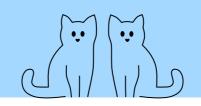
Document Retrieval with Fine-grained Relevance Cues

Ing. Antonín Jarolím Ing. Martin Fajčík Ph.D.





1. Let Human **Annotate Small** Dataset

2. Evaluate LLM on Fine-Grained **Extraction Task**

Ever disappointed by Google's highlighting?



Fine-Grained Cues Motivation

- highlight => get information faster
- lowering halucination in RAG
- token-cues without calling LLM

treatment of varicose veins in legs ca yen ne pepper . ca yen ne pepper is considered a miracle treatment ico se veins . being a very rich source of vitamin c and bio fl av it increases blood circulation and ease s the pain of cong swollen veins . add one tea sp oon of ca yen ne pepper powder to a cup of hot water and stir it well .

what is priority pass

provides you with access to their network of over 700 lounge s these lounge's are managed by a variety of airlines and companies, meaning that significant variation among them a but great international coverage

with thresholding

what is priority pass

membership provides you with access to their network of over 700 lounges these lounge's are managed by a variety of airlines and companies, meaning that significant variation among them a but great international coverage

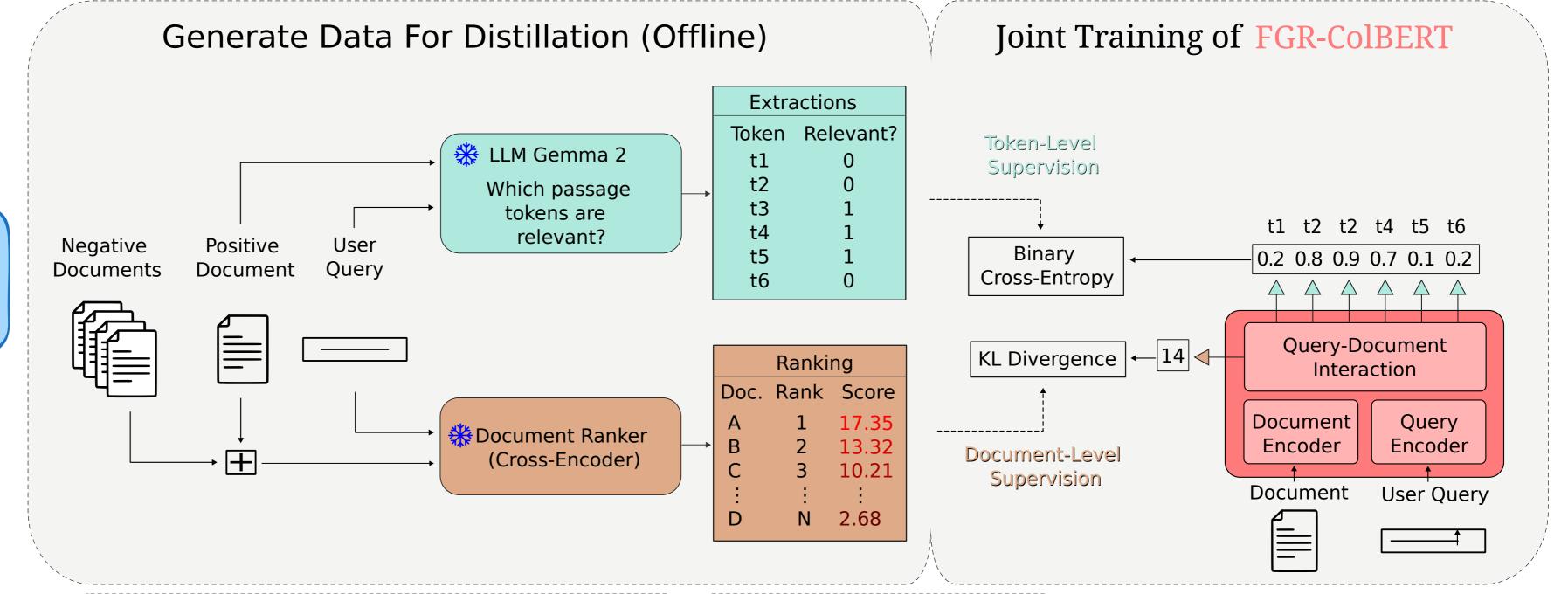
Model Precision Recall F1 Score human annotations (reference) 1.00001.0000select-all (baseline) 0.32491.0000 0.4905tf-idf (baseline) 0.52640.17160.2588gemma2:27b-instruct-fp16 0.73790.76350.7139gemma2:27b-instruct-q8 0.75920.70930.7334gemma2:9b-instruct-fp16 0.67170.50430.4037gemma2:9b-instruct-q8 0.68830.36540.47740.6854gpt-4o-2024-08-06 0.76670.61970.68230.6959gpt-4o-mini-2024-07-18 llama3.1:70b-instruct-q8 0.74950.65870.7012llama3.1:70b-instruct-q4 0.75870.63710.6926llama3.1:8b-instruct-fp16 0.59770.60420.6009

