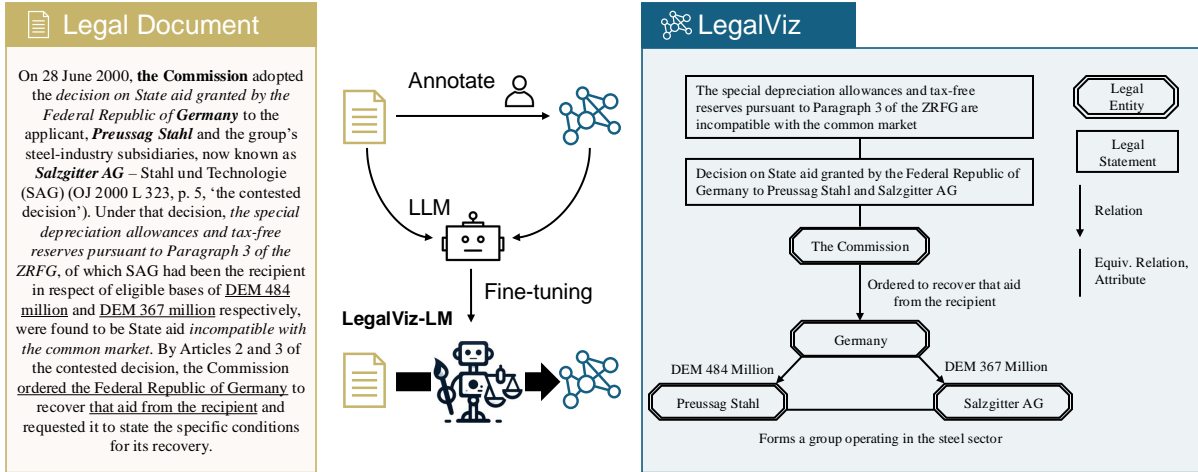


Figure 2: Legal text from EUR-LEX (left) to the resulting legal graph (right).

designed to support the visualization of legal documents for non-experts. In contrast, LegalViz offers legal specific annotations in 23 multilingual legal documents, specifically tailored for visualization. These annotations cover legal entities, their relationship, related rule, related facts of legal texts, thereby enhancing the clarity and interpretation of legal documents for judicial judgments.

**Text to graph generation.** Following the iconic successions of the GPT models, LLMs can generate not only contextual texts and program codes (Shi et al., 2022; Christopoulou et al., 2024) but also visualization codes (Bubeck et al., 2023), such as creation of scientific vector graphics with TiKZ (Belouadi et al., 2024) and diagram generation with refinements and diffusion process (Zala et al., 2023). Text-to-code generation studies are predominantly focused on mainstream programming languages like Python and shell scripts, and are typically examined with English text (Shi et al., 2022; Christopoulou et al., 2024). Both text-to-graph generation and graph-to-text generation studies are often conducted for clarifying paragraph structure and summarizing critical issues and relationship between words (Koncel-Kedziorski et al., 2019; Jin et al., 2020) of the input plain texts mainly in English. These text-to-graph approaches are suitable for free drawing based on text instructions, but they sacrifice the visualization of logical relationships within the visualized content. In comparison, our graph generation approach utilizes the DOT language of Graphviz, enabling models to focus specifically on visualizing the logical relationships within the content.

### 3 Building LegalViz Dataset

#### 3.1 Legal Visualization

The aim of legal visualization tasks is to generate an interpretable graph that clarifies the legal relationships embedded within the input legal texts. The constructed graph comprises legal entities and/or rules as nodes, connected by edges representing legal transactions and/or significant facts relevant for judicial determination. To effectively visualize these legal relationships, we utilize the DOT language of Graphviz, a widely adopted open-source tool for graph visualization. Figure 2 presents an overview of our proposed task, showcasing both the expected input and output.

**Legal entity.** Legal entities are applicants and respondents of judgment, courts, creditors, debtors, criminal suspects, or companies and employees. Legal entities are represented in **double octagons**. In contrast to grammatical general nouns, proper nouns, or objects, we concentrate on persons or organizations capable of exercising legal rights and engaging in transactions.

**Legal relationship & transactions.** Legal relationships encompass various form of relationships between legal entities, including the exercise of legal rights from one to another, legally significant transactions, the interrelations between legal statements made by entities and the underlying norms that support them, and relationships defined under law such as employment, contractual agreements, marriage, and family relationships. Legal transactions are specific types of relations among legal entities, such as purchases, notifications and any actions exercising rights. Both legal relationships

Split	# Instances	# Nodes	# Relations
Train	4,710	12,624	16,367
Validation	1,150	3,404	2,717
Test	1,150	3,128	3,589
Total	7,010	19,156	22,673

Table 1: Dataset splits.

Lang.	ISO	# Ins.	$L_{\text{word}}$	$L_{\text{char}}$	$L_{\text{code}}$
All	-	7,010	109.0	644.2	759.8
Bulgarian	BG	290	113.4	625.5	759.8
Spanish	ES	307	133.7	693.4	708.4
Czech	CS	307	102.8	582.9	832.5
Danish	DA	307	110.9	640.7	766.8
German	DE	312	108.9	683.0	759.1
Estonian	ET	307	83.9	588.8	809.4
Greek	EL	307	121.4	698.6	779.2
English	EN	312	122.6	629.2	623.0
French	FR	312	128.6	674.8	766.9
Croatian	HR	263	103.2	577.7	718.7
Italian	IT	312	123.3	705.1	788.7
Latvian	LV	307	94.4	598.8	725.7
Lithuanian	LT	307	94.6	609.4	749.8
Hungarian	HU	307	97.4	670.2	809.7
Maltese	MT	305	100.4	706.3	777.7
Dutch	NL	312	122.0	687.0	784.7
Polish	PL	307	106.7	655.0	759.2
Portuguese	PT	307	125.2	653.1	778.0
Romanian	RO	290	118.3	672.0	791.8
Slovak	SK	308	101.0	585.7	727.9
Slovenian	SL	308	106.6	580.0	730.5
Finnish	FI	308	78.5	649.6	808.7
Swedish	SV	308	108.9	639.1	748.2

Table 2: Dataset statistics.  $L_{\text{word}}$  and  $L_{\text{char}}$  are length of legal text.  $L_{\text{code}}$  is character length of Graphviz code.

and transactions are represented in the directed or undirected **edges** with various styles between legal entities. Some edges have textual relation labels.

**Legal source.** Legal sources are rules applied or referred by court and support legal statements. Here we concentrate on legal sources explicitly written in the legal document. They include constitutions, statutes, ordinances, and case laws. These extracted rules are represented in **trapeziums**. Each trapezium of the legal source is connected to a node interpreting rules supported by the legal source via undirected edges.

**Legal statement.** Legal statements are detailed explanations of transactions and factual descriptions of the case notable for the final judgment to summarize. Adding these summaries to diagrams help non-experts grasp the facts important for final judgments at a glance. Legal statements are represented in **squares**, connected to a node by an edge.

## 3.2 Dataset Creation

**Collection of legal document.** To construct the legal graph dataset, we collected legal documents in the following manner: (i) We sourced legal documents from the EUR-LEX website<sup>2</sup>, which provides public access to judgments, orders, and rules of EU countries in official EU languages. We specifically selected judgments from the years 2006 to 2019, available in translations across 23 languages, to capture the latest legal trends. (ii) We extracted the factual sections of the judgments that contain legal facts to be expressed in the graph. (iii) Finally, we gathered the corresponding sections of legal documents in 23 languages to ensure consistency across translations.

**Graphviz annotation.** We have manually annotated Graphviz code visualization from the legal documents by an annotator with expertise in the legal domain. The process involved several steps: (i) We broke down long judgment cases into short paragraphs so that DOT language can draw diagrams in units easily understandable at a glance. (ii) We extracted the legal entities and rules as nodes of the diagram, legal transactions as edge relations within the diagram, and the summary of statements and explanations as squared nodes. (iii) We created a Graphviz diagram to represent the extracted relations, using variations in node shape and relations, as defined by the rules in Appendix F.

**Translation of Graphviz annotation.** To cover the European Union’s official languages present at the time the judgment was written, we translated our English annotation to other languages as follows: (i) We first used GPT-4 to extract the legal words and sentences from the provided English judgements, aiming to save as many terms as possible from the EU’s officially translated variations of judgments. (ii) If GPT-4 fails the extraction task, we then apply GPT-4 translations from English to other languages. (iii) We manually checked the previously translated sentences and retranslated them using DeepL and the Azure GPT API if any translation errors were detected. The prompts used in the translation process are described in Appendix I.

**Dataset statistics.** Table 1 shows the statistics of the LegalViz dataset. We build a total of 7,010 pairs of legal texts and graphs, encompassing 23 language variations and more than 300 unique legal texts. The constructed legal graph consists of 19,159 nodes and 22,673 relations. We also summa-

<sup>2</sup><https://eur-lex.europa.eu>

rize the average word length, number of characters in legal sentences, and character length of Graphviz code for each language in Table 2.

