



**UNIPAVIA**  
DIPARTIMENTO  
DI INGEGNERIA  
DELL'INFORMAZIONE



# **SEUPD@CLEF: TEAM CLOSE** Temporal Persistence of IR Systems' Performance

Gianluca Antolini

Marco Martinelli

Nicola Boscolo Cegion

Seyedreza Safavi

Nicola Ferro

Mirco Cazzaro

Farzad Shami

CLEF 2023 – 14<sup>th</sup> Edition  
Thessaloniki – Greece  
19<sup>th</sup> September 2023

# TABLE OF CONTENTS

**01**

## INTRODUCTION

Problem, task & goal

**03**

## RESULTS

Evaluation measures and results

**02**

## OUR SYSTEM

Detailed information about methodologies used in our system

**04**

## CONCLUSIONS

Conclusions and future work

# 01 INTRODUCTION



## PROBLEM

Performance of IR systems can deteriorate over time because web contents and user search preferences change.



## TASK

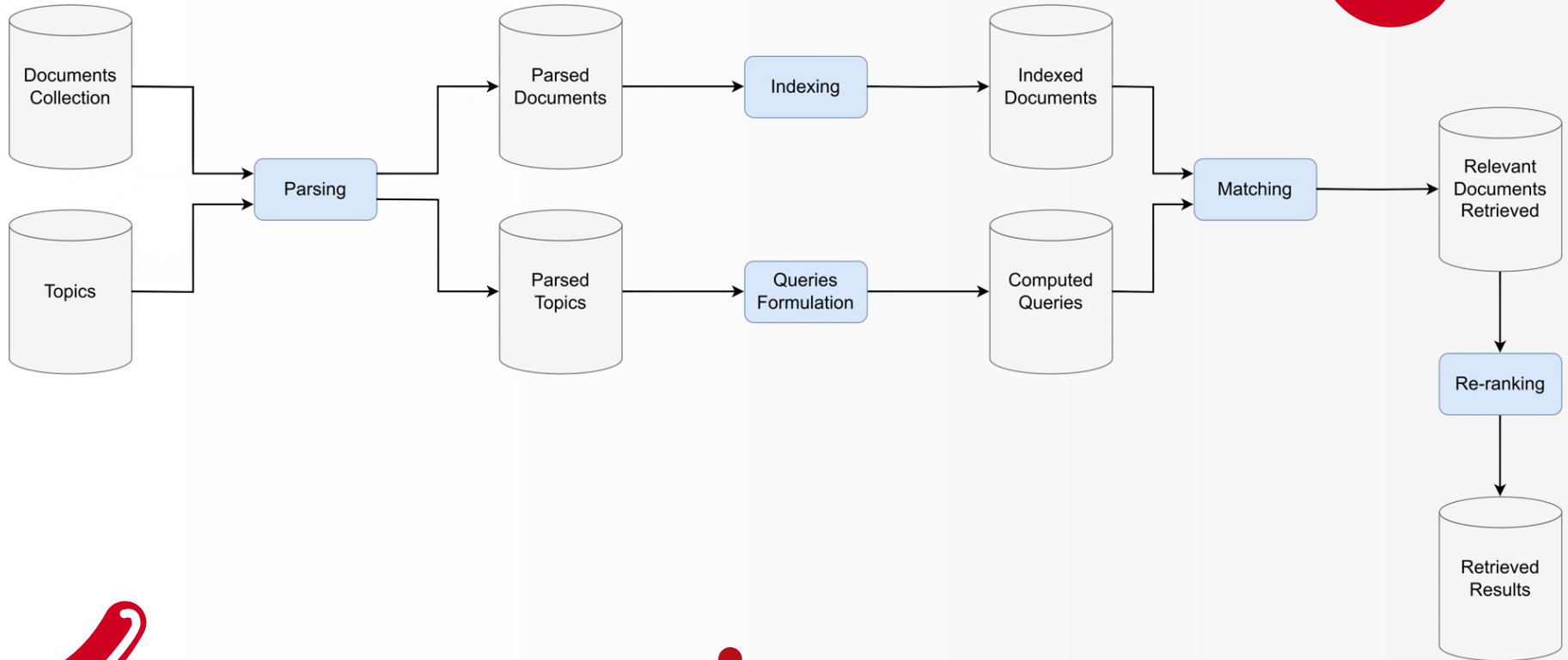
Develop an efficient IR system and training it using *Qwant* training data (French & English):  
user searches  
web documents



## GOAL

Test and measure the system's performance, ensuring that it doesn't decline with data coming from different/distant timestamps

## 02 OUR SYSTEM





# PARSER

Trial and error workflow/methodology (considering improvements in **MAP**):

