









# PAKTON: A MULTI-AGENT FRAMEWORK FOR QUESTION ANSWERING IN LONG LEGAL AGREEMENTS

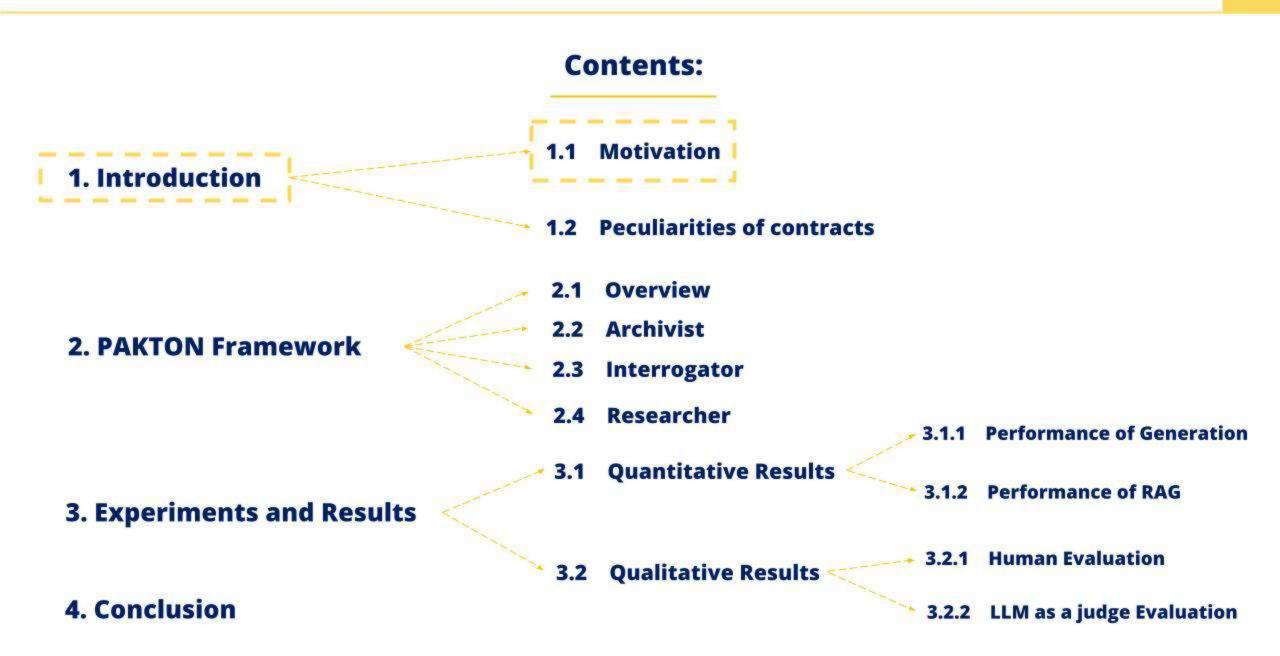
Petros Raptopoulos<sup>1</sup>, Giorgos Filandrianos<sup>1,2</sup>, Maria Lymperaiou<sup>1</sup>, Giorgos Stamou<sup>1</sup>

<sup>1</sup>School of Electrical and Computer Engineering, AILS Laboratory National Technical University of Athens, Greece <sup>2</sup>Instituto de Telecomunicações, Portugal petrosrapto@gmail.com, {geofila, marialymp}@islab.ntua.gr, gstam@cs.ntua.gr









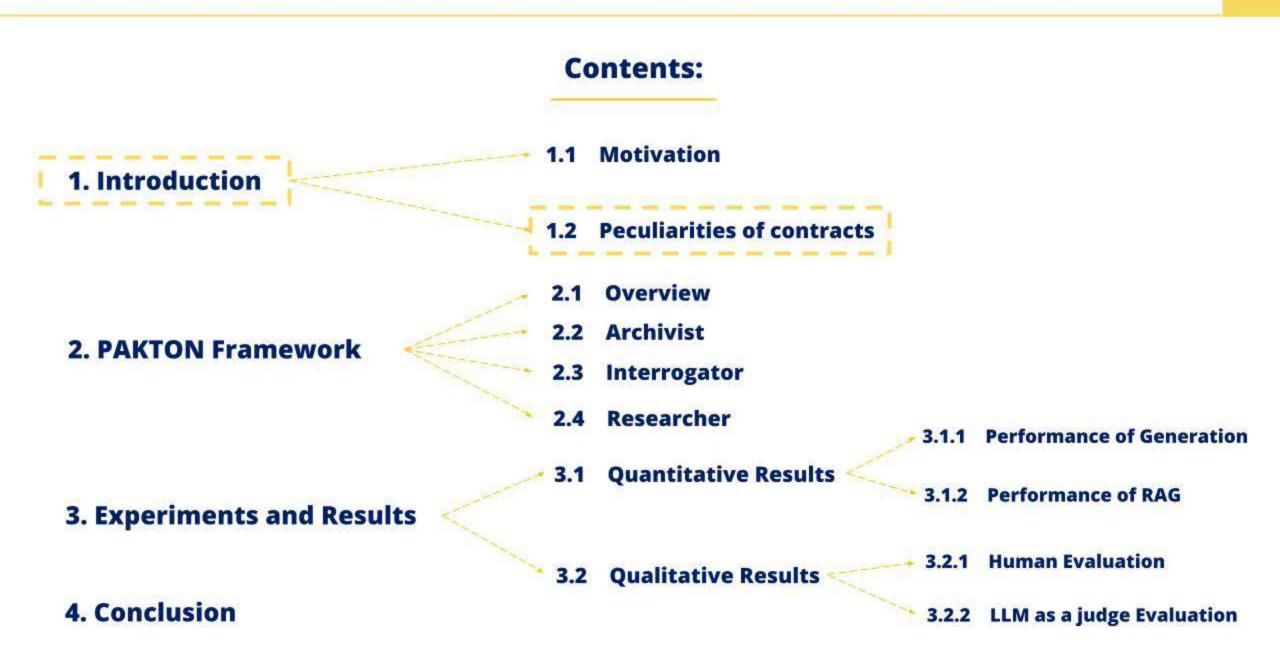
#### 1. Introduction - 1.1 Motivation

## **Motivation**

- \* Contract review: Complex, time-consuming, requiring specialized legal expertise
- \* Inaccessible to non-experts (general public), demanding even for professionals
- \* High economic impact: avg. 9.2% annual revenue loss
- \* High volume: large organizations manage 20k-40k active contracts
- \* Confidentiality concerns make proprietary models unsuitable