## 4 Evaluation

We compare the reference and hypothesis graphs to assess the quality of the generated. One straightforward way to achieve this is to directly compare two images visualized by GraphViz. However, this approach clearly ignores the textual structure of the legal documents. One other approach is directly comparing the DOT language codes in textual manner, ignoring numerous minor differences of the visualization codes that can result in different graphs. Therefore we propose an approach to first compose graphs, align the graph components, and then compare each component of graphs using textual metrics as described in this section.

### 4.1 Similarity of two graphs with texts

Formally, let  $\mathcal{G}_r$  and  $\mathcal{G}_h$  be the reference and hypothesis graphs. Each graph is composed of a set of edges  $E$  and nodes  $V$ . When an edge  $e \in E$  connects a starting node  $v_s$  and an end node  $v_e$ , it is represented by a tuple  $e = [v_s, v_e, l]$ , where  $l$  is a label of this edge. Nodes always include non-empty texts, while edge-label texts can be blank for edges without labels.

**DOT code validation.** First, we examine whether the generated Graphviz code forms a valid graph  $\mathcal{G}_h$  in terms of the DOT language. This is done by simply processing with the pydot library<sup>3</sup>.

**Nodes alignment by bipartite matching.** Second, we extract nodes  $V_h$  from  $\mathcal{G}_h$  and align them with nodes from the reference graph:  $V_r$  from  $\mathcal{G}_r$  using the similarity of the texts in nodes. For this node alignment, we apply the bipartite matching problem for the two sets of nodes  $V_h$  and  $V_r$ , using the similarity function  $\text{sim}(v_r, v_h)$  between two texts in the reference  $v_r \in V_r$  and hypothesis nodes  $v_h \in V_h$  with a textual similarity metric. Here we use BERTScore (Zhang et al., 2020) for the similarity metric because of its robustness and high human correlation. Given the textual similarity scores between all reference and hypothesis nodes, we apply a bipartite matching solver in NetworkX<sup>4</sup> for nodes and obtain the nodes alignment between the reference and hypothesis graphs that are used for later evaluation.

<sup>3</sup><https://github.com/pydot/pydot>

<sup>4</sup><https://networkx.org/>

**Graph, node & edge evaluation.** After we determined the node alignment, we performed the comparison of the two graphs based on the nodes, edges and their labels. We introduce the following three metrics with different depth: Graph, Graph&Node and Graph&Node&Edge.

Graph is an F1 metrics of the matched edges after the node alignment by bipartite matching. This metric is for the similarity measurement of the entire graph structure, ignoring the textual differences of nodes and edges after the alignment. Let the node set  $V_r$  and edge set  $E_r$  composes reference graph  $G_r$ , and the node set  $V_h$  and edge set  $E_h$  composes hypothesis (generated) graph  $G_h$ . Using a node alignment function  $a(\cdot) : V_h \rightarrow \{V_r, \phi\}$  from the generated graph nodes  $V_h$  to the reference graph nodes  $V_r$  and Kronecker delta  $\delta_{\mu\nu} = 1$  iif  $\mu = \nu$  otherwise  $\delta_{\mu\nu} = 0$ , which represents the agreement of the nodes here, we compute the agreement score of the nodes as

$$f_{\text{Graph}}(e_h, e_r) = \delta_{a(v_{s,h})v_{s,r}} \delta_{a(v_{e,h})v_{e,r}}$$

$v_{s,h}$  and  $v_{s,r}$  are the start nodes of each edge from the hypothesis (generated) and reference graphs, while  $v_{e,h}$  and  $v_{e,r}$  are the end nodes of it, respectively. The generated nodes can be aligned to  $\phi$  when they are not aligned to any reference nodes:  $v_h \xrightarrow{a(\cdot)} \phi$ . Here  $\phi$  is a null node, and we assume for any nodes  $\nu$ ,  $\delta_{\phi\nu} = 0$ . From this binary function  $f_{\text{Graph}}(e_h, e_r)$ , we can compute TP, FP and FN by

$$\begin{aligned} \text{TP} &= \sum_{e_h \in E_h, e_r \in E_r} f_{\text{Graph}}(e_h, e_r) \\ \text{FP} &= |E_h| - \text{TP}, \quad \text{FN} = |E_r| - \text{TP} \end{aligned}$$

and obtained F1 value from these for Graph.

Graph&Node is an F1-based metric where BERTScore penalize the dissimilar texts of the two aligned nodes pairs  $\{v_{s,h}, v_{s,r}\}$  and  $\{v_{e,h}, v_{e,r}\}$  for each edge. TP is calculated as

$$\text{TP} = \sum_{e_h \in E_h, e_r \in E_r} f_{\text{Graph}}(e_h, e_r) \cdot \text{sim}(v_{s,r}, v_{s,h}) \text{sim}(v_{e,r}, v_{e,h})$$

while FP and FN are calculated from TP. Because of the products of the start and end node similarities, the Graph&Node metric is sensitive to the difference of node texts compared with the Graph metric.

Model	Graph-based			Valid Graph		Legal Content			
	G	G-N	G-N-E	Top1	Top10	Entity	R & T	Source	Statement
<i>Few-shot results</i>									
CodeLlama 7B	12.88	9.10	3.70	16.78	85.22	48.94	5.17	10.11	1.24
CodeLlama 7B it.	15.67	11.78	6.07	37.65	89.39	55.10	8.01	11.00	1.29
CodeLlama 13B	15.33	10.90	5.23	17.30	85.04	51.46	7.34	11.89	2.38
CodeLlama 13B it.	16.37	12.35	6.47	33.39	88.70	55.00	8.54	10.76	2.21
Llama3.1 8B	26.10	20.32	11.18	30.00	83.22	64.06	14.21	16.85	2.84
Llama3.1 8B it.	24.47	17.91	10.32	24.00	84.00	62.96	13.95	16.22	2.21
Llama3.2 3B	22.20	17.06	8.65	27.13	80.52	57.35	11.18	12.24	2.28
Llama3.2 3B it.	25.64	19.80	11.38	56.26	92.09	54.11	14.51	10.93	2.78
Gemma2 9B	15.35	11.28	5.28	35.30	93.30	54.88	7.03	9.18	2.56
Gemma2 9B it.	27.22	22.44	12.64	70.70	94.17	73.27	15.16	17.21	1.82
GPT-3.5-Turbo	26.66	22.28	13.51	94.26	<b>100.0</b>	73.02	16.18	13.81	<b>3.88</b>
GPT-4	<b>33.46</b>	<b>28.70</b>	<b>19.96</b>	<b>99.13</b>	<b>100.0</b>	<b>75.31</b>	<b>23.24</b>	<b>21.52</b>	3.30
GPT-4o	23.58	20.10	13.42	95.22	<b>100.0</b>	75.15	15.82	19.93	2.97
<i>Finetuning results</i>									
CodeLlama 7B	30.56	23.04	16.34	94.43	99.57	76.73	21.54	39.81	8.59
CodeLlama 7B it.	33.47	25.85	18.68	96.61	99.65	76.90	24.00	34.61	9.03
CodeLlama 13B	34.44	25.94	17.70	97.13	99.83	76.73	23.23	42.23	7.43
CodeLlama 13B it.	35.61	27.75	19.65	96.17	99.65	77.68	24.87	46.32	9.85
Llama3.1 8B	30.09	19.86	13.25	94.70	<b>100.0</b>	68.22	19.75	29.01	9.39
Llama3.1 8B it.	29.59	20.32	13.42	87.91	99.83	70.57	18.98	31.51	9.28
Llama3.2 3B	33.38	24.29	17.56	92.78	99.83	73.29	23.83	47.47	9.89
Llama3.2 3B it.	30.37	21.51	14.70	87.22	99.83	71.93	20.38	43.24	10.08
Gemma2 9B	<b>43.38</b>	<b>36.47</b>	<b>27.52</b>	<b>98.00</b>	<b>100.0</b>	<b>81.85</b>	<b>32.53</b>	<b>50.97</b>	<b>12.75</b>
Gemma2 9B it.	42.30	34.26	25.95	96.17	<b>100.0</b>	81.02	31.80	42.05	11.92

Table 3: Overall results of the legal text visualization in the LegalViz test set. **G**, **G-N** and **G-N-E** denote Graph, Graph&Node and Graph&Node&Edge respectively. Valid Graph is success rate of creating valid DOT language codes in top-1 and top-10 generated results. “it.” means instruct tuning models. The highest scores of each column are in bold.

Similarly, Graph&Node&Edge is an F1-based metric that considers node and edge text similarity in terms of the BERTScore as following

$$TP = \sum_{e_h \in E_h, e_r \in E_r} f_{\text{Graph}}(e_h, e_r) \cdot \text{sim}(v_{s,r}, v_{s,h}) \text{sim}(v_{e,r}, v_{e,h}) \text{sim}(l_r, l_h)$$

by penalizing dissimilar texts of the edge texts.