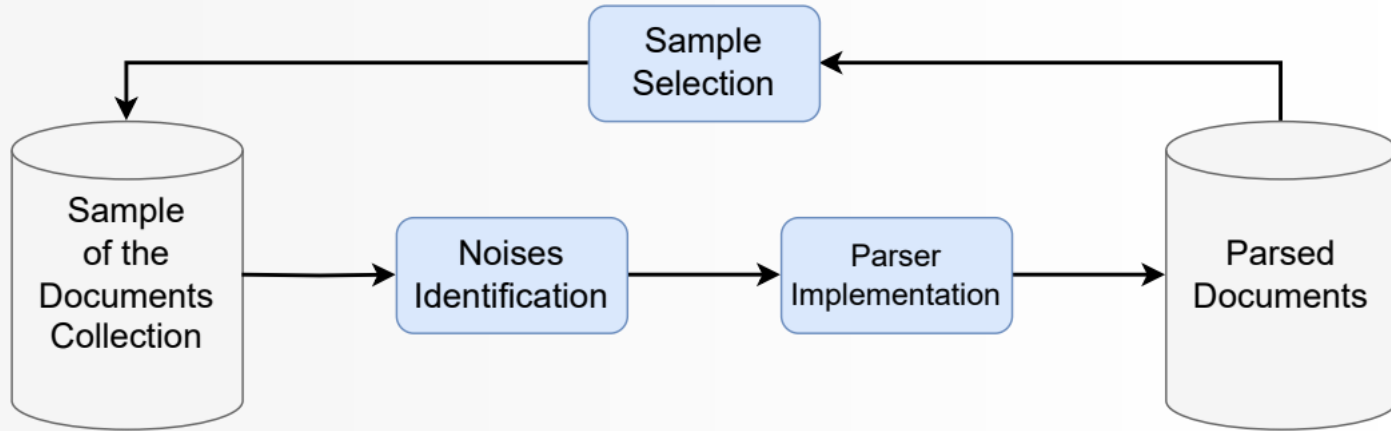
1. Examination of sizeable sample of collection documents in the collection to decide types of noises to be removed.
2. Implementation and run of the parser.
3. Results stored and sample of parsed documents analyzed to restart the procedure.

Final Parsed Document structure:

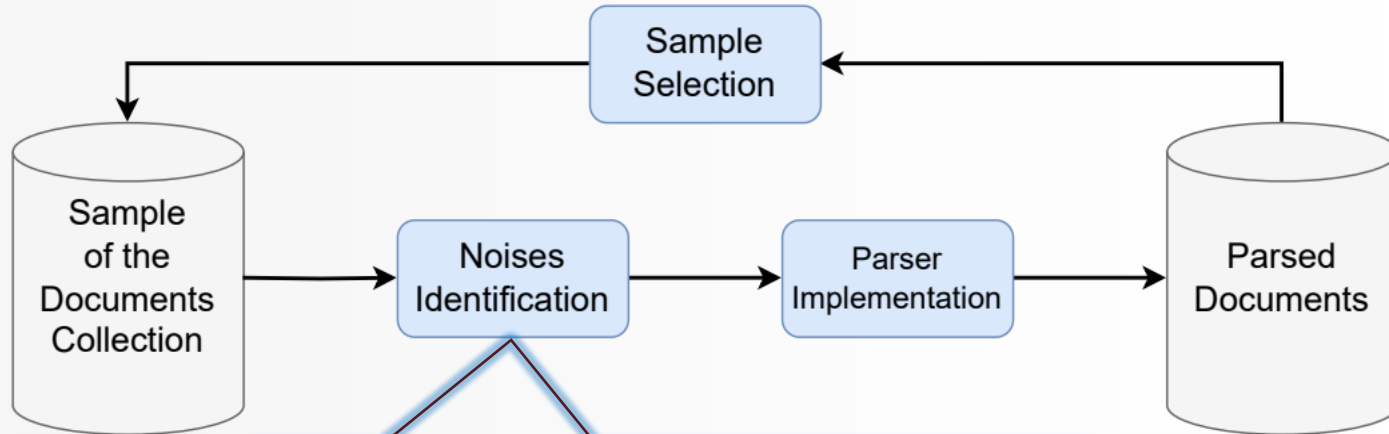
- *id* : document identifier
- *body* : parsed content of document



# PARSER



# PARSER



- JavaScript code
  - HTTP and HTTPS URIs
  - Word patterns  
(e.g. word1\_word2 or word1.word2 or word1:word2)
- + ~0.06 of MAP

- HTML tags and CSS stylesheets
- XML and JSON codes
- Meta tags and document properties
- Navigation menus
- Advertisements
- Footers
- Social media handlers
- Hashtags and mentions





# ANALYZER

Fully customizable class to analyze documents using different approaches.