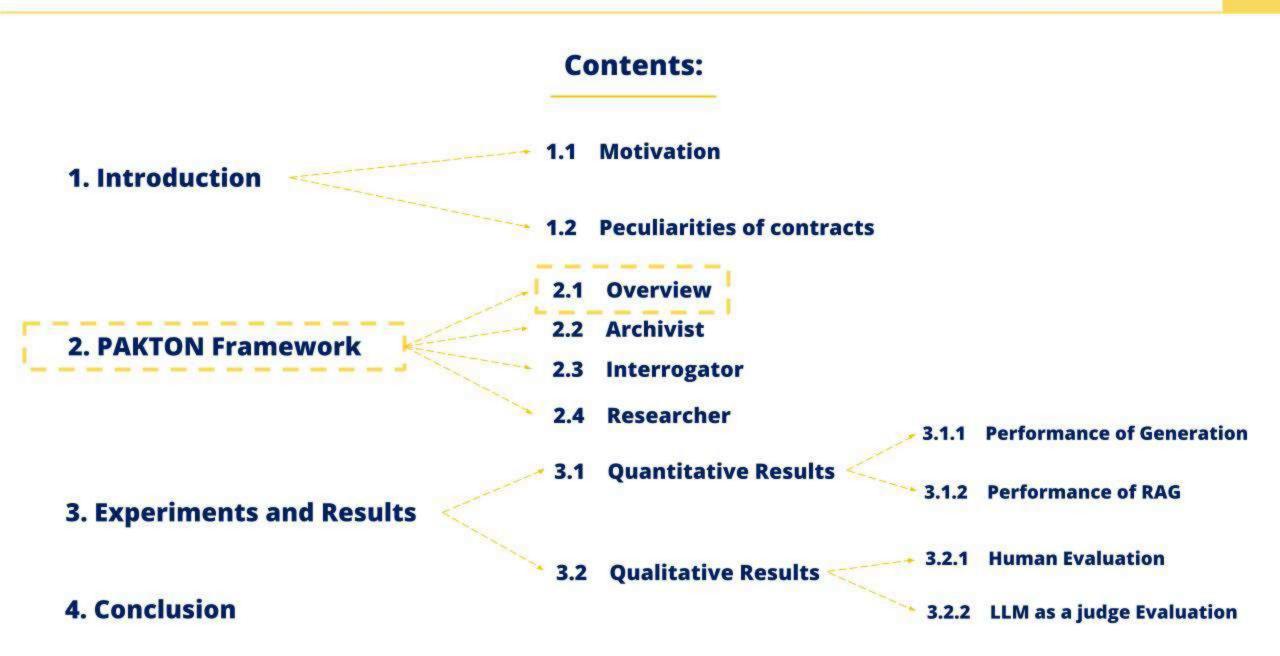
#### 1. Introduction - 1.2 Peculiarities of contracts

### **Peculiarities of contracts**

- \* Complex legal language & terminology
- \* Overlapping or contradictory clauses
- \* Frequent exceptions & cross-references
- \* Substantial length
- \* Ambiguous wording & multiple interpretations
- \* Jurisdictional differences







#### 2. PAKTON Framework - 2.1 Overview

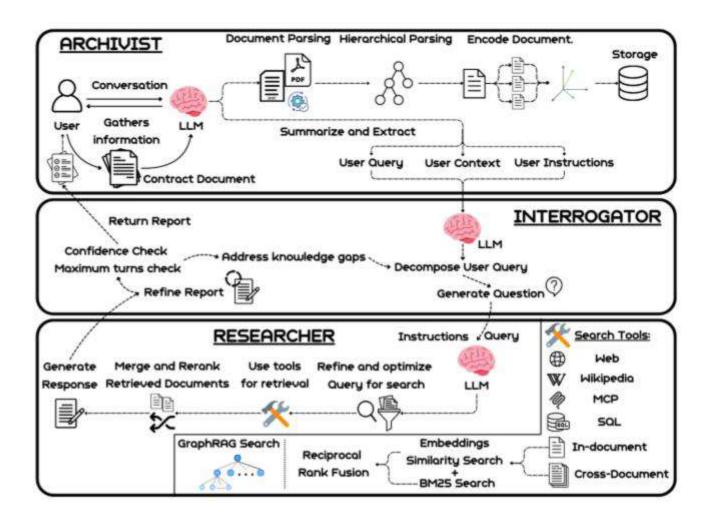
- ✓ Task: Answer the questions of the user based on the contract provided and integrated external knowledge.
- Output: Legal Report
- @ Goals:
  - Transparent and traceable reasoning.
  - Reference evidence spans from the contract.





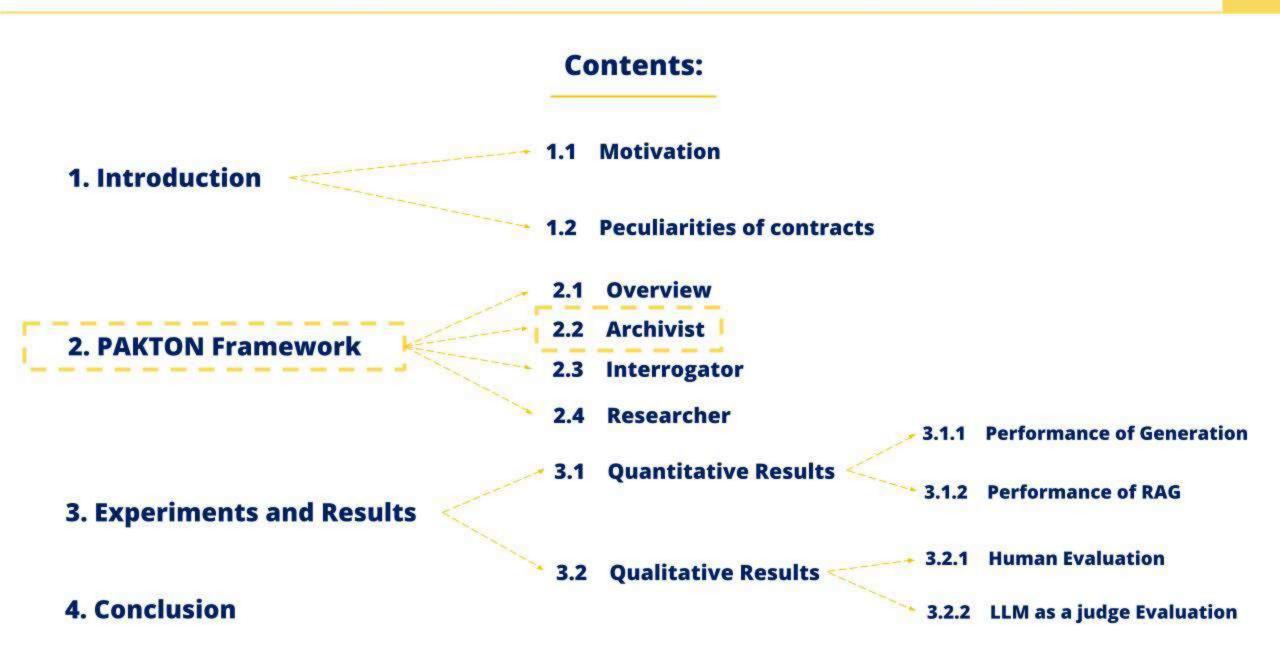


#### 2. PAKTON Framework - 2.1 Overview









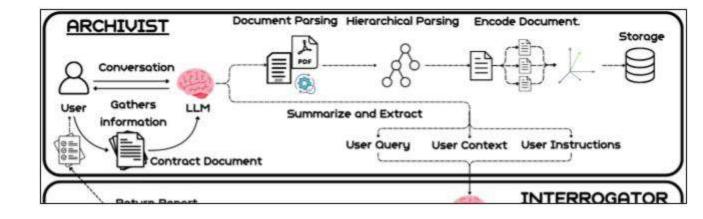
#### 2. PAKTON Framework - 2.2 Archivist