## 4.2 Evaluation of Legal Content

We also introduce the evaluation of legal contents in our visualizations. As described in Sec. 3.1, the legal contents in LegalViz are associated with specific diagrams in GraphViz. For diagrams of legal entities (**double octagon**), Legal source (**trapeziums**), Legal statement (**squares**), we extract these nodes from the reference ( $v_r$ ) and hypothesis ( $v_h$ ) graphs and check whether they are properly aligned in the alignment in the previous section. Then we measure the similarity of the node texts with BERTScore for successfully aligned nodes and compute micro averaged F1 score as

following TP, FP, FN:

$$TP = \sum_{v_h \in V_h, v_r \in V_r} \delta_{v_r, v_h} \text{sim}(v_r, v_h)$$

$$FP = |\{v_h\}| - TP$$

$$FN = |\{v_r\}| - TP$$

where  $\delta_{v_r, v_h} = 1$  iff  $v_r$  and  $v_h$  is aligned, and otherwise 0. For legal relations & transactions (**edges**), we extract the aligned edge label texts and compute F1 score from the similarity of labels.

## 5 Experiments

We evaluate the ability to visualize graphs from legal sentences with LegalViz. Overall, our finetuned models overperformed fewshot GPT models.

### 5.1 Experimental settings

We conduct the experiments of legal visualization in the manner of the DOT language code generation with the publicly available Llama and Gemma family models and OpenAI GPT APIs via Microsoft Azure. We use the GPT-3.5-Turbo (1106), GPT-4 (0613) and GPT-4o (2024-05-13) models. For

Model	BG	ES	CS	DA	DE	ET	EL	EN	FR	HR	IT	LV	LT	HU	MT	NL	PL	PT	RO	SK	SL	FI	SV
<i>Entity</i>																							
Gemma 2 9B fs.	59.22	59.20	52.53	54.47	56.24	53.51	53.22	59.95	53.77	55.32	55.42	51.33	55.04	42.26	47.55	60.06	51.17	61.68	59.21	52.85	50.72	55.69	57.21
Gemma 2 9B ft.	<b>80.62</b>	<b>84.33</b>	<b>80.06</b>	<b>82.53</b>	<b>83.31</b>	<b>78.57</b>	<b>80.44</b>	<b>86.98</b>	<b>82.90</b>	<b>81.14</b>	<b>80.94</b>	<b>81.11</b>	<b>77.16</b>	<b>81.66</b>	<b>82.59</b>	<b>82.70</b>	<b>82.58</b>	<b>85.64</b>	<b>83.97</b>	<b>81.27</b>	<b>79.16</b>	<b>79.23</b>	<b>83.22</b>
Avg. fs. models	60.78	63.32	57.12	59.67	61.30	57.67	53.59	65.77	60.61	60.48	59.06	56.85	55.73	56.64	61.37	60.18	61.36	60.89	59.30	57.20	59.68	61.04	
Avg. ft. models	73.19	77.05	73.59	74.14	74.28	72.62	72.31	78.78	75.5	74.87	75.82	71.87	72.08	72.16	73.36	75.2	74.25	76.72	76.22	74.08	72.77	72.47	75.17
<i>Relations&amp;Transactions</i>																							
Gemma 2 9B fs.	10.64	7.24	8.97	10.44	2.28	7.31	7.31	13.39	5.23	7.11	4.42	3.54	6.79	2.81	3.71	9.27	8.01	11.18	6.90	6.81	5.94	2.40	6.73
Gemma 2 9B ft.	<b>32.00</b>	<b>33.88</b>	<b>31.19</b>	<b>37.13</b>	<b>32.34</b>	<b>26.13</b>	<b>30.42</b>	<b>36.53</b>	<b>23.06</b>	<b>36.43</b>	<b>27.14</b>	<b>30.74</b>	<b>28.06</b>	<b>38.13</b>	<b>31.88</b>	<b>29.68</b>	<b>40.27</b>	<b>34.39</b>	<b>37.17</b>	<b>35.15</b>	<b>29.50</b>	<b>34.40</b>	<b>33.02</b>
Avg. fs. models	11.22	11.80	10.28	13.90	11.26	8.74	12.56	10.76	11.67	10.47	10.12	12.78	12.19	11.24	9.22	13.66	12.18	11.88	10.62	10.42	11.27	10.02	11.68
Avg. ft. models	23.66	23.18	22.96	25.47	22.98	23.0	23.01	27.72	22.09	25.56	20.87	24.4	25.79	21.45	23.09	20.95	24.13	24.08	23.2	23.36	23.3	21.94	24.86
<i>Source</i>																							
Gemma 2 9B fs.	0.00	0.00	15.15	10.37	13.85	17.95	15.91	17.96	0.00	9.92	0.00	4.97	0.00	6.89	0.00	6.20	19.31	12.47	6.05	14.77	0.00	7.51	23.25
Gemma 2 9B ft.	<b>58.98</b>	<b>49.07</b>	<b>41.06</b>	<b>52.69</b>	<b>56.19</b>	<b>54.21</b>	<b>47.31</b>	<b>61.32</b>	<b>61.87</b>	<b>51.14</b>	<b>49.90</b>	<b>64.56</b>	<b>38.86</b>	<b>42.34</b>	<b>48.10</b>	<b>44.85</b>	<b>53.68</b>	<b>42.08</b>	<b>54.92</b>	<b>48.00</b>	<b>47.03</b>	<b>54.46</b>	<b>54.68</b>
Avg. fs. models	12.19	11.78	14.39	13.89	13.58	12.11	12.38	14.07	15.02	15.15	13.04	14.51	14.84	13.35	11.65	16.77	14.65	13.66	12.89	15.30	9.52	11.08	15.20
Avg. ft. models	36.8	37.3	36.52	37.38	45.55	38.04	33.58	36.04	36.92	40.67	36.5	40.69	<b>45.54</b>	39.53	37.58	35.53	42.92	34.26	35.04	44.14	36.46	34.54	38.53
<i>Statement</i>																							
Gemma 2 9B fs.	0.00	3.25	3.27	3.74	4.42	2.30	1.53	1.54	1.59	3.85	0.00	0.00	1.95	2.41	4.40	3.94	3.03	4.86	1.40	1.32	4.43	4.54	1.45
Gemma 2 9B ft.	<b>16.78</b>	<b>17.10</b>	<b>9.17</b>	<b>12.76</b>	6.78	<b>9.86</b>	<b>20.95</b>	5.90	<b>11.09</b>	<b>11.13</b>	<b>16.32</b>	<b>11.21</b>	<b>15.90</b>	<b>15.22</b>	<b>11.10</b>	<b>8.92</b>	<b>12.03</b>	<b>11.90</b>	<b>16.09</b>	<b>17.65</b>	<b>13.17</b>	<b>7.46</b>	<b>14.89</b>
Avg. fs. models	2.10	2.66	1.57	2.55	2.50	1.62	2.15	3.75	2.98	1.99	1.38	2.50	3.04	1.75	1.71	2.31	3.63	2.24	2.49	1.82	2.34	1.25	1.86
Avg. ft. models	10.52	12.05	6.94	9.62	<b>8.64</b>	7.33	13.1	<b>11.23</b>	8.9	9.07	10.92	9.05	10.68	8.47	10.1	8.32	10.77	9.28	8.89	9.22	8.67	6.45	9.04

Table 4: Scores of Entity, Relations&Transactions, Source, and Statement by 23 languages. “fs.” means few-shot and “ft.” means finetuned with the LegalViz dataset. Avg. fs. models exclude GPTs for comparison.

Llama family models, we experimented with the models specialized for code generation of CodeLlama (Rozière et al., 2023) and the recently released Llama 3.1 & 3.2 models (Dubey et al., 2024) and Gemma 2-9B (Riviere et al., 2024) models. Experimental settings are two holds: few-shot generation and finetuning of the publicly available models. In few-shot experiments, we notice not only the GPT models but only publicly available models are capable of producing valid DOT language codes without finetuning to some extent. We follow the supervised finetuning of Hugging Face with the detailed finetuning parameters in Appendix B. In evaluation, we generate ten different Graphviz codes with each model and examine their quality in evaluation methods of graph and legal contents described in Sec. 4.

## 5.2 Result

### Graph-based Evaluation.

The Graph-based Evaluation section of Table 3 presents the experimental results of each model evaluated by Graph, Graph&Node, and Graph&Node&Edge metrics explained in Section 4. Among three evaluation metrics, Graph&Node&Edge is the most difficult because all three graph elements must be correct as shown in the evaluation. Most importantly, our finetuned models outperformed few-shot counterparts and even GPT models, which are assumed to be larger than the CodeLlama-13B models, suggesting the effectiveness of our dataset for finetuning. Also, finetuned Gemma-2-9B took the highest scores on Graph and Graph&Node, and Graph&Node&Edge. Surprisingly, Gemma-2-9B performed worse than Gemma-2-9B-it before

finetuning, suggesting the effectiveness of finetuning with LegalViz.