Table 1: Performanc of LLMs on fine-grained extraction task.

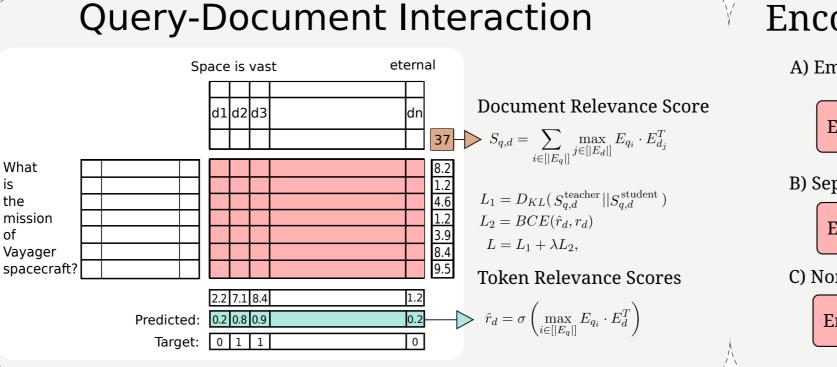
Dataset	1	2	3	4	5	6	7	>8
Train	497,922	250,182	33,030	7,561	2,238	816	335	364
Dorr	4.479	2 101	267	61	25	5	1	6

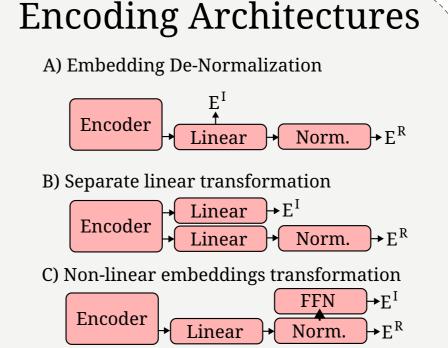
Table 2: Number of extracted spans in a train and dev. sets (24 is max).

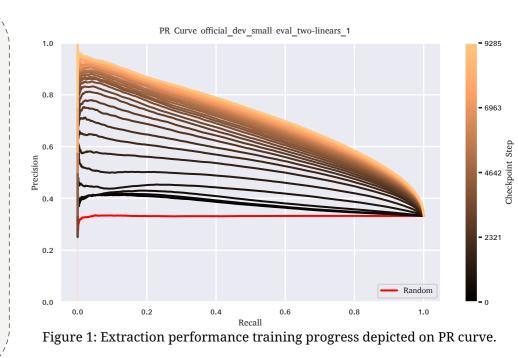
3. Use Winner LLM To Create Large-Scale Training Dataset