### in Role:

- Interaction with the user
- Gather/organize information
- Storage and embedding management of contracts

### \* Functionalities:

- Document processing
- Hierarchical Representation
- Document Segmentation and embedding







#### 2. PAKTON Framework - 2.2 Archivist

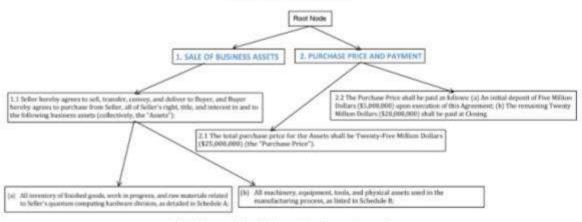
#### 1. SALE OF BUSINESS ASSETS

- 1.1 Seller hereby agrees to sell, transfer, convey, and deliver to Buyer, and Buyer hereby agrees to purchase from Seller, all of Seller's right, title, and interest in and to the following business assets (collectively, the "Assets"):
- (a) All inventory of finished goods, work in progress, and raw materials related to Seller's quantum computing hardware division, as detailed in Schedule A;
- (b) All machinery, equipment, tools, and physical assets used in the manufacturing process, as listed in Schedule B;

#### 2. PURCHASE PRICE AND PAYMENT

- 2.1 The total purchase price for the Assets shall be Twenty-Five Million Dollars (\$25,000,000) (the "Purchase Price").
- 2.2 The Purchase Price shall be paid as follows: (a) An initial deposit of Five Million Dollars (\$5,000,000) upon execution of this Agreement; (b) The remaining Twenty Million Dollars (\$20,000,000) shall be paid at Closing.

#### (a) Section Detection



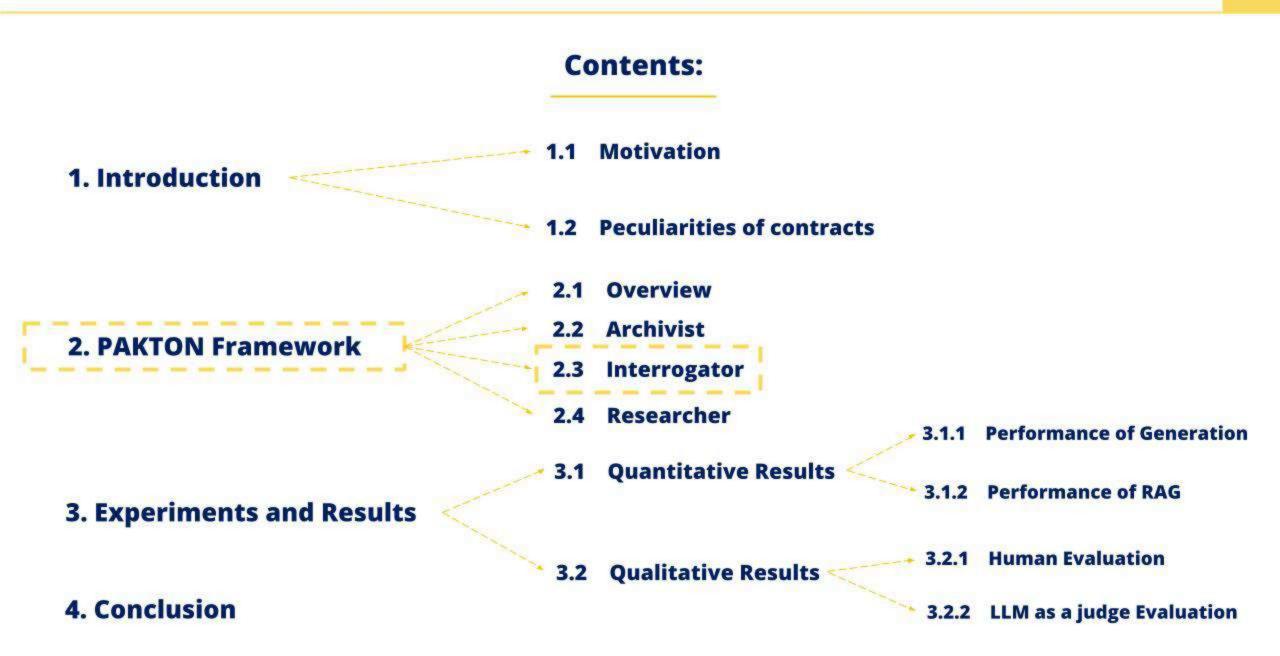
#### (b) Hierarchical Organization of sections.



(c) Contextual Embeddings for node "1.1 Seller ..."







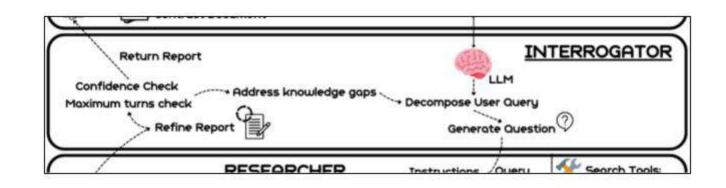
### 2. PAKTON Framework - 2.3 Interrogator

### in Role:

- Generate final report
- Orchestrate multi-step reasoning process
- Initiate interrogation of the Researcher

#### \* Functionalities:

- Decomposition of the original query
- Generate a series of questions
- Refine report at each step