**Valid DOT code ratio.** The Valid Graph section of Table 3 presents the success rate of forming a valid DOT language code without code syntax errors. In the first generation trial, GPT-4 is the most accurate to generate valid DOT language codes among all models in both few-shot and finetuned settings and the second best is Gemma2-9B finetuned model. When we let models generate ten variations, several finetuned models (Llama-3.1-8B, Gemma2-9B, Gemma2-9B-it) are able to generate valid DOT code in 100.0 percents of the test set. Comparing the publicly available few-shot and finetuned models, finetuning with our dataset strongly improves the valid graph creation of all models.

**Legal Content Evaluation.** The Legal Content section of Table 3 presents the legal aspect-wise evaluation as described in Sec. 4.2. By nature of the legal entities, Entity can be extracted from input sentences in many conditions, achieving high scores in the table. However, the other three aspects aren’t easily extracted from the input legal text. Statement includes the text generation for legal facts and tends to be lower scores than others. This is because legal statements appear in texts without some remarkable keywords, compared to legal Source, which is often mentioned in texts with terminology such as “Law” and “Act” and legal Relations & Transactions, which is found in texts with terminologies such as “contract,” “issue” with some warrants and orders, “notification” with notable as a legal act. Statement acts for summarizing notable facts related to rule and its interpretation in question, to describe the detail of other nodes especially legal relations and trans-

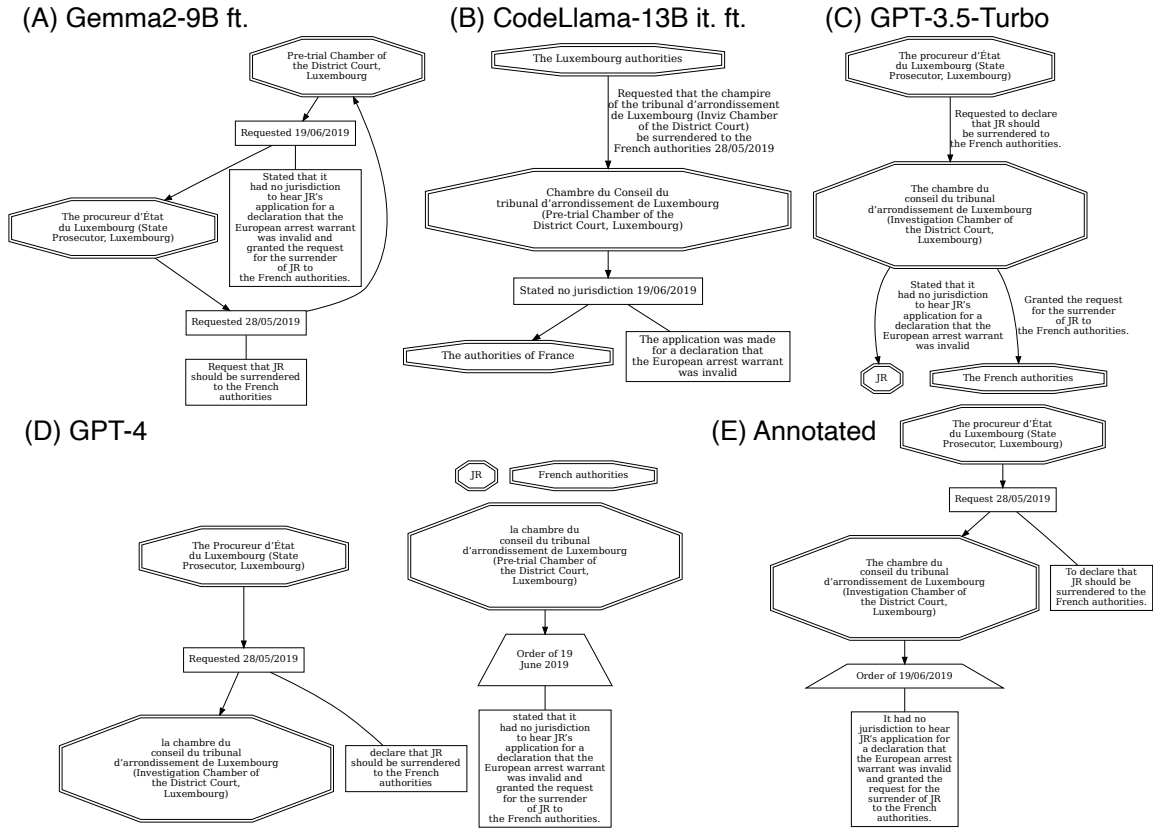


Figure 3: Qualitative analysis of diagrams by Graphviz code. Figures are generated by the finetuned models of Gemma2-9B, CodeLlama-13B-Instruct, few shot models of GPT-3.5-Turbo, GPT-4 and an annotated diagram.

actions, and to explain the facts applicable to legal requirements. Finetuned Gemma2-9B achieved highest score in all four aspects of the legal content evaluation. The scores in *Statement* are improved by approximately three times compared to the few-shot scores across all models, suggesting the effectiveness of finetuning with our dataset.

**Scores by languages.** Table 4 presents the results in legal contents by all 23 languages in EUR-LEX. We present the best performing model of Gemma2-9B in finetuned and fewshot conditions. We also present the averaged results of 10 models in Table 3 without GPTs to highlight the performance difference before and after finetuning with LegalViz across languages, while minimizing the influence of individual model characteristics. Among these languages, models perform relatively weakly in languages that have relatively fewer resources (Chalkidis et al., 2021a), such as Maltese, Latvian, Lithuanian, Hungarian. For languages that have relatively more resources such as English and French, models tend to have higher scores than others. From a linguistic point of view, Hungarian and Finnish, belonging to the same Uralic language group, have low scores in

each model. This may reflect their linguistic difference from other languages. Similarly, for the Romance language group, e.g., Romanian, French, Spanish, Italian, and Portuguese, models have moderate performances, seemingly better than those of the Uralic language group and languages that also have fewer resources than those of English and French. Among four legal aspects, the source and statement parts include the summarization task of the legal document for visualization. They are considerably difficult parts in the graphs and the performance in some languages becomes 0 without finetuning. It is also notable that finetuning contributes the performance in all aspects in all of these languages.

### 5.3 Qualitative Analysis

We conduct a qualitative analysis using the best performing model of the finetuned Gemma2-9B and CodeLlama-13B Instruct models, along with the few-shot GPT-3.5-Turbo and GPT-4. Figure 3 presents the graphs generated from English legal documents along with the annotated graph. This legal document used for the model input is on the Appendix A. This is a part of a criminal procedural



case in which the prosecutor requested the court to declare securing custody where the prosecutor and the court are legal entities. Square nodes in annotated data are describing intention of request and consequence of the request written in order.

For fewshot models of GPT-3.5-Turbo and GPT-4, GPT-3.5-Turbo wrongly recognize that all nodes are legal entities and illustrate them in double octagon. Its description of the legal relationship between the court and the authority is also incorrect because the court didn't take the legal action directly in this article. The GPT-4 model assumes that "JR" and "French authorities" are legal entities but fail to illustrate relationships between those and other entities as these entities lack connections to other graph parts. Also, GPT-4 extracted two different court names from given text as different entities. However they are indeed different divisions of the same court and fails to summarize the relationship between them. Gemma2-9B ft. successfully extracts the detailed description of court's request, while it extracted incorrect entity as double octagon shape and it couldn't extract the court order and its description. CodeLlama-13B-Instruct ft. model successfully forms the graph structure but the square node mentioned European arrest warrant in the right bottom is inconsistent with input text.

## 6 Conclusion

We have proposed LegalViz, the first manually annotated dataset to visualize legal text with DOT language Graphviz and introduced a novel evaluation method taking into account both diagram visualization quality and sentences of not only graph nodes and relations but also legal contents. We empirically confirmed the effectiveness of our dataset with wide-range of experiments including comparisons of few-shot and finetuning models and demonstrated trained models outperform the closed models of GPTs in all evaluation metrics.

## Limitation

LegalViz contains the same number of instances in 23 languages of EUR-LEX. However, this doesn't mean that the models with finetuned or few-shot have the same ability to treat all 23 languages equally. Especially models face difficulties in fewer language resources as we experimented. We cannot offer any warranty for using our dataset and models for real usages such as legal advice. We also consider that our dataset should be used with

appropriate supervision by experts. This can be a *potential risk* when our dataset is misused. We assume that results of automatic visualizations by models are still different from the annotated visualizations in most cases, suggesting the current limitation of the generation.