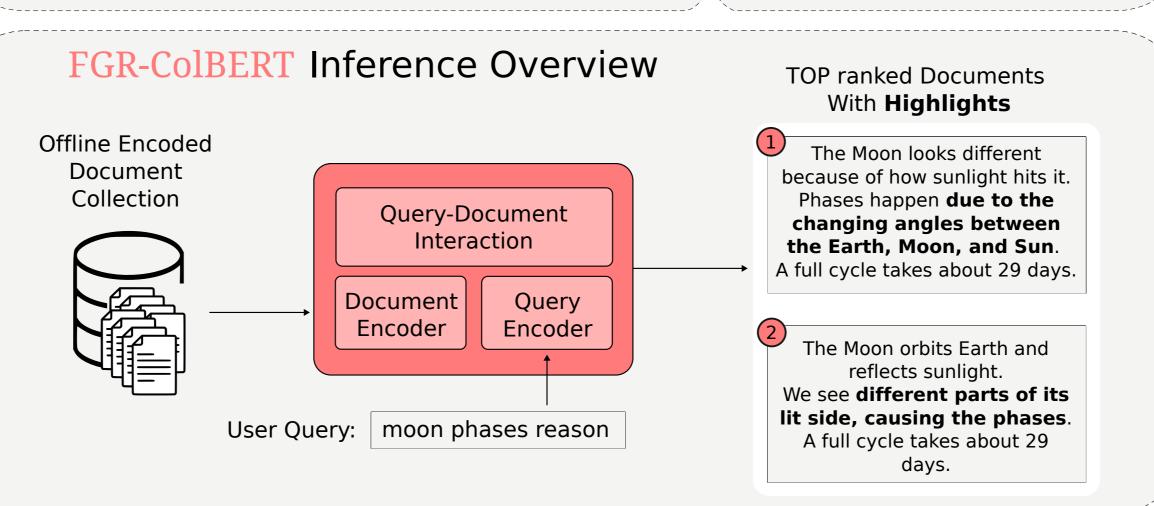
4. Modify Retrieval Model To Enable **Token Extraction**

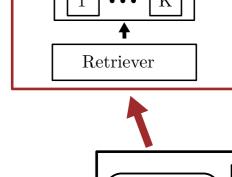






5. Get LLM Quality Token Relevance Cues **During Retrieval**





Gemma size = 27BColbert size = 120M

Holy Grail:

Language Lang. Vision Speech Model Model Model Vision

Output: Ranked list with

Fine Grained Relevancy

Model Database User Query Input **Quick Facts and Notes**

Interpretable

Retrieval

Algorithm

Model

Speech

Model

Retrieval

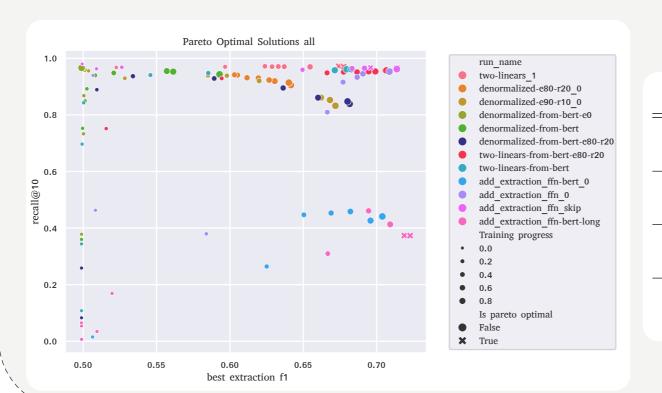
Conclusions

modified retrieval approach providing built-in relevance explainablity

cues obtained from our retrieval COLBERT (120 M) model even matched Gemma-2 (27 B)

three different approaches offering deployment flexibility

Training & Evaluation



Model Architecture	Initialization	F1-score				
Non-linear embeddings transformation	BERT ColBERT	0.7038 0.7093				
Separate linear transformation	BERT ColBERT	0.7067 0.6775				
Embedding de-normalization	BERT ColBERT	0.6816 0.6719				
Retrieval Only	BERT	0.5008				
Table 3: Human Match vs Doc Retrieval Performance						

Generation Errors: 792k / 808K total - (470K MS-MARCO + 317K Top-1 Retrieved) 98% generated correctly

Span Not Found (8400) + Nothing Selected (7800)

Huristacally Fixed 4600 Samples

Another Dataset of 160 unique samples 3 annotators inter agreement fleiss kappa 0.445

Previous evaluation on MD2D (Grad-SAM, AttCAT, Attention-Rollout)