Based on experiment results with different parameters and by trying both English and French dataset, the following are the parameters used to get the best results:

- French dataset
- FrenchLightStemFilter
- StandardTokenizer
- LowerCaseFilter
- Minimum token length: 2
- Maximum token length: 15
- French word stoplist (662 French words): built upon popular French stoplist and most frequent stopwords in the collection.





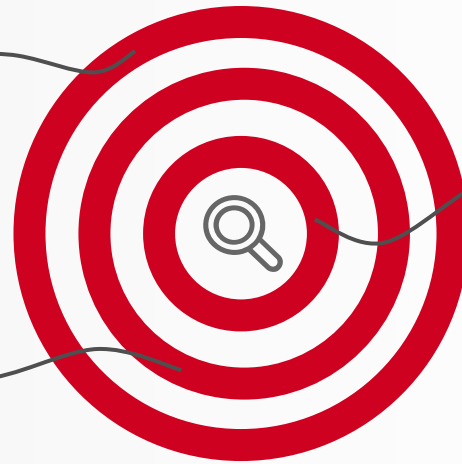
# SEARCHER

Retrieve relevant information by analyzing user queries and searching through indexed documents, returning a ranked list of matching documents.

**Query Expansion**

**Query Boosting**

**Document  
Re-Ranking**





# QUERY EXPANSION

During the search function, Query Expansion performed by generating new queries from the original ones.

Python script that utilizes OpenAI's Text Completion Endpoints to generate expanded terms for each query.



We used the DaVinci model with a *temperature parameter*  $T$  of 0.6 for optimal results.

$$p(x_i) = \frac{e^{x_i}}{\sum_{j=1}^V e^{x_j}}$$

PROBABILITY DISTRIBUTION USING SOFTMAX  
FUNCTION WITHOUT TEMPERATURE PARAMETER

$$p(x_i) = \frac{e^{\frac{x_i}{T}}}{\sum_{j=1}^V e^{\frac{x_j}{T}}}$$

PROBABILITY DISTRIBUTION USING SOFTMAX  
FUNCTION WITH A TEMPERATURE PARAMETER  $T$





# QUERY BOOSTING

To assign higher relevance to specific query terms or queries.

Our approach consists of building **BooleanQueries** in the search function in the following way:

- Each query has its expansions added with **SHOULD** clause;
- A main query with the **MUST** clause;
- Main query boosted using Lucene's BoostQuery with:

$$14,68 * |lq|$$

where  $lq$  is the list of the expansions of the query  $q$

- Boost value fine-tuned to 14.68 through a trial and error for parameter optimization.





# DOCUMENT RE-RANKING



- Rank documents retrieved by the Searcher using **all-MiniLM-L6-v2**, a 384-dimensional Sentence Transformer model.
- Initialize Re-Ranker in the Searcher's constructor and create a predictor for inference.
- During search function, embeddings created for documents using predictor, similarity calculated between query and documents.
- Scores multiplied by document's **BM25Similarity** and **Cosine Similarity**.

