

# Multi-stage Large Language Model Pipelines Can **Outperform GPT-40 in Relevance Assessment**



Julian A. Schnabel<sup>1</sup>, Johanne R. Trippas<sup>2</sup>, Falk Scholer<sup>2</sup> and Danula Hettiachchi<sup>2</sup>

1. Heinrich Heine University, Düsseldorf, Germany 2. RMIT University, Melbourne, Australia Correspondence: julian.schnabel@hhu.de



#### Abstract

Generating relevance labels for documents to indicate the usefulness for specific queries and users. Obtaining these labels from real users is costly and scaling is challenging. It has been shown [Thomas et al., 2024] that Large Language Models can be used to predict these relevance labels. We propose a novel multi-LLM pipeline that divides the relevance assessment task into multiple stages, each utilising different prompts and models of varying sizes and capabilities. This approach beats the baseline performance of GPT-40 on TREC-DL23.

Mo	odel	Pro	mpt	Binary	4-sc	cale	Cost
1	2	1	2	$\kappa$	$\kappa$	lpha	USD
40	-	Normal	-	0.453	0.296	0.408	5.00
mini	-	Normal	-	0.400	0.254	0.359	0.15
mini	mini	Binary	Relevant	0.437	0.284	0.422	0.21
4o	40	Binary	Relevant	0.428	0.280	0.450	6.57
mini	40	Binary	Relevant	0.437	0.286	0.432	2.05
4o	mini	Binary	Relevant	0.428	0.279	0.443	5.05
mini	mini	Binary	Normal	0.439	0.281	0.425	0.21
40	4o	Binary	Normal	0.429	0.280	0.452	6.57
mini	4o	Binary	Normal	0.450	0.295	0.446	2.05
4o	mini	Binary	Normal	0.430	0.276	0.445	5.05
mini	4o	Normal	Normal	0.400	0.260	0.367	2.87
4o	mini	Normal	Normal	0.462	0.294	0.411	5.05





**Figure 1:** Visual overview of the pipeline approach where different models can judge at different stages of the relevance judgement labelling.

We use the TREC 2023 Deep Learning Track (TREC-DL 23). The original judgments were made by NIST assessors, who, given a query, assigned relevance scores to passages based on the following scale:

Label	Count
Label 0 (Irrelevant)	13,866
Label 1 (Related)	4,372
Label 2 (Highly relevant)	2,259
Label 3 (Perfectly relevant)	1,830

Table 2: Accuracy for different GPT model/prompt combinations on TREC-DL23. Cost in USD per million input tokens.

- Four-scale Krippendorff's  $\alpha$  is higher in all multi-model approaches compared to single-model single-stage with GPT-40.
- The four-scale  $\kappa$ -score of GPT-40 remains the highest overall (see Table 2).
- GPT-40 mini/GPT-40 with binary/normal prompts achieves similar accuracy at significantly lower cost (see Figure 3).



Table 1: Label distribution across the dataset.

- 0. **Irrelevant**: The passage has nothing to do with the query.
- 1. **Related**: The passage seems related to the query but does not answer it.
- 2. **Highly relevant**: The passage has some answer for the query, but the answer may be a bit unclear, or hidden amongst extraneous information.
- 3. **Perfectly relevant**: The passage is dedicated to the query and contains the exact answer.

# Results

#### **Reproducing Existing Baselines**

- Reproduction of UMBRELA Baseline
- Similar accuracy to Upadhyay et al. [2024]
- -Comparable misjudgement pattern
- GPT-40 slightly over-optimistic
- Model Comparison
- -GPT-40: Best performance among tested stand-alone models
- Confirmed benchmarks by Alaofi et al. [2024]
- GPT-40 Mini
- -Only 3% of GPT-4o's cost
- High agreement with GPT-40
- Higher binary accuracy with custom prompt (see Tab. 2)

GPT-40 TREC-DL 23

GPT-40 mini TREC-DL 23

## Summary

#### **Pipeline Comparison & Results**

- Accuracy Increase: All pipelines, except the Multi-model Single-stage approach, improved Krippendorff's  $\alpha$  compared to the baseline.
- GPT-40 mini: Achieved the largest accuracy increase with a significant cost reduction.
- Single-model Multi-stage Approach: High accuracy at low cost, outperforming GPT-40 (over 20x more expensive).
- Cost Efficiency: GPT-40 mini enables more complex pipelines, e.g., multistage relevance classification or specialised prompts.

## **Key Insights**

- Specialised Binary Classification: Improved accuracy, especially in multistage approaches with binary relevance decision and classification.
- Limitations:
  - -No duplicate filtering (due to complexity) as noted in Upadhyay et al. [2024].
- -Could expand tests to other TREC-DL datasets and broader tasks.

#### **Future Directions**

0	- 9787	3224	689	162 -	- 8000	- 8606	3892	1026	333 -	- 8000 -
label	- 1467	2025	708	172 -	- 6000	- 1269	1761	981	359 -	- 6000 - -
True ] 2	- 330	1023	686	220 -	- 4000 L ~	- 227	828	854	350 -	- 4000 - -
3	- 156	567	674	433 -	- 2000 - °	- 84	514	673	558 -	- 2000 - -
	0	<sup>1</sup> Predicte	<sup>2</sup> ed label	3		0	<sup>1</sup> Predicte	<sup>2</sup> ed label	3	

Figure 2: Reproduction of baseline (UMBRELA) with GPT-40 and GPT-40 mini.

#### **Relevance Judgement Pipeline Approaches**

Table 2 summarises the evaluation outcomes of the baselines and all proposed pipelines. We evaluate two models (GPT-40 and 40-mini) under both homogeneous (40–40, mini–mini) and heterogeneous (mini–40, 40–mini) pairings, alongside two prompt types: Binary–Relevant and Binary–Normal.

- Prompt Optimization: Techniques like chain-of-thought Wei et al. [2022] and narratives Sadiri Javadi et al. [2024] could further improve performance.
- Spam Filtering: With 75% zero-weighted data in TREC-DL, small, affordable spam-filters could reduce assessment costs significantly.

# Acknowledgments

This research is partially supported by the Australian Research Council (CE200100005) and RMIT AWS Cloud Supercomputing Hub.

#### References

Paul Thomas, Seth Spielman, Nick Craswell, and Bhaskar Mitra. Large language models can accurately predict searcher preferences. 2024. ISBN 9798400704314.

Shivani Upadhyay, Ronak Pradeep, Nandan Thakur, Nick Craswell, and Jimmy Lin. Umbrela: Umbrela is the (open-source reproduction of the) bing relevance assessor. *arXiv preprint arXiv:2406.06519*, 2024.

Marwah Alaofi, Paul Thomas, Falk Scholer, and Mark Sanderson. Llms can be fooled into labelling a document as relevant. 2024.

- Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Fei Xia, Ed Chi, Quoc V Le, Denny Zhou, et al. Chain-of-thought prompting elicits reasoning in large language models. Advances in neural information processing systems, 35:24824–24837, 2022.
- Vahid Sadiri Javadi, Johanne R Trippas, and Lucie Flek. Unveiling information through narrative in conversational information seeking. In Proc. CUI, 2024.