

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



Petros Raptopoulos¹, Giorgos Filandrianos^{1,2}, Maria Lymperaiou¹, Giorgos Stamou¹

¹School of Electrical and Computer Engineering,

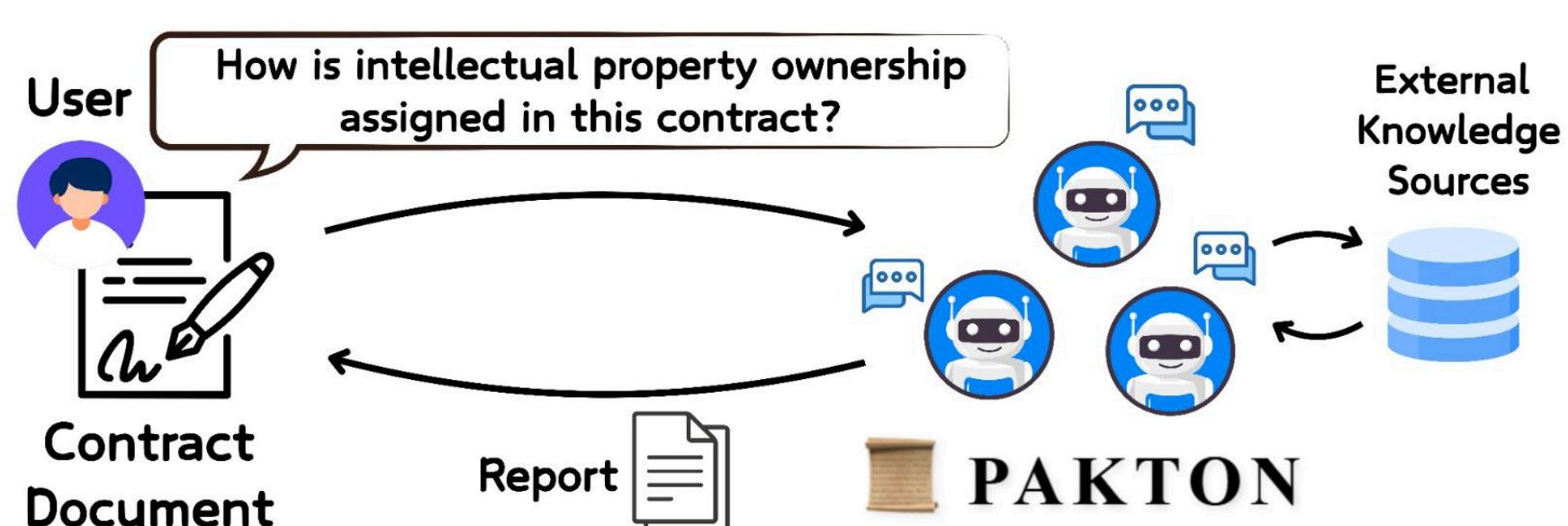
AILS Laboratory National Technical University of Athens, Greece

²Instituto de Telecomunicações, Portugal

petrosraptos@gmail.com, {geofila, marialymp}@islab.ntua.gr, g.stam@cs.ntua.gr



THE TASK



GENERATE ANSWERS TO THE USER'S QUESTIONS USING INFORMATION FROM THE GIVEN CONTRACT AND SUPPLEMENTED BY EXTERNAL KNOWLEDGE SOURCES.