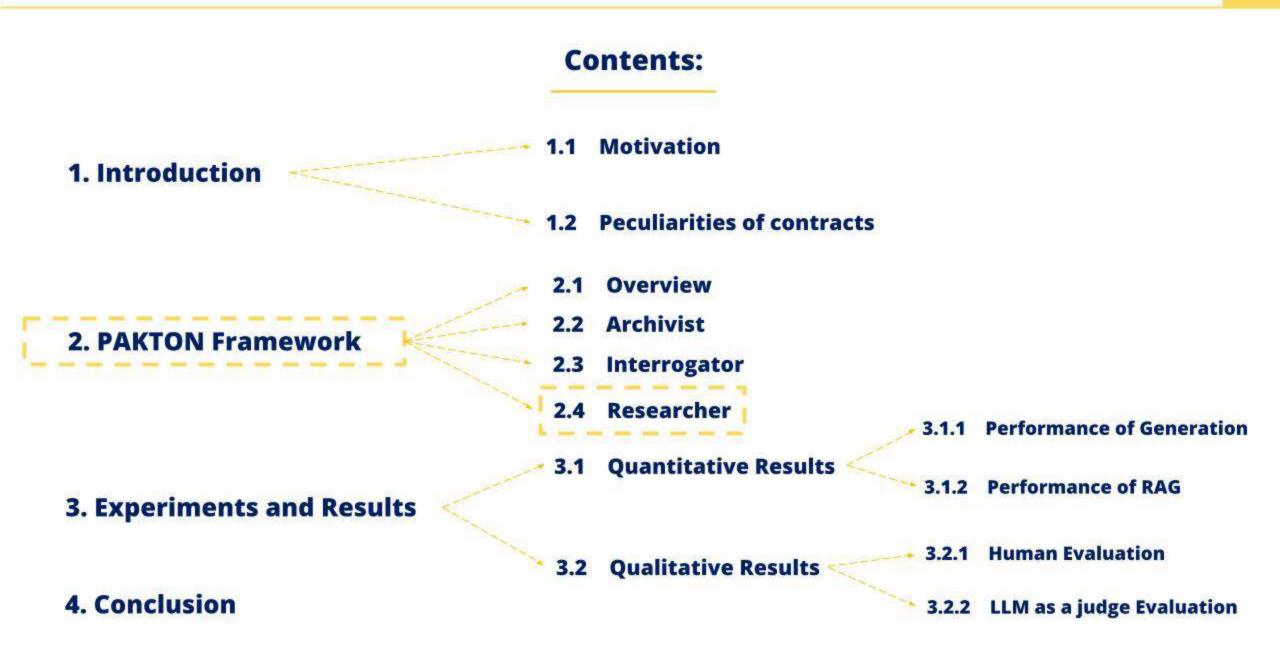
# Report Structure:

- 1) Title and Topic Summary
- 2) Legal Reasoning and Key findings
- 3) Preliminary answer and suggested research directions
- 4) Knowledge gaps and follow-up questions
- 5) Cited sources and evidentiary support





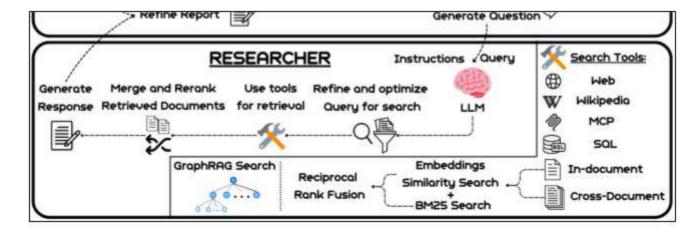




#### 2. PAKTON Framework - 2.4 Researcher

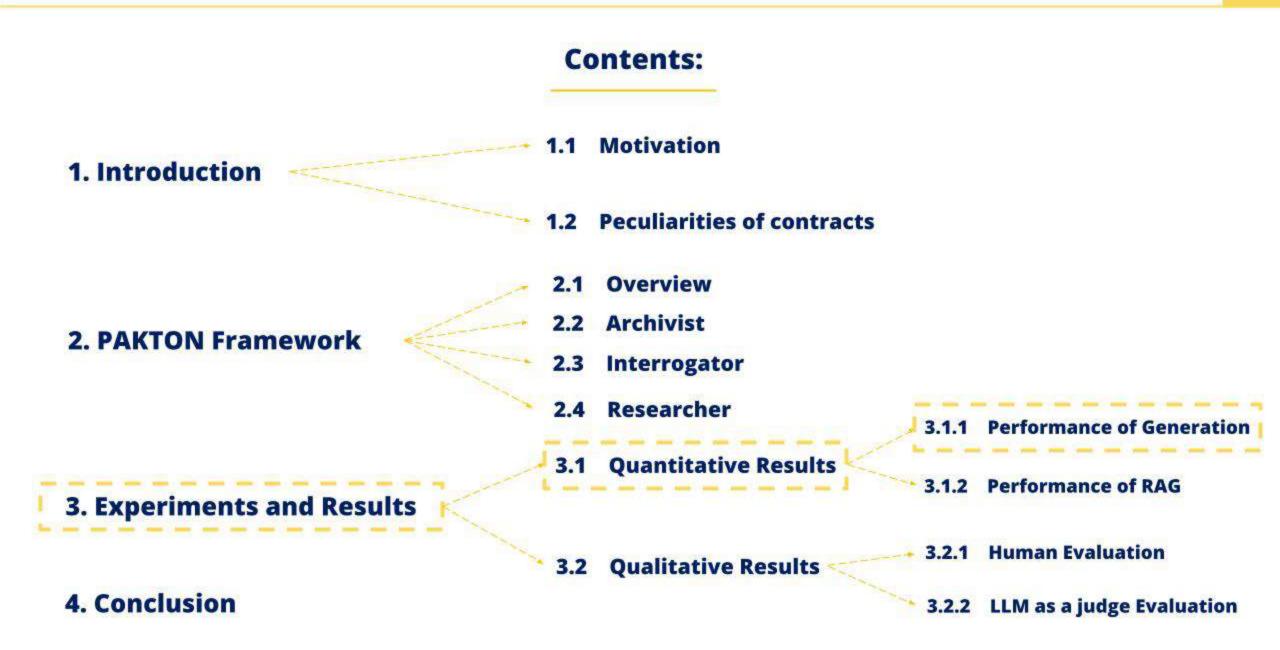


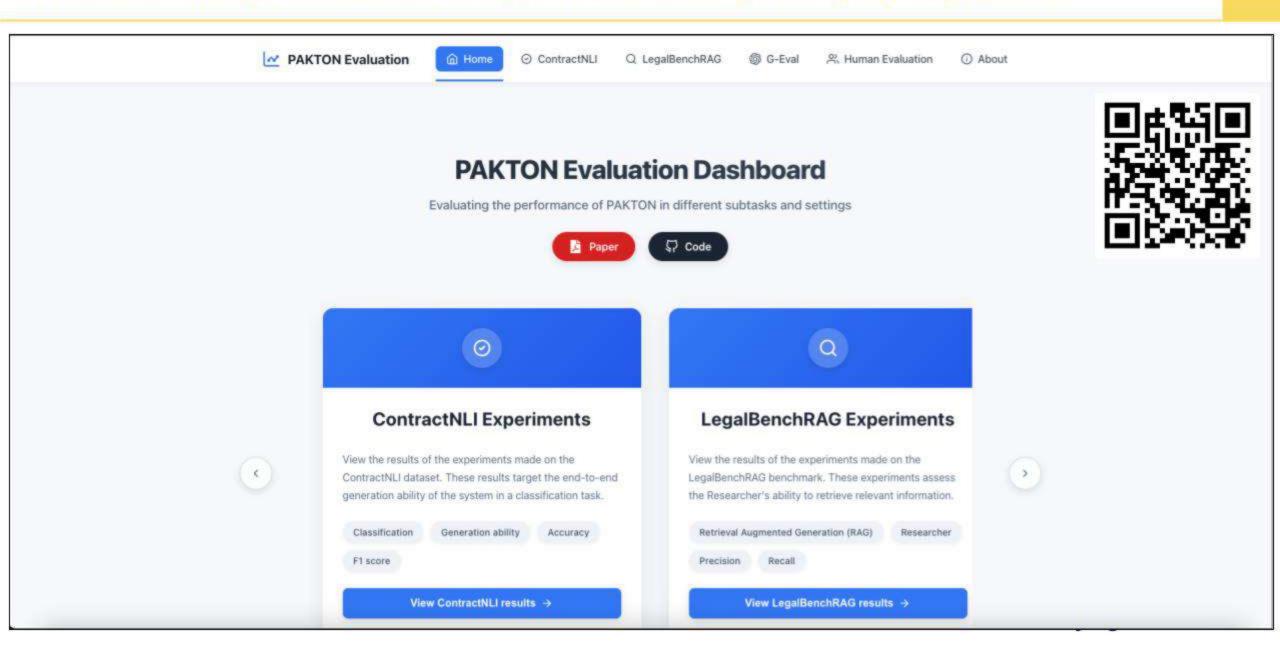
- Retrieve relevant information given a query
- Select the most suitable retrieval method or combination of methods