## Ethic Statements

The annotation material of this dataset is publicly available EU legal materials including judgments and orders, which do not include personal or sensitive information, with the exception of trivial information presented by consent, e.g., the names of the active presidents of the European Parliament, European Council, or other official administration bodies. The copyright for the editorial content of this website, the summaries of EU legislation, and the consolidated texts, which are owned by the EU, is licensed under the Creative Commons Attribution 4.0 International license.<sup>5</sup>

## References

- Iosif Angelidis, Ilias Chalkidis, and Manolis Koubarakis. 2018. [Named entity recognition, linking and generation for greek legislation](#). In *International Conference on Legal Knowledge and Information Systems*.
- Dennis Aumiller, Satya Almasian, Sebastian Lackner, and Michael Gertz. 2021. [Structural text segmentation of legal documents](#). In *Proceedings of the Eighteenth International Conference on Artificial Intelligence and Law, ICAIL '21*, page 2–11, New York, NY, USA. Association for Computing Machinery.
- Dennis Aumiller, Ashish Chouhan, and Michael Gertz. 2022. [EUR-lex-sum: A multi- and cross-lingual dataset for long-form summarization in the legal domain](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 7626–7639, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Claire Barale, Mark Klaisoongnoen, Pasquale Minervini, Michael Rovatsos, and Nehal Bhuta. 2023. [AsyLex: A dataset for legal language processing of refugee claims](#). In *Proceedings of the Natural Legal Language Processing Workshop 2023*, pages 244–257, Singapore. Association for Computational Linguistics.
- Jonas Belouadi, Anne Lauscher, and Steffen Eger. 2024. [Automatikz: Text-guided synthesis of scientific vector graphics with tikz](#). In *International Conference on Learning Representations (ICLR)*.

<sup>5</sup><https://eur-lex.europa.eu/content/legal-notice/legal-notice.html>

- Tom Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared D Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel Ziegler, Jeffrey Wu, Clemens Winter, Chris Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. [Language models are few-shot learners](#). In *Advances in Neural Information Processing Systems*, volume 33, pages 1877–1901. Curran Associates, Inc.
- Sébastien Bubeck, Varun Chandrasekaran, Ronen Eldan, John A. Gehrke, Eric Horvitz, Ece Kamar, Peter Lee, Yin Tat Lee, Yuan-Fang Li, Scott M. Lundberg, Harsha Nori, Hamid Palangi, Marco Tulio Ribeiro, and Yi Zhang. 2023. [Sparks of artificial general intelligence: Early experiments with gpt-4](#). *ArXiv*, abs/2303.12712.
- Ilias Chalkidis, Emmanouil Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2019. [Extreme multi-label legal text classification: A case study in EU legislation](#). In *Proceedings of the Natural Legal Language Processing Workshop 2019*, pages 78–87, Minneapolis, Minnesota. Association for Computational Linguistics.
- Ilias Chalkidis, Manos Fergadiotis, and Ion Androutsopoulos. 2021a. [MultiEURLEX - a multi-lingual and multi-label legal document classification dataset for zero-shot cross-lingual transfer](#). In *Proceedings of the 2021 Conference on Empirical Methods in Natural Language Processing*, pages 6974–6996, Online and Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Ilias Chalkidis, Manos Fergadiotis, Dimitrios Tsarapatanis, Nikolaos Aletras, Ion Androutsopoulos, and Prodromos Malakasiotis. 2021b. [Paragraph-level rationale extraction through regularization: A case study on European court of human rights cases](#). In *Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pages 226–241, Online. Association for Computational Linguistics.
- Ilias Chalkidis, Nicolas Garneau, Catalina Goanta, Daniel Katz, and Anders Søgaard. 2023. [LeXFiles and LegalLAMA: Facilitating English multinational legal language model development](#). In *Proceedings of the 61st Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 15513–15535, Toronto, Canada. Association for Computational Linguistics.
- Ilias Chalkidis, Abhik Jana, Dirk Hartung, Michael Bommarito, Ion Androutsopoulos, Daniel Katz, and Nikolaos Aletras. 2022a. [LexGLUE: A benchmark dataset for legal language understanding in English](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4310–4330, Dublin, Ireland. Association for Computational Linguistics.
- Ilias Chalkidis, Tommaso Pasini, Sheng Zhang, Letizia Tomada, Sebastian Schwemer, and Anders Søgaard. 2022b. [FairLex: A multilingual benchmark for evaluating fairness in legal text processing](#). In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 4389–4406, Dublin, Ireland. Association for Computational Linguistics.
- Daniel L. Chen and Jess Eigel. 2017. [Can machine learning help predict the outcome of asylum adjudications?](#) In *Proceedings of the 16th Edition of the International Conference on Artificial Intelligence and Law, ICAIL '17*, page 237–240, New York, NY, USA. Association for Computing Machinery.
- Huajie Chen, Deng Cai, Wei Dai, Zehui Dai, and Yadong Ding. 2019. [Charge-based prison term prediction with deep gating network](#). In *Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing (EMNLP-IJCNLP)*, pages 6362–6367, Hong Kong, China. Association for Computational Linguistics.
- Jonathan H Choi, Kristin E Hickman, Amy B Monahan, and Daniel Schwarcz. 2021. [Chatgpt goes to law school](#). *J. Legal Educ.*, 71:387.
- Fenia Christopoulou, Guchun Zhang, and Gerasimos Lampouras. 2024. [Text-to-code generation with modality-relative pre-training](#). In *Proceedings of the 18th Conference of the European Chapter of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1194–1208, St. Julian’s, Malta. Association for Computational Linguistics.
- Ona de Gibert Bonet, Aitor García Pablos, Montse Cuadros, and Maite Melero. 2022. [Spanish datasets for sensitive entity detection in the legal domain](#). In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pages 3751–3760, Marseille, France. European Language Resources Association.
- Kasper Drawzeski, Andrea Galassi, Agnieszka Jablonowska, Francesca Lagioia, Marco Lippi, Hans Wolfgang Micklitz, Giovanni Sartor, Giacomo Tagiuri, and Paolo Torroni. 2021. [A corpus for multilingual analysis of online terms of service](#). In *Proceedings of the Natural Legal Language Processing Workshop 2021*, pages 1–8, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, et al. 2024. [The llama 3 herd of models](#). *Preprint*, arXiv:2407.21783.
- Matt Dunn, Levent Sagun, Hale Şirin, and Daniel Chen. 2017. [Early predictability of asylum court decisions](#). In *Proceedings of the 16th Edition of the International Conference on Artificial Intelligence and Law*,