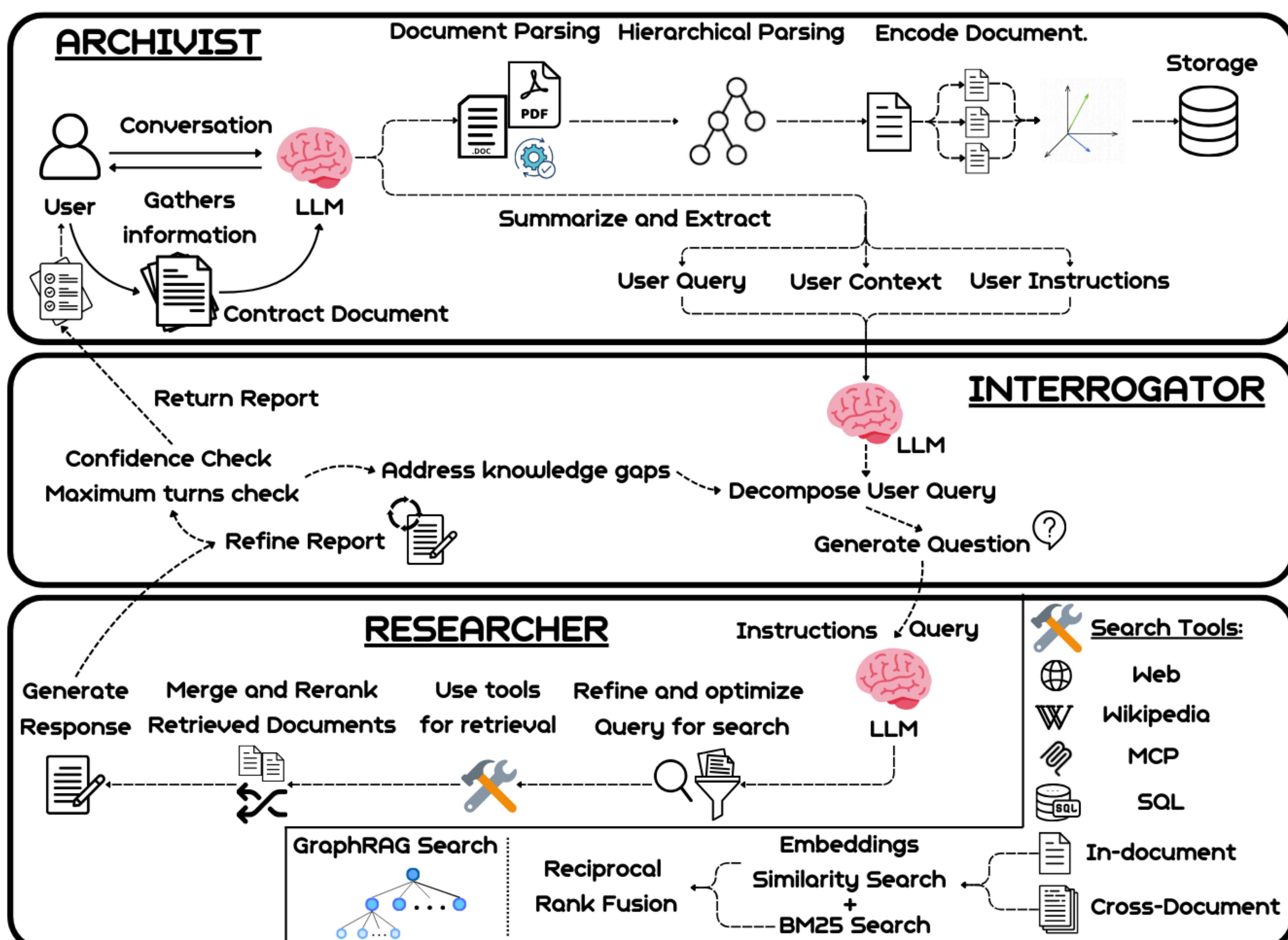
GOALS: ENSURE TRANSPARENT AND TRACEABLE REASONING, WHILE EXPLICITLY REFERENCING THE RELEVANT EVIDENCE SPANS FROM THE CONTRACT.

So, why PAKTON?

- 💪 SUPERIOR GENERATION PERFORMANCE MEASURED, COMPARED TO BASELINES
- 💪 STATE-OF-THE-ART RETRIEVAL PERFORMANCE ACHIEVED
- 💪 PAKTON WAS PREFERRED OVER CHATGPT BY HUMANS FOR CONTRACT ANALYSIS
- 💪 LLM EVALUATION SHOWED PREFERENCE FOR PAKTON VALIDATING HUMANS

- ☑ BETTER EXPLAINABILITY THAN BLACK BOX MODELS
- ☑ OPEN-SOURCE PIPELINE
- ☑ PAKTON BRIDGES THE PERFORMANCE GAPS BETWEEN SMALLER AND LARGER MODELS
- ☑ EASY ON-PREMISE DEPLOYMENT ALLEVIATING PRIVACY CONCERN FROM SENDING SENSITIVE CONTRACTUAL INFORMATION TO PROPRIETARY SYSTEMS