$$similarity = \frac{t \cdot d_i}{|t| \cdot |d_i|}$$

**COSINE SIMILARITY FUNCTION**

$$rank = BM25\_score \cdot similarity$$

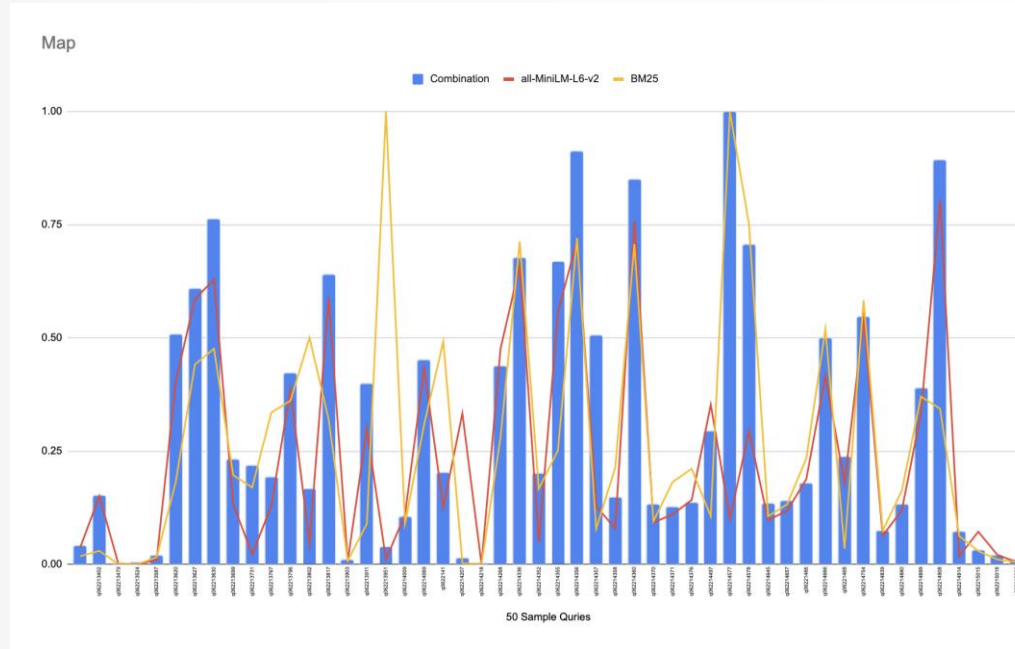
**FINAL RE-RANKING SCORE**

- Tested parameter combinations and determined that multiplying document's score by BM25Similarity with cosine similarity yields the best results.





# DOCUMENT RE-RANKING



RE-RANKING PERFORMED ON A SAMPLE OF 50 QUERIES (MAP SCORES).



# 03 RESULTS

The following slides provide the results of our **best** performing runs, with both **French** and **English** dataset

Parameter	Run 1	Run 2	Run 3	Run 4	Run 5
Token Filter	Porter-StemFilter	FrenchLight-StemFilter	FrenchLight-StemFilter	PorterStem-Filter	FrenchLight-StemFilter
Tokenizer	Standard	Standard	Standard	Standard	Standard
Length Filter	2-15	2-15	2-15	2-15	2-15
Stop Filter	" <i>long-stoplist.txt</i> "	" <i>long-stoplist-fr.txt</i> "	" <i>long-stoplist-fr.txt</i> "	" <i>long-stoplist.txt</i> "	" <i>new-long-stoplist-fr.txt</i> "
Lower Case Filter	Yes	Yes	Yes	Yes	Yes
Similarity	BM25	BM25	BM25	BM25	BM25
Query Expansion	No	Yes	Yes	Yes	Yes
Re-ranking	No	No	Yes	Yes	Yes