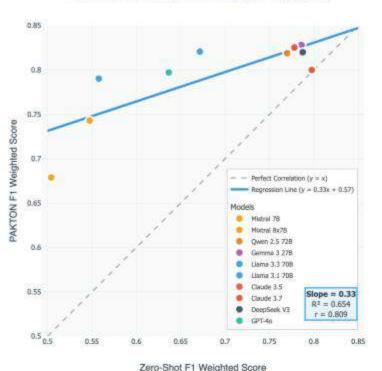




### 3. Experiments and Results - 3.1 Quantitative Results - 3.1.1 Performance of Generation

- \* Assessment of the Generation ability of the end-to-end system
- \* ContractNLI dataset (premise, hypothesis, classification)
- \* PAKTON consistently outperforms the baseline methods across all models

#### Zero-Shot vs PAKTON F1 Weighted Performance

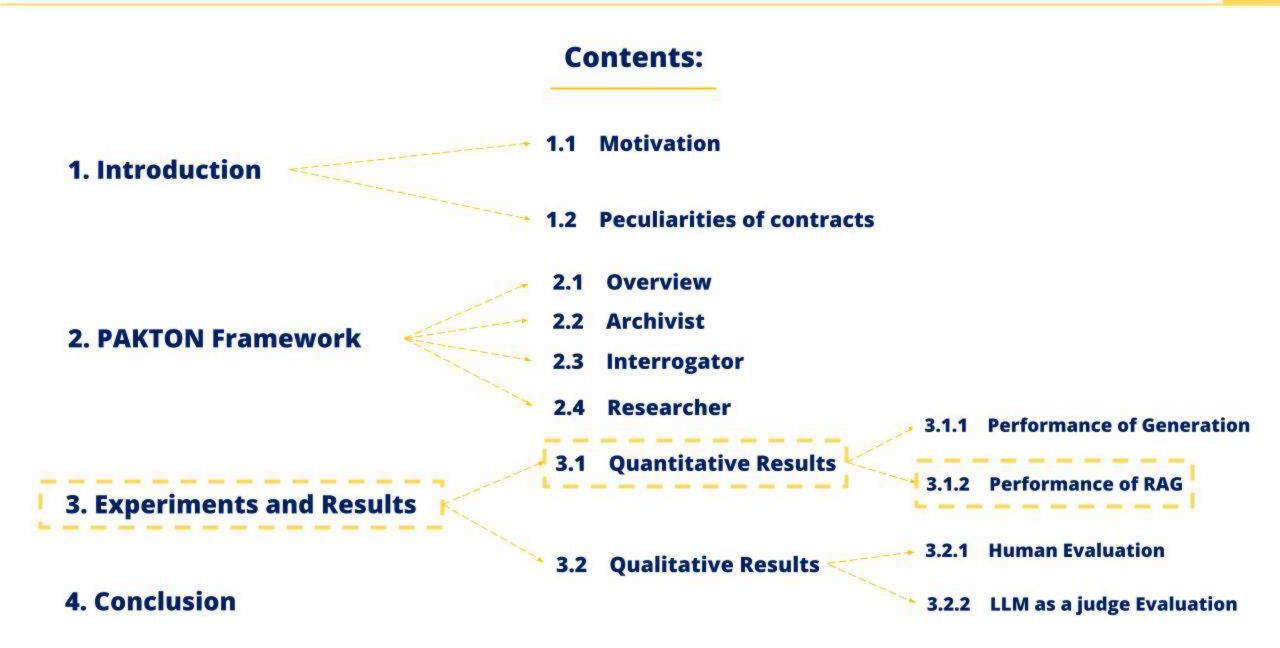


#### **Robustness analysis:**

- Coefficient of variation (CV): 5.7% across all models for PAKTON
- ZS scores exhibit CV: 16.86%, PAKTON reduces CV by ~66%
- ANOVA: p ~= 0.075, F statistic ~= 4.19

Model	Method	Acc.		F1 [E]		
Saul7B	Inst. Tun.	0.4196	0.2900	0.0589	100000000000000000000000000000000000000	Decision In terms on
Saul54B	Inst. Tun.	0.7020	0.6792	0.7727	0.1729	0.7024
	ZS	0.5364	0.5042	0.5279	0.0248	0.5951
Mistral	FS	0.5065	0.4702	0.6053	0.0082	0.4379
And a Contract of the Contract	FS+Spans	0.4940	0.4576	0.6085	0.0076	0.4053
7B	PAKTON	0.7032	0.6789	0.7782	0.2469	0.6828
	ZS	0.5423	0.5475	0.6445	0.4103	0.4770
Mixtral	FS	0.6002	0.5804	0.6836	0.1931	0.5642
State of the same of	FS+Spans	0.6150	0.6017	0.6901	0.1951	0.6060
8x7B	PAKTON	0.7423	0.7429	0.7864	0.6655	0.7187
-	ZS	0.7728	0.7699	0.8248	0.5776	0.7579
Owen	FS	0.7351	0.7241	0.8094	0.4920	0.6892
2.5 72B	FS+Spans	0.7484	0.7432	0.8196	0.4378	0.7357
2.5 728	PAKTON	0.8192	0.8188	0.8353	0.7737	0.8132
	ZS	0.7886	0.7860	0.8316	0.6348	0.7739
Commo	FS	0.7191	0.7049	0.7815	0.4608	0.6891
Gemma	FS+Spans	0.7720	0.7639	0.8287	0.4728	0.7662
3 27B	PAKTON	0.8287	0.8283	0.8487	0.7546	0.8255
	ZS	0.6767	0.6716	0.7366	0.5378	0.6346
Llama	FS	0.6657	0.6565	0.7326	0.4431	0.6268
	FS+Spans	0.6915	0.6879	0.7382	0,4244	0.6982
3.3 70B	PAKTON	0.8217	0.8207	0.8422	0.7488	0.8165
	ZS	0.5811	0.5577	0.5216	0.3152	0.6555
WWGGGGG	FS	0.5729	0.5506	0.5421	0.2381	0.6358
Llama	FS+Spans	0.5538	0.5180	0.4471	0.3014	0.6468
3.1 70B	PAKTON	0.7916	0.7903	0.8097	0.6846	0.7960
	ZS	0.7916	0.7977	0.8757	0.5722	0.7691
Claude	FS	0.7778	0.7816	0.8588	0.5702	0.7505
	FS+Spans	0.7999	0.8034	0.8678	0.6046	0.7826
3.5	PAKTON	0.7990	0.8000	0.8157	0.7046	0.8072
	ZS	0.7704	0.7781	0.8633	0.5602	0.7398
Claude	FS	0.7590	0.7602	0.8463	0.5607	0.7165
3.7	FS+Spans	0.7724	0.7766	0.8538	0.5805	0.7417
3.7	PAKTON	0.8247	0.8254	0.8386	0.7495	0.8304
	ZS	0.7886	0.7875	0.8487	0.6117	0.7648
Deep-	FS	0.7681	0.7607	0.8346	0.6104	0.7182
seek	FS+Spans	0.7743	0.7714	0.8377	0.5812	0.7465
V3	PAKTON	0.8192	0.8200	0.8315	0.7615	0.8224
	ZS	0.6121	0.6366	0.7490	0.4162	0.5698
	FS	0.6640	0.6789	0.7372	0.4734	0.6666
GPT-40	FS+Spans	0.6482	0.6574	0.6664	0.4636	0.6950
011-10	PAKTON	0.7966	0.7972	0.7964	0.7592	0.8068