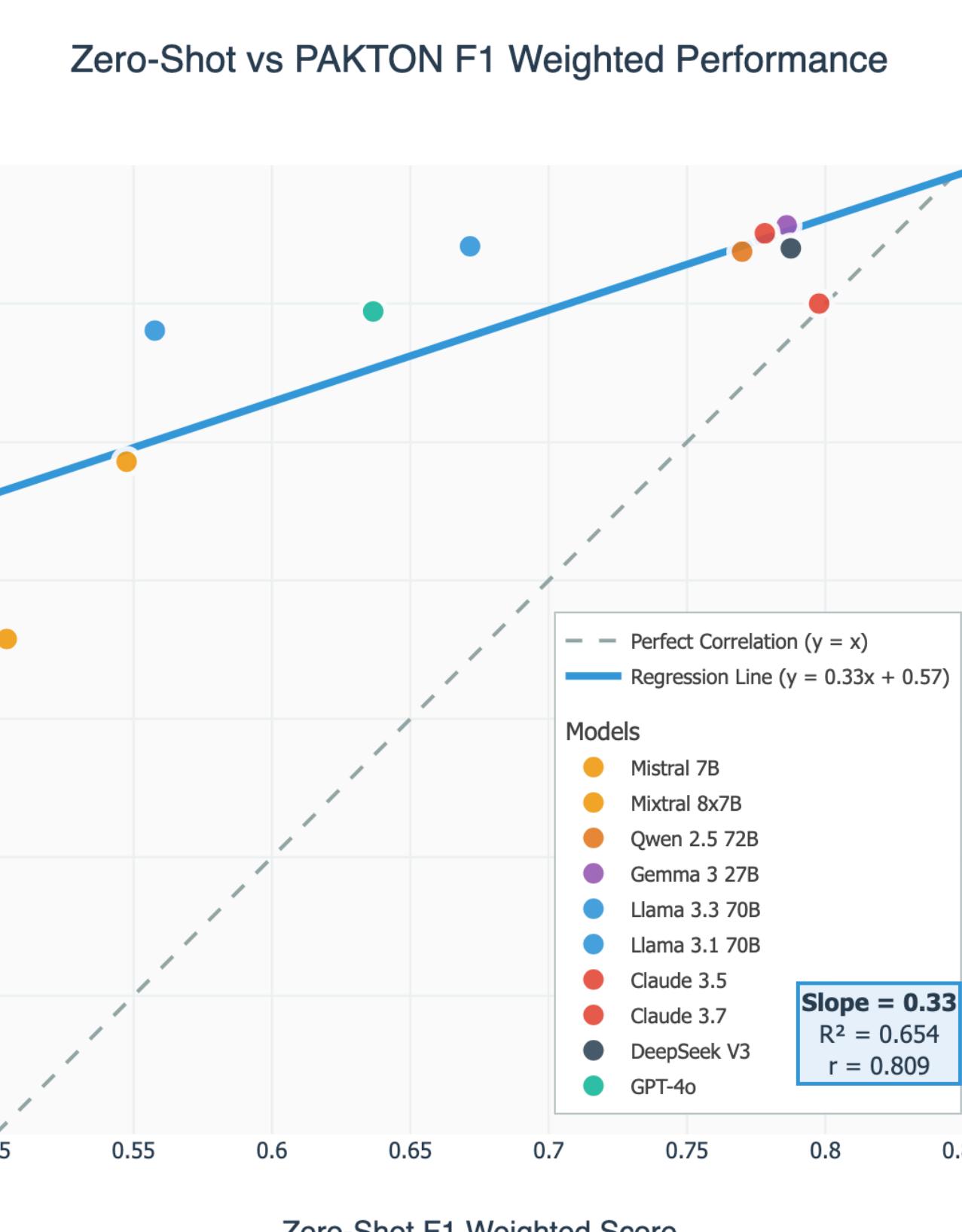