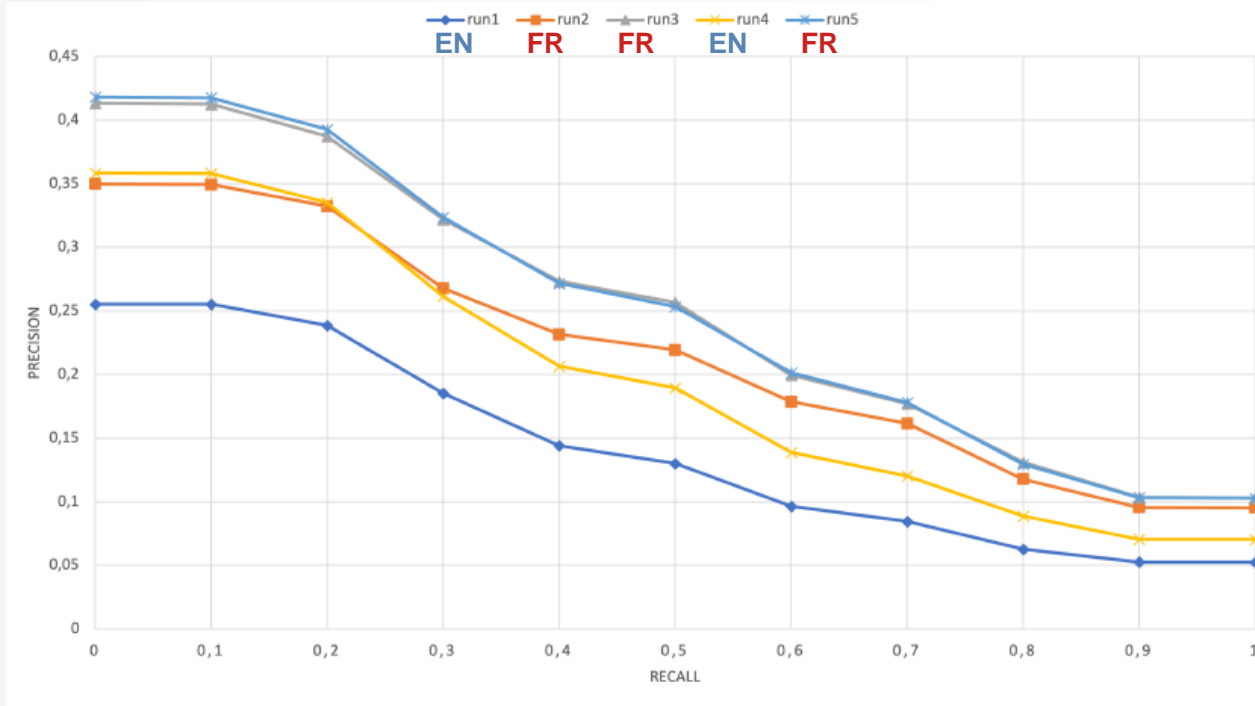


# RESULTS

metrics	run1	run2	run3	run4	run5
num_q	657	669	669	667	667
num_ret	646525	658446	658347	652222	657903
num_rel	2550	2611	2611	2603	2600
num_rel_ret	1772	2182	2191	1866	2232
map	0.1307	0.2022	0.2335	0.1856	0.2351
gm_map	0.0117	0.046	0.061	0.0239	0.0629
Rprec	0.1041	0.1697	0.1989	0.1654	0.2022
bpref	0.3142	0.3734	0.3869	0.3466	0.3861
recip_rank	0.2436	0.3287	0.3891	0.3441	0.3945
iprec_at_recall_0.00	0.2553	0.3499	0.4134	0.3584	0.4182
iprec_at_recall_0.20	0.2387	0.3324	0.3873	0.3353	0.3927
iprec_at_recall_0.40	0.1441	0.2316	0.2732	0.2066	0.2716
iprec_at_recall_0.60	0.0965	0.1786	0.1996	0.1388	0.2014
iprec_at_recall_0.80	0.0628	0.1178	0.1311	0.0887	0.1295
iprec_at_recall_1.00	0.0525	0.0954	0.1031	0.0704	0.1028
P_10	0.0848	0.1296	0.1435	0.1126	0.1432
P_100	0.0186	0.0256	0.0268	0.0222	0.0268
P_1000	0.0027	0.0033	0.0033	0.0028	0.0033
recall_10	0.2166	0.3352	0.367	0.2849	0.3621
recall_100	0.4718	0.6426	0.6714	0.5536	0.6723
recall_1000	0.6816	0.8192	0.8218	0.7004	0.8392
infAP	0.1307	0.2022	0.2335	0.1856	0.2351
gm_bpref	0.0152	0.0387	0.0405	0.022	0.038
utility	-978.6621	-977.701	-977.5262	-972.2489	-979.6687
ndcg	0.2719	0.3655	0.3924	0.3291	0.3982
ndcg_rel	0.236	0.3119	0.3416	0.2939	0.3471
Rndcg	0.1708	0.2387	0.2657	0.2271	0.2714
ndcg_cut_5	0.1285	0.1908	0.2232	0.1854	0.2269
ndcg_cut_10	0.1609	0.2426	0.2739	0.2227	0.2758
ndcg_cut_100	0.2351	0.3349	0.3652	0.3016	0.3678
ndcg_cut_1000	0.2719	0.3655	0.3924	0.3291	0.3982
map_cut_10	0.1046	0.1665	0.1975	0.1556	0.1993
map_cut_100	0.1284	0.2	0.2315	0.1836	0.2328
map_cut_1000	0.1307	0.2022	0.2335	0.1856	0.2351



# RESULTS





# RESULTS

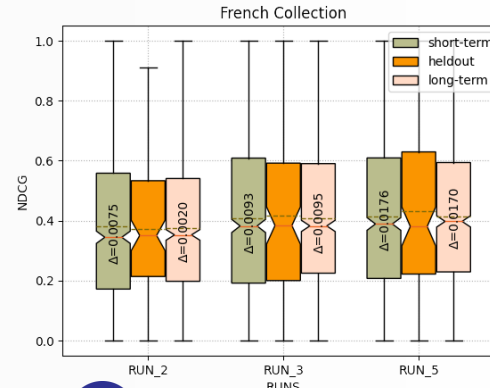
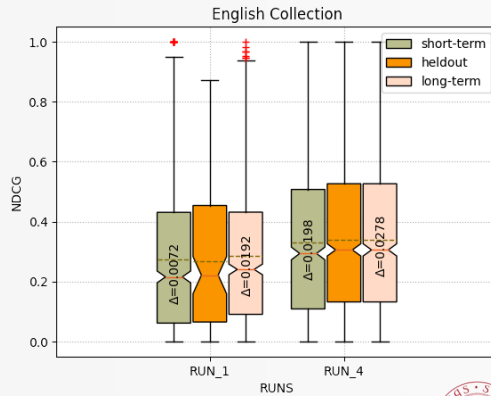
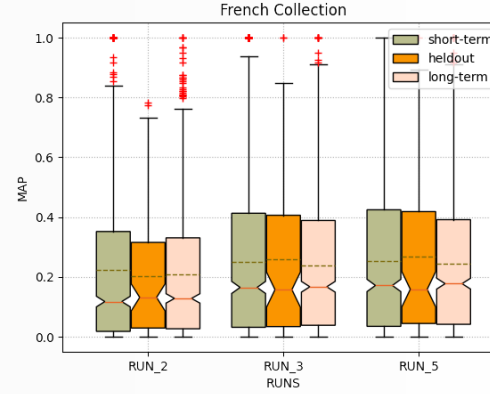
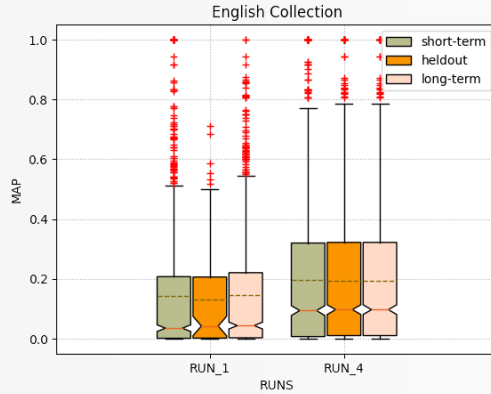
heldout						
run	language	type	map	p@10	NDCG	recall
run2	FR	QUEREXPANSION	0.2029	0.1367	0.3725	0.8312
run3	FR	RERANKING	0.2595	0.1541	0.4166	0.8348
run5	FR	SBERT_BM25	0.2675	0.1561	0.4318	0.8726
run1	EN	JSCLEANER_BM25	0.1299	0.0897	0.2674	0.6381
run4	EN	RERANKING_ENGLISH	0.1822	0.1122	0.3113	0.6279