Table 1: Performance comparison of PAKTON and other methods across models on the ContractNLI test set. The highest accuracy and F1[w] are shown in **bold**.



#### 3. Experiments and Results - 3.1 Quantitative Results - 3.1.2 Performance of RAG

				Precis	ion @	b k			Recall @ k							
Dataset	Method	1	2	4	8	16	32	64	1	2	4	8	16	32	64	
	Naive	7.86	7.31	6.41	5.06	3.58	2.41	1.54	7.45	12.53	20.88	32.38	42.45	54.27	66.07	
D.:O.4	RCTS	14.38	13.55	12,34	9.03	6.06	4.17	2.81	8.85	15.21	27.92	42.37	55.12	71.19	84.19	
PrivacyQA	Naive + Cohere	14.38	13.55	12.34	9.02	6.06	4.17	2.81	8.85	15.21	27.92	42.37	55.12	71.19	84.19	
	RCTS + Cohere	13.94	15.91	13.32	9.57	6.88	4.68	3.28	7.32	16.12	25.65	35.60	51.87	64.98	79.61	
	PAKTON	19.94	16.84	11.44	8.62	7.38	6.42	6.08	13.34	22.43	32.67	43.39	61.65	82.30	89.42	
	Naive	16.45	14.80	12.53	9.73	6.70	4.65	3.04	11.32	19.10	29.79	45.59	56.75	69.88	86.57	
C	RCTS	6.63	5.29	3.89	2.81	1.98	1.29	0.90	7.63	11.33	17.34	24.99	35.80	46.57	61.72	
ContractNLI	Naive + Cohere	6.63	5.28	3.89	2.81	1.98	1.29	0.90	7.63	11.34	17.34	24.99	35.80	46.57	61.72	
	RCTS + Cohere	5.08	5.59	5.04	3.67	2.52	1.75	1.17	4.91	9.33	16.09	25.83	35.04	46.90	62.97	
	PAKTON	33.02	30.34	17.33	9.98	5.87	4.68	4.52	53.14	67.47	80.06	89.71	95.50	99.56	99.82	
MAUD	Naive	3.36	2.65	2.18	1.89	1.48	1.06	0.75	2.54	3.12	4.53	8.75	13.16	18.36	25.62	
	RCTS	2.65	1.77	1.96	1.40	1.39	1.15	0.82	1.65	2.09	4.59	6.18	12.93	21.04	28.28	
MAUD	Naive + Cohere	2.64	1.77	1.96	1.40	1.38	1.15	0.82	1.65	2.09	5.59	6.18	12.93	21.04	28.28	
	RCTS + Cohere	1.94	2.63	2.05	1.77	1.79	1.55	1.12	0.52	2.48	4.39	7.24	14.03	22.60	31.46	
	PAKTON	25.47	17.45	10.51	7.24	5.08	3.18	1.85	23.99	30.09	34.49	46.42	59.74	74.96	82.80	
	Naive	9.27	8.05	5.98	4.33	2.77	1.77	1.09	12.60	19.47	27.92	40.70	51.02	64.38	75.71	
CUAD	RCTS	1.97	4.03	4.83	4.20	2.94	1.99	1.25	1.62	8.11	17.72	31.68	44.38	60.04	74.70	
CUAD	Naive + Cohere	1.97	4.03	4.83	4.20	2.94	1.99	1.25	1.62	8.11	17.72	31.68	44.38	60.04	74.70	
	RCTS + Cohere	3.53	4.18	6.18	5.06	3.93	2.74	1.66	3.17	7.33	18.26	28.67	42.50	55.66	70.19	
	PAKTON	11.02	8.83	6.81	4.72	2.78	2.07	1.62	16.52	24.76	33.34	46.67	59.53	77.08	86.23	
ALL	Naive	2.40	3.76	4.97	4.33	3.39	2.17	1.29	3.37	8.44	21.30	34.51	48.88	64.47	76.39	
	RCTS	6.41	6.16	5.76	4.36	3.09	2.15	1.45	4.94	9.19	16.90	26.30	37.06	49.71	62.22	
	Naive + Cohere	6.41	6.16	5.76	4.36	3.09	2.15	1.45	4.94	9.19	16.90	26.30	37.05	49.71	62.22	
	RCTS + Cohere	6.13	7.08	6.65	5.02	3.78	2.68	1.81	3.98	8.82	16.10	24.34	35.86	47.54	61.06	
	PAKTON	22.34	18.37	11.52	7.63	5.26	4.08	3.52	26.77	36.32	45.26	56.66	69.17	83.50	89.58	

Table 2: Precision and Recall @  $k \in \{1, 2, 4, 8, 16, 32, 64\}$  for four retrieval pipelines on five legal-text datasets.

\* Assessment of the Retrieval ability \* LegalBench-RAG benchmark

**Task:** Given a contract and a question, did the system retrieve the "golden truth" spans? How precise and complete was the retrieval?