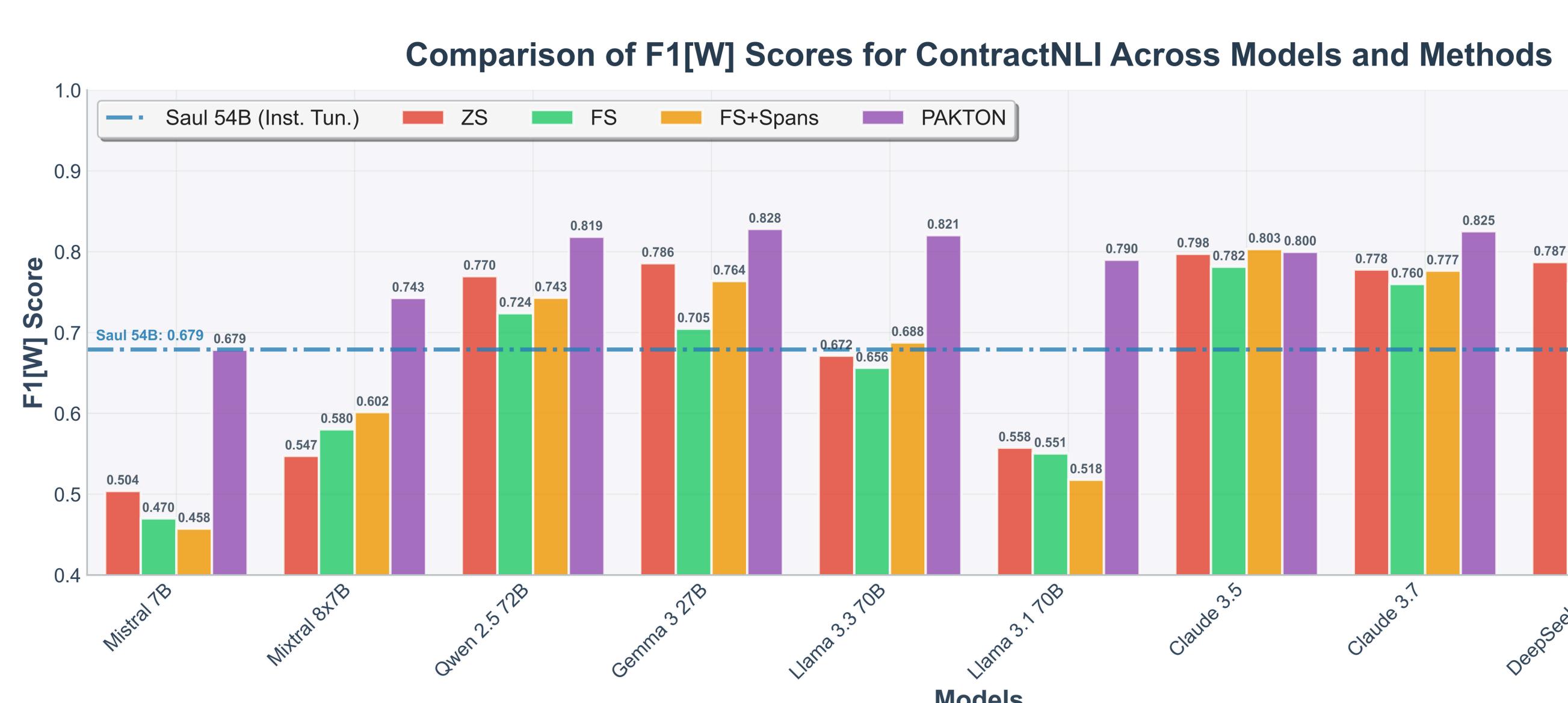