Short term						
run	language	type	map	p@10	NDCG	recall
run2	FR	QUEREXPANSION	0.2215	0.1326	0.3800	0.8164
run3	FR	RERANKING	0.2511	0.2171	0.4073	0.8142
run5	FR	SBERT_BM25	0.2540	0.1497	0.4142	0.8360
run1	EN	JSCLEANER_BM25	0.1438	0.0902	0.2746	0.6566
run4	EN	RERANKING_ENGLISH	0.1956	0.1145	0.3311	0.6804

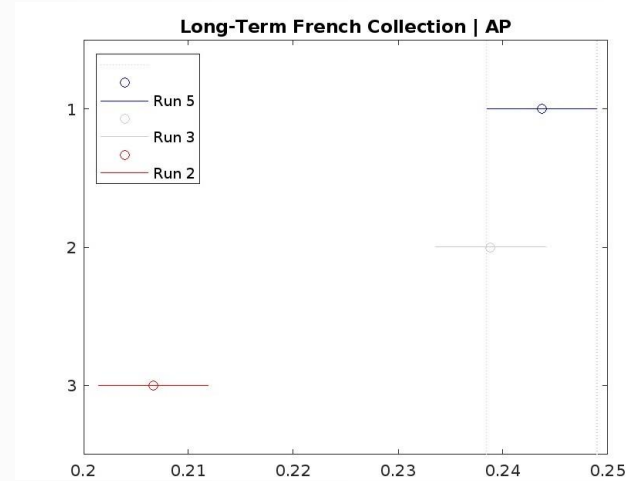
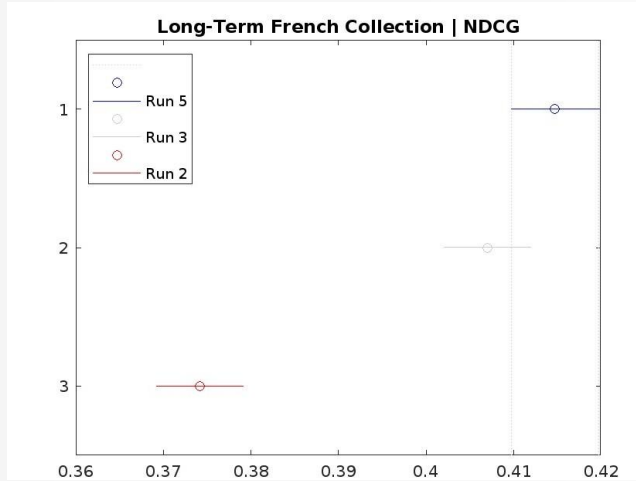
Long term						
run	language	type	map	p@10	NDCG	recall
run2	FR	QUEREXPANSION	0.2067	0.1423	0.3745	0.8312
run3	FR	RERANKING	0.2388	0.1555	0.4071	0.8336
run5	FR	SBERT_BM25	0.2437	0.1594	0.4148	0.8540
run1	EN	JSCLEANER_BM25	0.1450	0.0975	0.2866	0.6906
run4	EN	RERANKING_ENGLISH	0.1930	0.1258	0.3391	0.7119