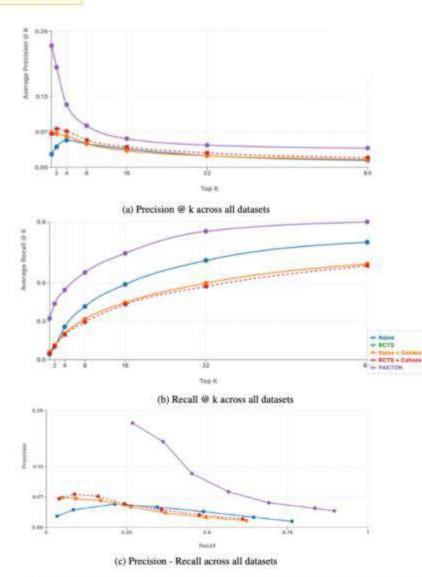


Figure 6: Precision and Recall values for different k across all datasets for all methods

#### 3. Experiments and Results - 3.1 Quantitative Results - 3.1.2 Performance of RAG

1	1	Precision @ k								Recall @ k							
Dataset	LLM filtering	1	2	4	8	16	32	64	1	2	4	8	16	32	64		
PrivacyQA	without	19.94	16.84	11.44	8.62	7.38	6.42	6.08	13.34	22.43	32.67	43.39	61.65	82.30	89.42		
riivacyQA	with	29.11	30.33	25.99	23.35	22.64	22.33	22.33	10.51	18.13	23.33	26.35	27.65	28.19	28.19		
ContractNLI	without	33.02	30.34	17.33	9.98	5.87	4.68	4.52	53.14	67.47	80.06	89.71	95.50	99.56	99.82		
Contractive	with	59.59	51.36	46.32	45.00	45.00	44.87	44.87	38.53	45.25	51.95	54.94	58.00	58.69	58.69		
MATTE	without	25.47	17.45	10.51	7.24	5.08	3.18	1.85	23.99	30.09	34.49	46.42	59.74	74.96	82.80		
MAUD	with	38.87	36.99	33.54	33.12	32.77	32.33	32.29	19.06	22.60	24.06	26.52	27.51	27.64	27.64		
CTIAD	without	11.02	8.83	6.81	4.72	2.78	2.07	1.62	16.52	24.76	33.34	46.67	59.53	77.08	86.23		
CUAD	with	29.14	29.53	29.31	28.86	28.89	28.79	28.77	25.31	30.30	34.33	37.68	38.24	38.68	38.68		
	without	22.34	18.37	11.52	7.63	5.26	4.08	3.52	26.77	36.32	45.26	56.66	69.17	83.50	89.58		
ALL	with	39.17	37.03	33.78	32.58	32.26	32.08	32.05	23.37	29.07	33.42	36.36	37.84	38.29	38.29		

Table 3: Performance comparison on different datasets for Precision and Recall at various k values for PAKTON's Researcher and Archivist under configuration 1.

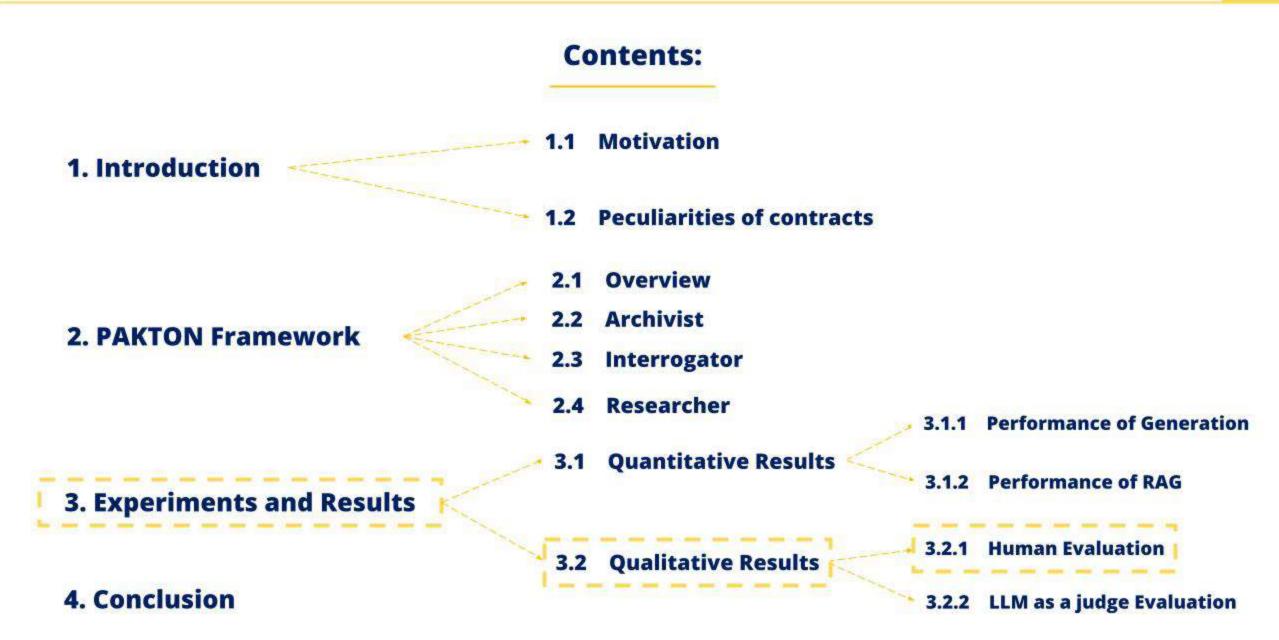
				Pre	cision	@ k			Recall @ k								
Dataset	LLM filtering	1	2	4	8	16	32	64	1	2	4	8	16	32	64		
PrivacyQA	without with	0.000	the second		10.39 26.66	1,000					450000000000000000000000000000000000000		79.24 28.20				
ContractNLI	without with								Dec. 200				<b>98.19</b> 65.07				
MAUD	without with				7 - 3 - 5 - 5 - 5 - 5 - 5								68.82 29.90		10.75 (0.75)		
CUAD	without with	VT1000			4.18 31.13	S-01500-							69.13 33.60				
ALL	without with								2000				78.845 39.19				

Table 4: Performance comparison across different datasets in terms of Precision and Recall at various k values, using PAKTON's Researcher and Archivist under Configuration 2.

#### **Ablations**

- \* Effect of Reranker on Performance
- \* Addition of LLM-filter step





### 3. Experiments and Results - 3.2 Qualitative Results - 3.2.1 Human Evaluation

Criterion	Instructions						
Explainability and Reasoning	Evaluate whether the report clearly and transparently explains not only the final conclusion, but also the reasoning process and supporting evidence in a step-by-step, understandable manner. The explanation should guide the reader through the logic in a way that supports comprehension, avoiding unexplained jumps in logic.						
Justification with Evidence	Determine whether the statements and claims are explicitly justified with relevant, specific, and clearly cited evidence (e.g., direct quotations, clause references). The justification should be traceable, allowing the reader to locate the original source material.						
Contextual and Legal Understanding	Assess whether the report demonstrates a deep and accurate understanding of the document, its legal terminology, and the broader context. Consider whether it correctly interprets clauses and captures implied assumptions or legal concerns behind the question.						
Handling Ambiguity	Determine whether the report identifies and handles ambiguities in the source material appropriately, such as by presenting multiple interpretations or justifying a chosen one clearly.						
Acknowledgment of Knowledge Gaps	Evaluate whether the report explicitly acknowledges when available information is insufficient to support a conclusion, avoiding speculation or overconfidence.						
Conciseness and Precision	Assess whether the report communicates clearly and efficiently, avoiding unnecessary repetition or verbosity, while still covering all key points.						
Coherence and Organization	Check whether the report is logically structured, flows smoothly, and maintains clarity across sections. Transitions between ideas should be natural and helpful.						
Relevance and Focus	Evaluate whether the report stays on topic and maintains focus on answering the core question, avoiding tangents or irrelevant content.						
Completeness	Assess whether the report addresses all important aspects of the question and offers a contextually broad and holistic answer. It should not omit any major points or perspectives.						