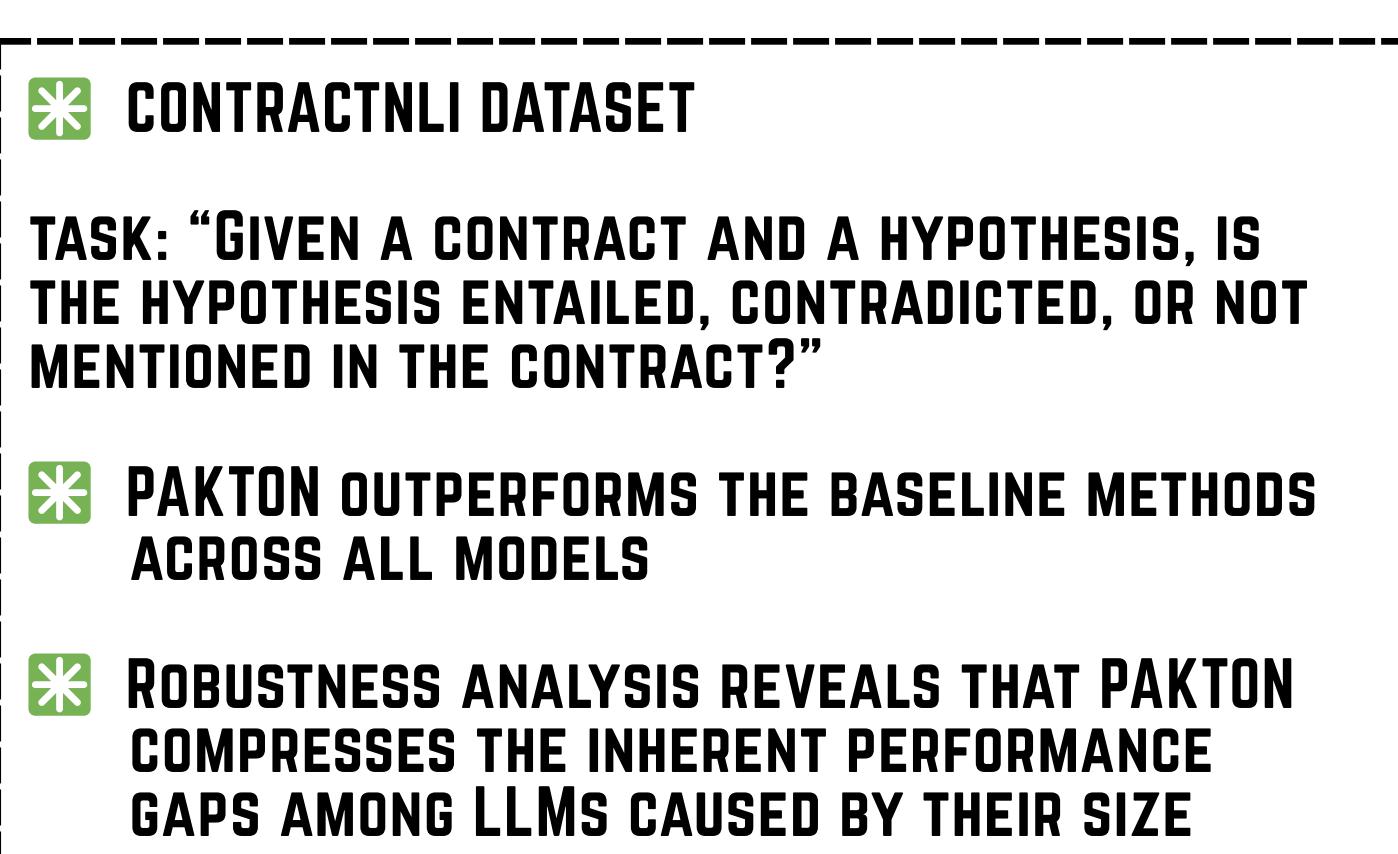
PAKTON FRAMEWORK