# Statistical Analysis



# Statistical Analysis | FR



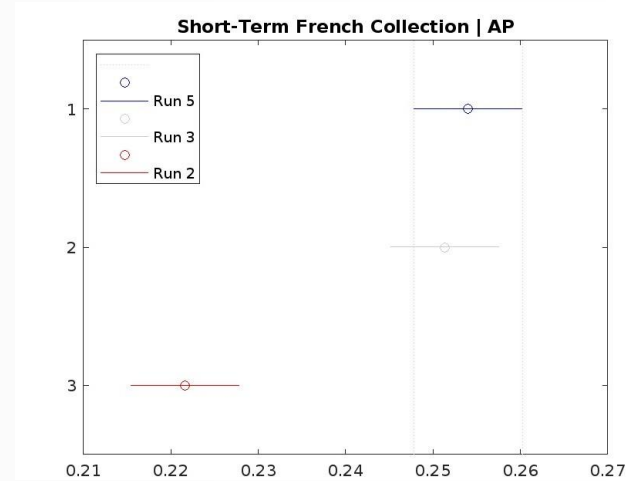
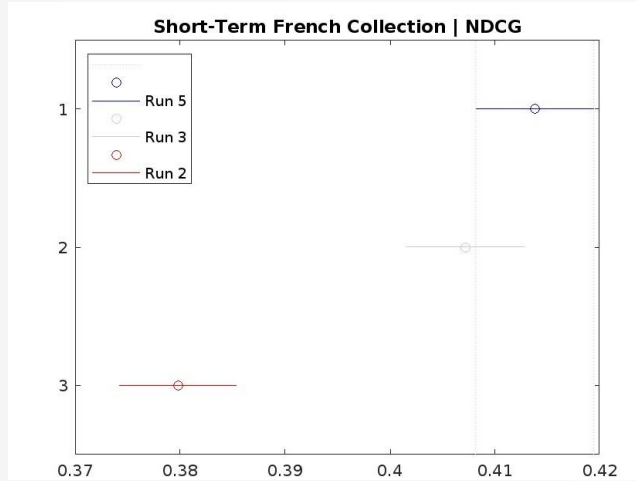
Run A	Run B	Lower Limit	A-B	Upper Limit	P-value
1	2	-0.0023	0.0077	0.0177	0.1661
1	3	0.0305	0.0405	0.0505	$1.16 \cdot 10^{-21}$
2	3	0.0228	0.0328	0.0428	$3.84 \cdot 10^{-14}$

Run A	Run B	Lower Limit	A-B	Upper Limit	P-value
1	2	-0.0057	0.0049	0.0154	0.523
1	3	0.0265	0.037	0.0476	$3.57 \cdot 10^{-16}$
2	3	0.0216	0.0322	0.0427	$2.50 \cdot 10^{-12}$





# Statistical Analysis | FR



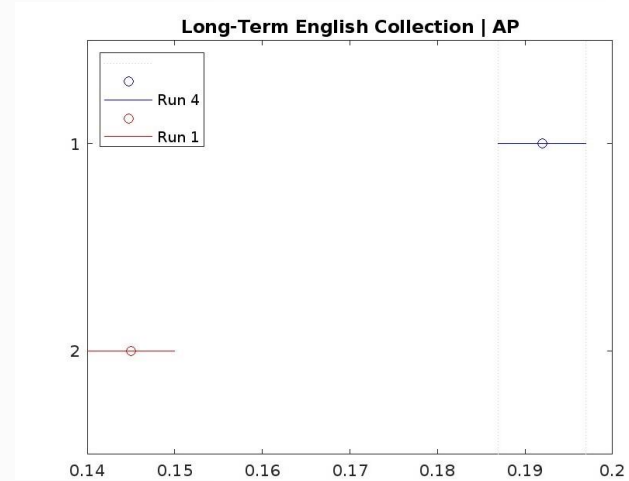
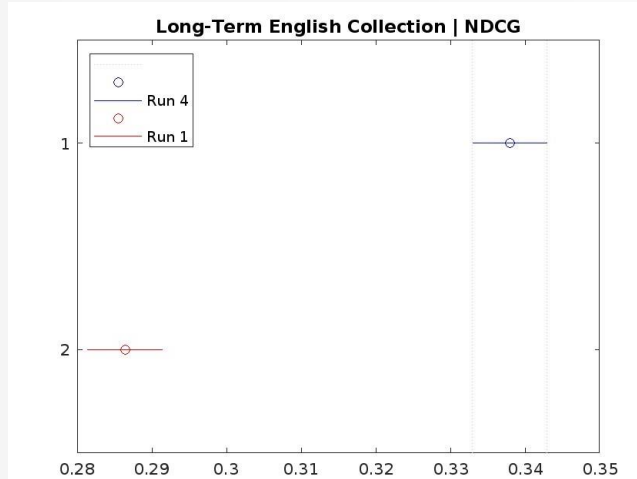
Run A	Run B	Lower Limit	A-B	Upper Limit	P-value
1	2	-0.0042	0.0067	0.0177	0.3215
1	3	0.0232	0.0341	0.0451	$7,97 \cdot 10^{-13}$
2	3	0.0164	0.0274	0.0384	$1.39 \cdot 10^{-8}$