Table 16: Instructions given to human annotators for each evaluation criterion used in the PAKTON vs. ChatGPT comparison. Similar instructions were given to the G-EVAL framework.

- \* Assessment of the quality using Human Evaluators
- \* Head-to-Head comparison PAKTON vs GPT-4o (RAG)
- \* Nine criteria, 15 legal questions, 540 individual judgments
- \* PAKTON is favored over ChatGPT on the majority of evaluation dimensions.

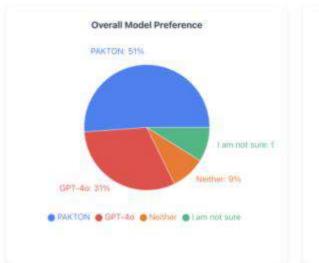


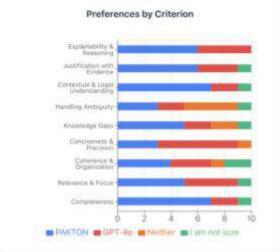




### 3. Experiments and Results - 3.2 Qualitative Results - 3.2.1 Human Evaluation

"If an application recognizes the emotions of a natural person, how is it classified according to the Regulation?"





\* PAKTON is favored over ChatGPT on the majority of evaluation dimensions.

\* ChatGPT is preferred for conciseness.

(a) Preference based on responses for a single question



(b) Overall Model Preference aggregated across all criteria and all questions

Figure 7: Comparative analysis of PAKTON vs. GPT-40 based on human evaluator judgments across different criteria



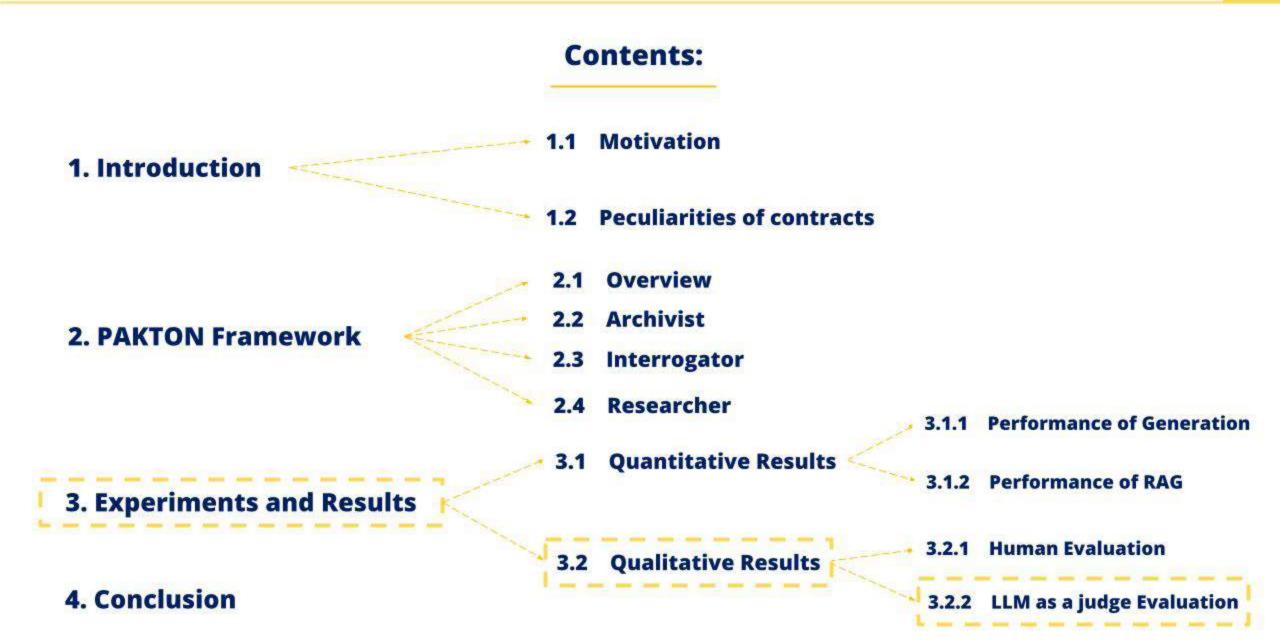
### 3. Experiments and Results - 3.2 Qualitative Results - 3.2.1 Human Evaluation





Figure 8: The user interface (UI) of PAKTON employed during the human evaluation with study participants.





#### 3. Experiments and Results - 3.2 Qualitative Results - 3.2.2 LLM as a judge Evaluation



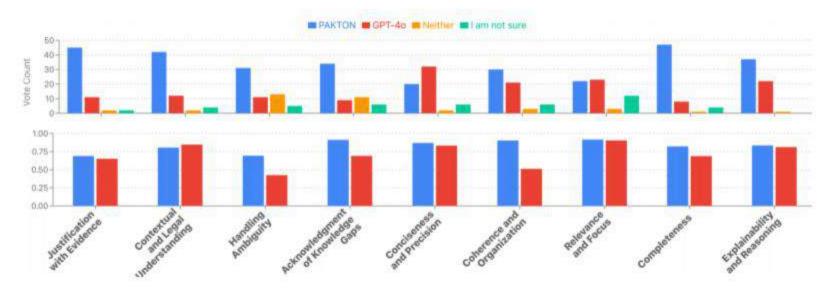


Figure 3: Comparison analysis of PAKTON and GPT-4o. Top plot presents human preferences across nine evaluation criteria aggregated for all questions. Bottom plot shows G-EVAL scores for the same criteria, aggregated across all ContractNLI outputs.

- \* Quality Assessment using G-Eval (LLM as a judge) \* Absolute scores across nine criteria (same as humans)
- \* Consistent results with humans, PAKTON is favored over ChatGPT on the majority of evaluation dimensions.





#### 4. Conclusion

### So, why PAKTON?

- Superior generation performance measured, compared to baseline approaches
- State-of-the-art performance achieved by the Researcher component (RAG).
- 6 PAKTON was preferred over ChatGPT by human evaluators for contract analysis.
- LLM-based evaluators demonstrated a clear preference for PAKTON over GPT-4o.

- Open-source pipeline
- Plug-and-play capabilities, end-to-end system
- Easy on-premise
- Explainability vs black box models













# Thanks for your attention!

For feedback, inquiries, or potential collaborations, feel free to reach out at <a href="mailto:petrosrapto@gmail.com">petrosrapto@gmail.com</a>