- ICAIL '17, page 233–236, New York, NY, USA. Association for Computing Machinery.
- Mohamed Elaraby and Diane Litman. 2022. [ArgLegal-Summ: Improving abstractive summarization of legal documents with argument mining](#). In *Proceedings of the 29th International Conference on Computational Linguistics*, pages 6187–6194, Gyeongju, Republic of Korea. International Committee on Computational Linguistics.
- Jens Frankenreiter and Julian Nyarko. 2022. Natural language processing in legal tech. *Legal Tech and the Future of Civil Justice*.
- Neel Guha, Julian Nyarko, Daniel E. Ho, Christopher Ré, Adam Chilton, Aditya Narayana, Alex Chohlas-Wood, Austin Peters, Brandon Waldon, Daniel N. Rockmore, Diego Zambrano, Dmitry Talisman, Enam Hoque, Faiz Surani, Frank Fagan, Galit Sarfaty, Gregory M. Dickinson, Haggai Porat, Jason Hegland, Jessica Wu, Joe Nudell, Joel Niklaus, John Nay, Jonathan H. Choi, Kevin Tobia, Margaret Hagan, Megan Ma, Michael Livermore, Nikon Rasumov-Rahe, Nils Holzenberger, Noam Kolt, Peter Henderson, Sean Rehaag, Sharad Goel, Shang Gao, Spencer Williams, Sunny Gandhi, Tom Zur, Varun Iyer, and Zehua Li. 2023. Legalbench: A collaboratively built benchmark for measuring legal reasoning in large language models. In *Proceedings of the 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Dan Hendrycks, Collin Burns, Anya Chen, and Spencer Ball. 2021. Cuad: An expert-annotated nlp dataset for legal contract review. *NeurIPS*.
- Nils Holzenberger, Andrew Blair-Stanek, and Benjamin Van Durme. 2020. [A dataset for statutory reasoning in tax law entailment and question answering](#). In *NLLP@KDD*.
- Zikun Hu, Xiang Li, Cunchao Tu, Zhiyuan Liu, and Maosong Sun. 2018. [Few-shot charge prediction with discriminative legal attributes](#). In *Proceedings of the 27th International Conference on Computational Linguistics*, pages 487–498, Santa Fe, New Mexico, USA. Association for Computational Linguistics.
- Wonseok Hwang, Dongjun Lee, Kyoungyeon Cho, Hanuhl Lee, and Minjoon Seo. 2024. A multi-task benchmark for korean legal language understanding and judgement prediction. In *Proceedings of the 36th International Conference on Neural Information Processing Systems, NIPS '22*, Red Hook, NY, USA. Curran Associates Inc.
- Michael J. Bommarito II, Daniel Martin Katz, and Eric M. Detterman. 2021. *Chapter 11: LexNLP: Natural language processing and information extraction for legal and regulatory texts*. Edward Elgar Publishing, Cheltenham, UK.
- Zhijing Jin, Qipeng Guo, Xipeng Qiu, and Zheng Zhang. 2020. [GenWiki: A dataset of 1.3 million content-sharing text and graphs for unsupervised graph-to-text generation](#). In *Proceedings of the 28th International Conference on Computational Linguistics*, pages 2398–2409, Barcelona, Spain (Online). International Committee on Computational Linguistics.
- Daniel Martin Katz, Michael J. Bommarito, II, and Josh Blackman. 2017. [A general approach for predicting the behavior of the supreme court of the united states](#). *PLOS ONE*, 12(4):1–18.
- Daniel Martin Katz, Dirk Hartung, Lauritz Gerlach, Abhik Jana, and Michael James Bommarito. 2023. [Natural language processing in the legal domain](#). *ArXiv*, abs/2302.12039.
- Arshdeep Kaur and Bojan Bozic. 2019. [Convolutional neural network-based automatic prediction of judgments of the european court of human rights](#). In *Irish Conference on Artificial Intelligence and Cognitive Science*.
- Rik Koncel-Kedziorski, Dhanush Bekal, Yi Luan, Mirella Lapata, and Hannaneh Hajishirzi. 2019. [Text Generation from Knowledge Graphs with Graph Transformers](#). In *Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long and Short Papers)*, pages 2284–2293, Minneapolis, Minnesota. Association for Computational Linguistics.
- Tiffany H Kung, Morgan Cheatham, Arielle Medenilla, Czarina Sillos, Lorie De Leon, Camille Elepaño, Maria Madriaga, Rimel Aggabao, Giezel Diaz-Candido, James Maningo, et al. 2023. Performance of chatgpt on usmle: potential for ai-assisted medical education using large language models. *PLoS digital health*, 2(2):e0000198.
- Pan Lu, Swaroop Mishra, Tony Xia, Liang Qiu, Kai-Wei Chang, Song-Chun Zhu, Oyvind Taffjord, Peter Clark, and Ashwin Kalyan. 2022. Learn to explain: Multimodal reasoning via thought chains for science question answering. In *The 36th Conference on Neural Information Processing Systems (NeurIPS)*.
- Pedro Henrique Luz de Araujo, Teófilo E. de Campos, Renato R. R. de Oliveira, Matheus Stauffer, Samuel Couto, and Paulo Bermejo. 2018. Lener-br: A dataset for named entity recognition in brazilian legal text. In *Computational Processing of the Portuguese Language*, pages 313–323, Cham. Springer International Publishing.
- Masha Medvedeva, Michel Vols, and Martijn Wieling. 2020. [Using machine learning to predict decisions of the european court of human rights](#). *Artificial Intelligence and Law*, 28(2):237–266.
- Joel Niklaus, Ilias Chalkidis, and Matthias Stürmer. 2021. [Swiss-judgment-prediction: A multilingual legal judgment prediction benchmark](#). In *Proceedings of the Natural Legal Language Processing Workshop 2021*, pages 19–35, Punta Cana, Dominican Republic. Association for Computational Linguistics.

- Joel Niklaus, Veton Matoshi, Pooja Rani, Andrea Galassi, Matthias Stürmer, and Ilias Chalkidis. 2023. [LEXTREME: A multi-lingual and multi-task benchmark for the legal domain](#). In *Findings of the Association for Computational Linguistics: EMNLP 2023*, pages 3016–3054, Singapore. Association for Computational Linguistics.
- OpenAI. 2023. GPT-4 technical report. Technical report.
- Vasile Pais, Maria Mitrofan, Carol Luca Gasan, Vlad Coneschi, and Alexandru Ianov. 2021. [Named entity recognition in the Romanian legal domain](#). In *Proceedings of the Natural Language Processing Workshop 2021*, pages 9–18, Punta Cana, Dominican Republic. Association for Computational Linguistics.
- Gemma Team Morgane Riviere, Shreya Pathak, Pier Giuseppe Sessa, Cassidy Hardin, Surya Bhupatiraju, L'eonard Hussenot, Thomas Mesnard, Bobak Shahriari, Alexandre Ram'e, Johan Ferret, Peter Liu, Pouya Dehghani Tafti, Abe Friesen, Michelle Casbon, Sabela Ramos, Ravin Kumar, Charline Le Lan, Sammy Jerome, Anton Tsitsulin, Nino Vieillard, Piotr Stańczyk, Sertan Girgin, Nikola Momchev, Matt Hoffman, Shantanu Thakoor, Jean-Bastien Grill, Behnam Neyshabur, Alanna Walton, Aliaksei Severyn, Alicia Parrish, Aliya Ahmad, Allen Hutchison, Alvin Abdagic, Amanda Carl, Amy Shen, Andy Brock, Andy Coenen, Anthony Laforge, Antonia Paterson, Ben Bastian, Bilal Piot, Boxi Wu, Brandon Royal, Charlie Chen, Chintu Kumar, Chris Perry, Christopher A. Welty, Christopher A. Choquette-Choo, Danila Sinopalnikov, David Weinberger, Dimple Vijaykumar, Dominika Rogozi'nska, D. Herbison, Elisa Bandy, Emma Wang, Eric Noland, Erica Moreira, Evan Senter, Evgenii Eltyshev, Francesco Visin, Gabriel Rasskin, Gary Wei, Glenn Cameron, Gus Martins, Hadi Hashemi, Hanna Klimczak-Pluci'nska, Harleen Batra, Harsh Dhand, Ivan Nardini, Jacinda Mein, Jack Zhou, James Svensson, Jeff Stanway, Jetha Chan, Jin Zhou, Joana Carrasqueira, Joana Il-jazi, Jocelyn Becker, Joe Fernandez, Joost R. van Amersfoort, Josh Gordon, Josh Lipschultz, Joshua Newlan, Junsong Ji, Kareem Mohamed, Kartikeya Badola, Kat Black, Katie Millican, Keelin McDonell, Kelvin Nguyen, Kiranbir Sodhia, Kish Greene, Lars Lowe Sjoesund, Lauren Usui, L. Sifre, L. Heuermann, Leticia Lago, Lilly McNealus, Livio Baldini Soares, Logan Kilpatrick, Lucas Dixon, Luciano Martins, Machel Reid, Manvinder Singh, Mark Iverson, Martin Gerner, Mat Velloso, Mateo Wirth, Matt Davidow, Matt Miller, Matthew Rahtz, Matthew Watson, Meg Risdal, Mehran Kazemi, Michael Moynihan, Ming Zhang, Minsuk Kahng, Minwoo Park, Mofi Rahman, Mohit Khatwani, Natalie Dao, Nenshad Bardoliwalla, Nesh Devanathan, Neta Dumai, Nilay Chauhan, Oscar Wahltinez, Pankil Botarda, Parker Barnes, Paul Barham, Paul Michel, Pengchong Jin, Petko Georgiev, Phil Culliton, Pradeep Kuppala, Ramona Comanescu, Ramona Merhej, Reena Jana, Reza Rokni, Rishabh Agarwal, Ryan Mullins, Samaneh Saadat, S. Mc Carthy, Sarah Perrin, S'ebastien Arnold, Sebastian Krause, Shengyang Dai, Shruti Garg, Shruti Sheth, Sue Ronstrom, Susan Chan, Timothy Jordan, Ting Yu, Tom Eccles, Tom Hennigan, Tomás Kociský, Tulsee Doshi, Vishan Jain, Vikas Yadav, Vilobh Meshram, Vishal Dharmadhikari, Warren Barkley, Wei Wei, Wenming Ye, Woohyun Han, Woosuk Kwon, Xiang Xu, Zhe Shen, Zhitao Gong, Zichuan Wei, Victor Cotruta, Phoebe Kirk, Anand Rao, Minh Giang, Ludovic Peran, Tris Brian Warkentin, Eli Collins, Joelle Barral, Zoubin Ghahramani, Raia Hadsell, D. Sculley, Jeanine Banks, Anca Dragan, Slav Petrov, Oriol Vinyals, Jeffrey Dean, Demis Hassabis, Koray Kavukcuoglu, Cl'ement Farabet, Elena Buchatskaya, Sebastian Borgeaud, Noah Fiedel, Armand Joulin, Kathleen Kenealy, Robert Dadashi, and Alek Andreev. 2024. [Gemma 2: Improving open language models at a practical size](#). *ArXiv*, abs/2408.00118.
- Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten Sootla, Itai Gat, Xiaoqing Tan, Yossi Adi, Jingyu Liu, Tal Remez, Jérémy Rapin, Artyom Kozhevnikov, I. Evtimov, Joanna Bitton, Manish P Bhatt, Cristian Cantón Ferrer, Aaron Grattafori, Wenhan Xiong, Alexandre D'efosse, Jade Copet, Faisal Azhar, Hugo Touvron, Louis Martin, Nicolas Usunier, Thomas Scialom, and Gabriel Synnaeve. 2023. [Code llama: Open foundation models for code](#). *ArXiv*, abs/2308.12950.
- Freda Shi, Daniel Fried, Marjan Ghazvininejad, Luke Zettlemoyer, and Sida I. Wang. 2022. [Natural language to code translation with execution](#). In *Proceedings of the 2022 Conference on Empirical Methods in Natural Language Processing*, pages 3533–3546, Abu Dhabi, United Arab Emirates. Association for Computational Linguistics.
- Don Tuggener, Pius von Däniken, Thomas Peetz, and Mark Cieliebak. 2020. [LEDGAR: A large-scale multi-label corpus for text classification of legal provisions in contracts](#). In *Proceedings of the Twelfth Language Resources and Evaluation Conference*, pages 1235–1241, Marseille, France. European Language Resources Association.
- Stefanie Urchs., Jelena Mitrović., and Michael Granitzer. 2021. [Design and implementation of german legal decision corpora](#). In *Proceedings of the 13th International Conference on Agents and Artificial Intelligence - Volume 2: ICAART*, pages 515–521. INSTICC, SciTePress.
- Chaojun Xiao, Haoxiang Zhong, Zhipeng Guo, Cunchao Tu, Zhiyuan Liu, Maosong Sun, Yansong Feng, Xianpei Han, Zhen Hu, Heng Wang, and Jianfeng Xu. 2018. [Cail2018: A large-scale legal dataset for judgment prediction](#). *ArXiv*, abs/1807.02478.
- Hai Ye, Xin Jiang, Zhunchen Luo, and Wenhan Chao. 2018. [Interpretable charge predictions for criminal cases: Learning to generate court views from fact descriptions](#). In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pages 1854–1864,