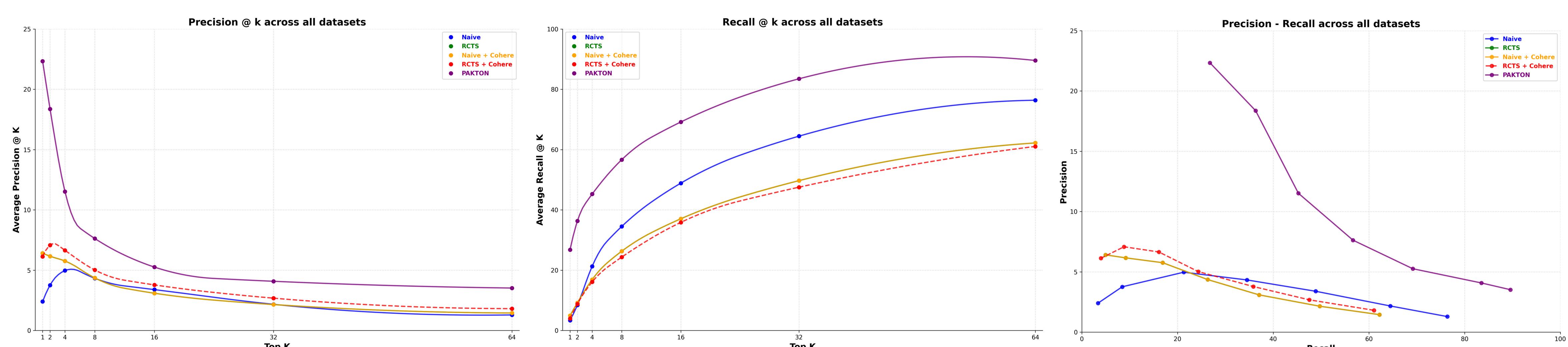
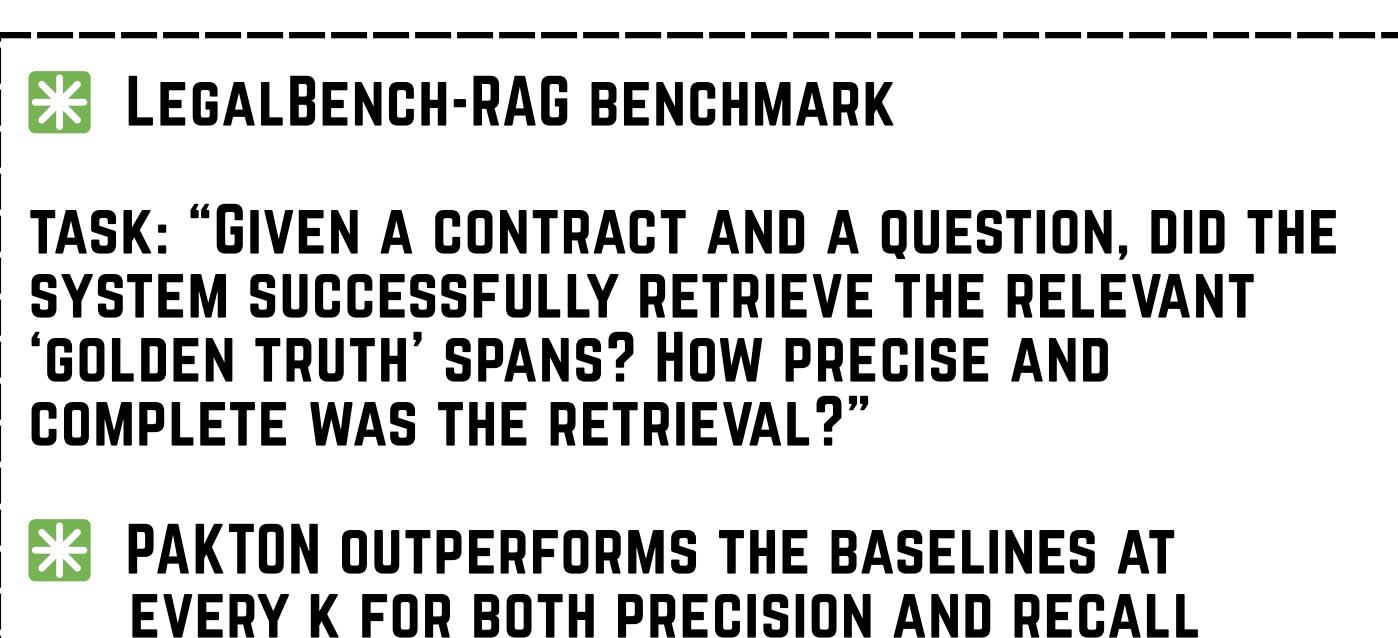
EXPERIMENTS AND RESULTS