Run A	Run B	Lower Limit	A-B	Upper Limit	P-value
1	2	-0.0095	0.0027	0.0148	0.8631
1	3	0.0203	0.0324	0.0445	$1.18 \times 10^{-9}$
2	3	0.0176	0.0297	0.0419	$2.87 \times 10^{-8}$





# Statistical Analysis | EN

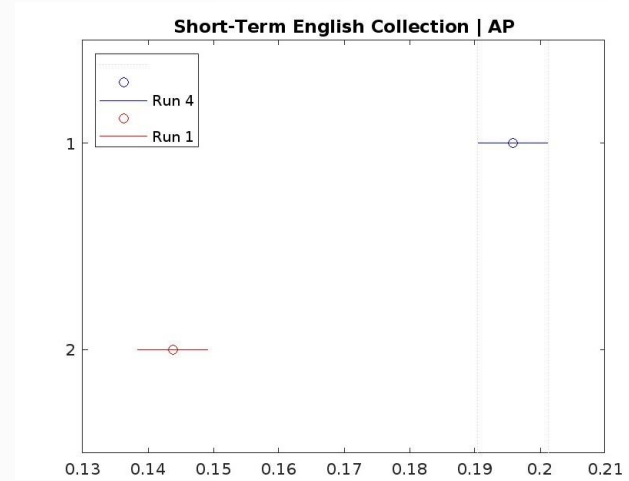
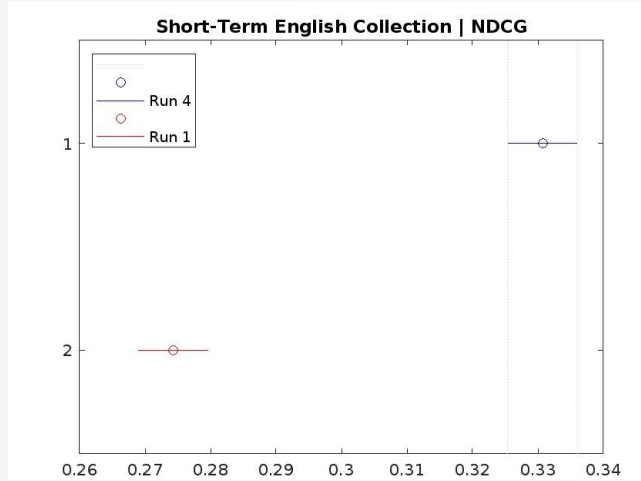


Run A	Run B	Lower Limit	A-B Difference	Upper Limit	P-value
1	2	0.0409	0.0508	0.0608	$1.49 \times 10^{-24}$

Run A	Run B	Lower Limit	A-B Difference	Upper Limit	P-value
1	2	0.0364	0.0464	0.0565	$4.44 \times 10^{-20}$



# Statistical Analysis | EN



Run A	Run B	Lower Limit	A-B Difference	Upper Limit	P-value
1	2	0.0458	0.0564	0.0671	0.00

Run A	Run B	Lower Limit	A-B Difference	Upper Limit	P-value
1	2	0.0412	0.052	0.0628	$1.05 \times 10^{-21}$



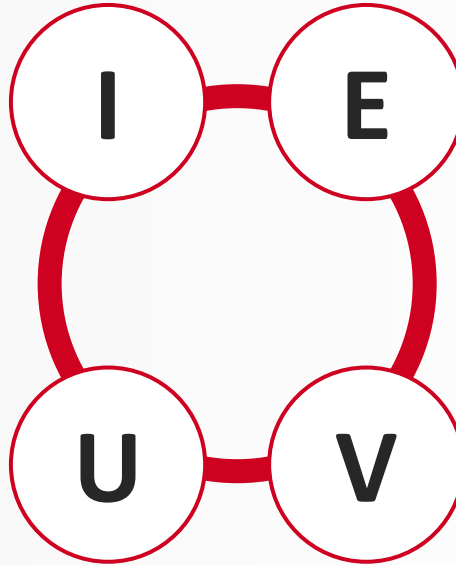
# 04 CONCLUSIONS & FUTURE WORK

## Improving Re-ranking

Diversify scores, explore BERT models, fine-tune SBERT for relevance and document retrieval enhancement.

## Utilizing Document Links

Enhance search results with link details and domain authority from URL keywords.



## Enhancing Query Expansion

Better ChatGPT prompts and explore alternative Large Language Models techniques.

## Vector-based Document Indexing

Store the documents embeddings in the indexing section to speed up the searching process



# THANKS!

Do you have any questions?