New Orleans, Louisiana. Association for Computational Linguistics.

Abhay Zala, Han Lin, Jaemin Cho, and Mohit Bansal. 2023. Diagrammergpt: Generating open-domain, open-platform diagrams via llm planning.

Tianyi Zhang, Varsha Kishore, Felix Wu, Kilian Q. Weinberger, and Yoav Artzi. 2020. [Bertscore: Evaluating text generation with BERT](#). In *8th International Conference on Learning Representations, ICLR 2020, Addis Ababa, Ethiopia, April 26-30, 2020*. OpenReview.net.

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

## A Qualitative analysis input

The English legal text used the qualitative analysis in Section 5.3 is the following:

On 28 May 2019, the procureur d'État du Luxembourg (State Prosecutor, Luxembourg) requested that the chambre du conseil du tribunal d'arrondissement de Luxembourg (Investigation Chamber of the District Court, Luxembourg) declare that JR should be surrendered to the French authorities. By order of 19 June 2019, the chambre du conseil du tribunal d'arrondissement de Luxembourg (Pre-trial Chamber of the District Court, Luxembourg) stated that it had no jurisdiction to hear JR's application for a declaration that the European arrest warrant was invalid and granted the request for the surrender of JR to the French authorities.

## B Detailed experimental settings

For training of LLMs, we follow the default setting of Hugging Face supervised finetuning of the `trl`<sup>6</sup> library for the optimizers and schedulers. We use the mini-batch size of 32. We use the max token length of 4096 for training as we notice some languages, e.g., Greek, require longer tokens than other languages depending on Llama tokenizers. In finetuning, we use FP32 precision and all trainable parameters are updated. All Llama-family experiments are done on a single node with four NVIDIA A100 GPUs.

## C Results of the Validation split

Table 9 shows the results of the validation split of the main performance table of Table 3.

## D Additional Multilingual Experiments

We experimented several models listed in the tables below and selected models with great scores are discussed in main paper.

**Multilingual results of all models.** We conducted experiments using the following models, namely, Llama3.1-8B, Llama3.1-8B-Instruct, Llama3.2-3B, Llama3.2-3B-Instruct, CodeLlama-7B, CodeLlama-7B-Instruct, CodeLlama-13B, CodeLlama-13B-Instruct, Gemma2-9B, and Gemma2-9B-it. Firstly, the results evaluated with legal content evaluation point of view are shown in Table 5 and Table 6. Table 5 shows few-shot Legal Content evaluation of all models conducted and Table 6 shows finetuning results of Legal Content evaluation.

Secondly, the results evaluated with graph-based point of view are summarized in Table 7 and Table 8. Table 7 shows few-shot results of graph-based evaluation and Table 8 represents.

For comparison of each language's score, we calculated the average of all few-shot models and the average of few-shot models despite GPT models in Table 5 and Table 7. In the same way, the average score of all finetuned models is calculated in Table 6 and Table 8.

## E Applications of traditional NLP tasks for legal domain

**Legal information extraction.** Named Entity Recognition (NER) is a fundamental information extraction task that has been developed for several languages, including Greek (Angelidis et al., 2018), Brazilian (Luz de Araujo et al., 2018), Romanian (Pais et al., 2021), and Spanish (de Gibert Bonet et al., 2022). Those NER approaches extract mainly the same objects as those in non-legal domains, while some efforts try to extract legal entities from court documents (II et al., 2021). Once NER identified entities, Relation Extraction in the legal domain (Chalkidis et al., 2021b) takes this information further by identifying and classifying the relationships between these entities, such as facts and allegedly violated articles, specific articles and paragraphs, and case references, as well as relevant facts and allegations.

**Legal classification.** The classification task of legal texts has been proposed with a focus on practical applications. For example, to enhance the interpretation of complex legal information, multi-label classification of legal texts assigns multiple conceptual class labels to words appearing in legal sentences (Chalkidis et al., 2019). Notably, FairLex (Chalkidis et al., 2022b) aims to ensure the fair application of the law by classifying attributes such as age, gender, region, and state.

<sup>6</sup><https://github.com/huggingface/trl>









Model	Graph-based			Valid Graph		Legal Content			
	G	G-N	G-N-E	Top1	Top10	Entity	R & T	Source	Statement
<i>Few-shot result of Validation split</i>									
CodeLlama 7B	14.46	10.14	4.73	18.35	86.96	47.03	6.51	12.44	1.78
CodeLlama 7B it.	16.70	12.21	6.64	43.30	91.91	52.18	8.91	17.67	1.71
CodeLlama 13B	14.99	10.53	5.17	18.09	84.96	47.58	7.14	14.93	2.76
CodeLlama 13B it.	18.02	13.33	6.82	38.26	89.74	54.30	8.98	18.65	3.75
Llama3.1 8B	24.87	19.05	10.07	36.17	87.04	60.77	12.74	23.86	4.17
Llama3.1 8B it.	26.88	20.00	11.42	26.96	88.00	60.86	15.02	29.35	3.77
Llama3.2 3B	21.04	16.00	7.92	31.91	85.83	54.39	10.25	21.69	2.89
Llama3.2 3B it.	26.45	20.21	11.45	61.04	93.83	53.56	14.78	20.63	3.69
Gemma2 9B	12.22	8.80	4.08	42.61	96.35	52.40	5.38	14.83	2.84
Gemma2 9B it.	27.66	21.92	13.37	71.04	96.87	69.11	16.81	30.67	3.36
GPT-3.5-Turbo	31.45	25.62	16.41	96.87	<b>100.0</b>	69.80	19.78	<b>29.43</b>	<b>5.48</b>
GPT-4	<b>33.09</b>	<b>27.34</b>	<b>19.15</b>	<b>98.96</b>	<b>100.0</b>	<b>75.31</b>	<b>23.24</b>	21.52	3.30
GPT-4o	30.98	25.79	17.24	95.74	<b>100.0</b>	71.61	20.40	40.72	4.77
<i>Finetuning result of Validation split</i>									
CodeLlama 7B	36.38	29.66	21.57	95.39	<b>100.0</b>	75.47	26.10	48.51	7.22
CodeLlama 7B it.	35.86	29.40	21.91	96.87	99.91	76.07	26.22	46.18	5.92
CodeLlama 13B	35.94	29.06	20.05	97.48	<b>100.0</b>	74.38	24.36	50.47	8.15
CodeLlama 13B it.	33.81	27.63	19.91	97.57	<b>100.0</b>	75.96	24.07	44.32	6.27
Llama3.1 8B	30.30	21.78	14.77	94.52	<b>100.0</b>	64.19	20.29	31.44	3.75
Llama3.1 8B it.	29.14	20.92	13.53	85.30	99.74	66.77	18.21	41.31	5.39
Llama3.2 3B	31.56	24.34	17.48	93.57	<b>100.0</b>	71.02	22.31	47.09	5.65
Llama3.2 3B it.	31.70	24.06	16.13	90.00	99.83	68.65	20.94	43.29	4.85
Gemma2 9B	43.62	37.32	<b>28.37</b>	<b>99.13</b>	<b>100.0</b>	<b>79.31</b>	<b>32.65</b>	<b>58.96</b>	8.99
Gemma2 9B it.	<b>43.88</b>	<b>36.57</b>	27.36	98.17	<b>100.0</b>	77.72	32.41	53.20	<b>10.51</b>

Table 9: Scores of the legal text visualization. **G**, **G-N** and **G-N-E** denote Graph, Graph&Node and Graph&Node&Edge respectively. Valid Graph Ratio is success rate of creating valid graphs in top-1 and top-10 generated results. The highest scores of each column are in bold.