END-TO-END GENERATION PERFORMANCE

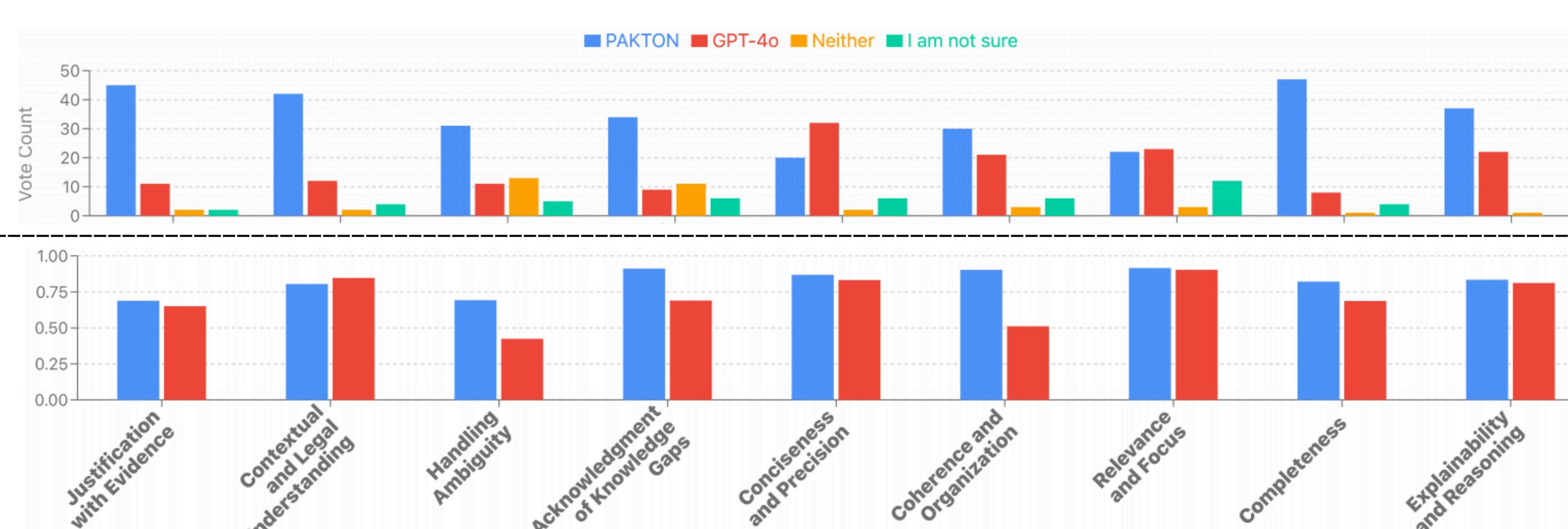


RETRIEVAL (RAG) PERFORMANCE

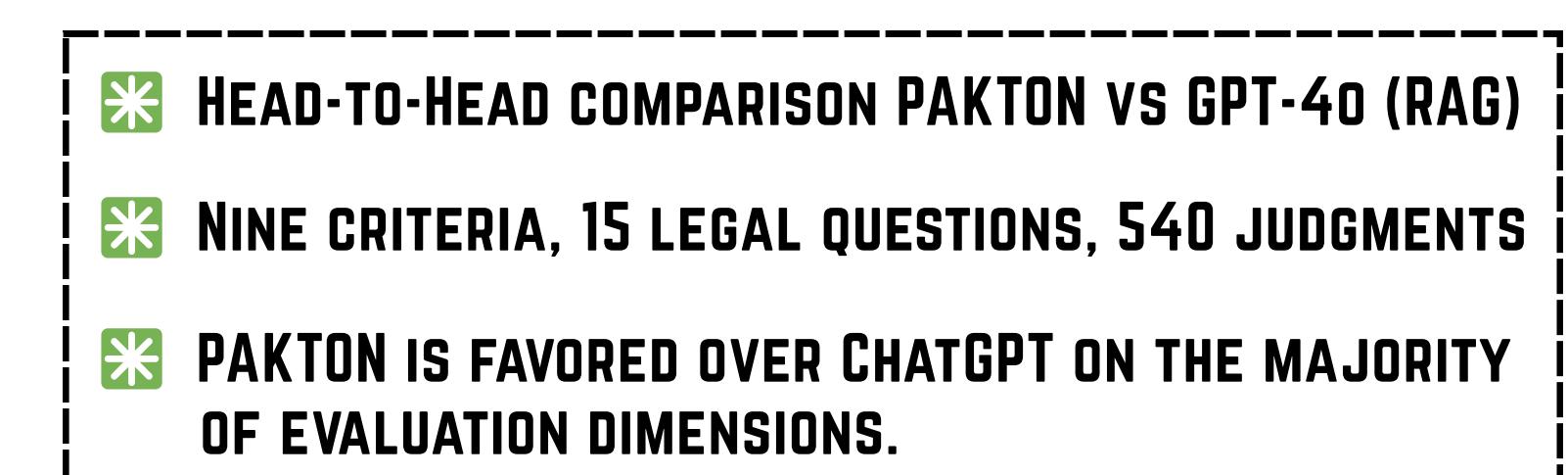


QUALITATIVE RESULTS

HUMAN EVALUATION



LLM-AS-A-JUDGE



- ✳ G-EVAL SCORES FOR PAKTON AND GPT-4O (RAG)
- ✳ SAME NINE CRITERIA, 102 CONTRACT SAMPLES