**Legal summarization.** As a more complex and application-oriented task, legal summarization is also prominent in the field, which aims to generate a summary of legal sentences. Existing summarization studies address Canadian legal cases (Elaraby and Litman, 2022) and EU legislations (Aumiller et al., 2022).

## F Legal Diagram Formalism

Here we define several rules to express legal relations within the DOT language grammar.

**Graph node rules.** Legal entities are represented by nodes (vertices) in DOT languages with the shape of double octagons except legally deceased persons who are presented in the shape of ellipses. Legal norms that are effective in the present case are represented by graph nodes with trapezium shapes.

**Graph edge rules.** Legal transactions and the explanatory relationships between legal entities are represented by directed edges. The family or marital relationships established under civil law are represented by an undirected bold edge. The legal rights that cannot be exercised are represented by dashed edges. Dotted edges denote relationships of the legal succession between legal entities.

To illustrate the equivalent relationship between diagram nodes, undirected edges are used to connect entities and their status explanations, rules and statements, legal transactions, and their explanations.

We also note that legal relations can also be represented by graph nodes when legal relations have some relations with other entities. Figure 4 explains how to draw graphs when additional description is required for graph relations. In Graphviz, we cannot draw lines directly to the graph relations. Hence we change graph labels to nodes and connect to other nodes for adding explanation. Further details of the DOT language grammar for representations of legal entity relations and an actual dataset example are provided in Appendix G & J.

## G Graphviz annotation

The following is an example of the Graphviz code annotation rules.

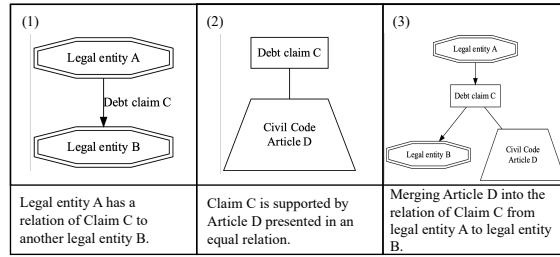


Figure 4: Annotation rule when adding explanation to graph relations.

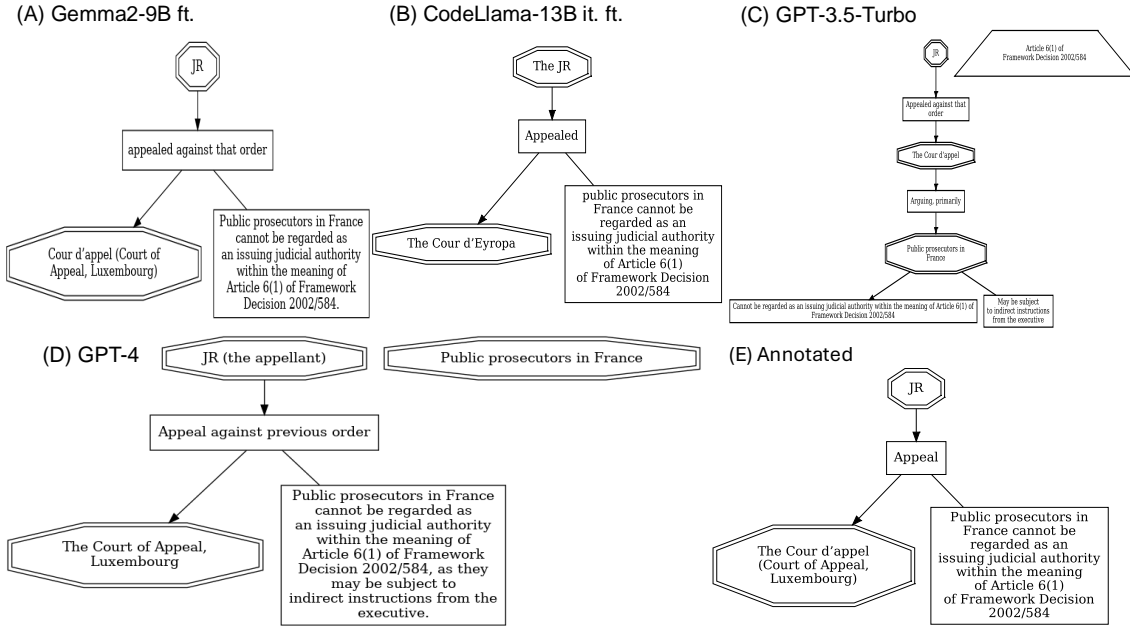


Figure 5: Additional qualitative analysis.

- 1 [shape=doubleoctagon]: Entities which are capable to act as legal entity.
- 2 [shape=trapezium]: Any kinds of rules which are legally effective, applied to the present case or supporting legal statements.
- 3 [style=dotted]: Relationship of succession between 2 entities.
- 4 [dir=none]: Equivalent relationship, agreements, or connecting detailed explanation of other nodes.
- 5 [dir=none, style=bold]: Marital relationships or family relationships which have been established under civil law.
- 6 [style=dashed]: Expressing a legal right that cannot be exercised or not existed.
- 7 [shape=ellipse]: Expressing a person who is legally deceased.

## H Qualitative Analysis

Figure 5 shows an additional case of qualitative analysis. Here, GPT-3.5-Turbo and GPT-4 failed to generate correct relational graph since some nodes (“Public prosecutors in France” and “Article 6(1) of Framework Decision 2002/584”) are not connected to other nodes lacking understandings of legal relations. On the other hand, Gemma2-9B ft. and CodeLlama13B-Instruct ft. models fine-tuned with LegalViz output almost the same diagram as annotated data, with the same legal entities (“JR” and “The Court of Appeal in Luxembourg”) and correctly extracted reason of appeal.

## I Prompt

The prompt for LLMs used in training, generation and dataset creation is presented in Table 10.

Method	Prompt
Prompt used for train and generation	Using the DOT language of Graphviz, draw a graph to explain legal entity nodes, legal relationships, legal statements and legal basis of them from given text, written in {language} text. Use “shape=trapezium” to represent a legally effective material and use “shape=doubleoctagon” to represent a legal entity in Graphviz code with {language}. At any time, reply only with the graphviz code.
Prompt for extraction	From legal text below of {language} language, extract the same meaning word or sentence as given English word to language. Please output only extracted result. Legal text: {legal text} Word or sentence to extract:
Prompt for translation	Translate below words or text from English to {language} Text:

Table 10: The prompts used in the experiment and data processing. {legal text} and {language} indicate the place to insert.

## J Train dataset examples

### Dataset Example

```
{'ID': '45',
 'category': 'EU law',
 'diagram_number': '7',
 'case_name': 'Case T-207/02: Nicoletta Falcone v Commission of the\nEuropean
Communities',
 'case_number': 'C2005/006/64',
 'document_url': 'https://eur-lex.europa.eu/legal-content/EN/TXT/PDF/?uri=CELEX:C200
5/006/64&qid=1713891140330',
 'year': '2004',
 'text': 'In Case T-207/02: Nicoletta Falcone, a candidate in Competition COM/A/10/0
1, represented by M. Condinanzi, against Commission of the European Communities
(Agent: J. Currall, assisted by A. Dal Ferro, with an address for service in
Luxembourg) - application for annulment of the decision of 2 May 2002 of the
selection board in Competition COM/A/10/01 to exclude the applicant from the
written tests on the ground that she did not obtain sufficient marks to be
included among the 400 best candidates - the Court of First Instance (Second
Chamber), composed of J. Pirrung, President, A.W.H. Meij and N. Forwood, Judges;
H. Jung, Registrar, has given a judgment on 26 October 2004, in which it:',
 'Graphviz': 'digraph {\n  rankdir=LR;\n  node [shape=box];\n\n  "Nicoletta
Falcone" -> "M. Condinanzi" [label="represent" dir=none];\n  "The Comission of
the European Comminties" -> "Nicoletta Falcone" [label="application for
annulment of the decision of 2 May 2002 of the selection board in Competition
COM/A/10/01 to exclude the applicant from the written tests on the ground that
she did not obtain sufficient marks to be included among the 400 best candidates
"]; \n}',
 'language': 'English'
}
```