

SEUPD@CLEF: Team FADERIC on A Query Expansion and Reranking Approach for the Longeval Task

Notebook for the LongEval Lab at CLEF 2023

Enrico Bolzonello¹, Christian Marchiori¹, Daniele Moschetta¹, Riccardo Trevisiol¹, Fabio Zanini¹ and Nicola Ferro¹

¹University of Padua, Italy

Abstract

This report explains and analyzes the system developed by Team FADERIC for the LongEval Lab at CLEF 2023, Task 1 - LongEval-Retrieval. The team members are all students following the Search Engines course a.y. 2022/23 at the Computer Engineering master degree at University of Padua. The system developed is a search engine that has to retrieve documents from a corpus, composed of original files in French language and automatically translated files in English language. The produced IR system exploits the query expansion technique, such as word N-grams and synonyms, and also the use of a reranking to improve the overall performance. Evaluating the longitudinal effectiveness of the system using the multiple collections provided by CLEF, we show that the performances remain satisfactory over time.

Keywords

CLEF, LongEval, Information retrieval, Search engines, Query expansion, Reranking

1. Introduction

Search engines have become an indispensable tool for people to retrieve various kinds of information in their daily routine. However, recent research has shown that *Information Retrieval (IR)* systems tend to perform poorly over time as the test data becomes more distant from the training data. This issue is particularly critical in the field of computer science, where data is constantly updated and information quickly becomes obsolete. Therefore, *Conference and Labs of the Evaluation Forum (CLEF) 2023* LongEval task [1] has gained interest in evaluating the temporal persistence of IR systems. The aim of this report is to present the solution of Team FADERIC to this challenge. The team members are all students following the Search Engines course a.y. 2022/23 at the Computer Engineering master degree at University of Padua. The paper is organized as follows: Section 2 shows the related work we have started from; Section 3 describes our approach; Section 4 explains our experimental setup; Section 5 discusses our main findings; finally, Section 6 draws some conclusions and outlooks for future work.

CLEF 2023: Conference and Labs of the Evaluation Forum, September 18–21, 2023, Thessaloniki, Greece

✉ enrico.bolzonello@studenti.unipd.it (E. Bolzonello); christian.marchiori@studenti.unipd.it (C. Marchiori); daniele.moschetta@studenti.unipd.it (D. Moschetta); riccardo.trevisiol.1@studenti.unipd.it (R. Trevisiol); fabio.zanini@studenti.unipd.it (F. Zanini); ferro@dei.unipd.it (N. Ferro)

🌐 <http://www.dei.unipd.it/~ferro/> (N. Ferro)

🆔 0000-0001-9219-6239 (N. Ferro)

© 2023 Copyright for this paper by its authors.
Use permitted under Creative Commons License Attribution 4.0 International (CC BY 4.0).



CEUR Workshop Proceedings (CEUR-WS.org)

2. Related Work

To understand the task and the collections provided by CLEF we have used the paper by LongEval organizers Galuščáková et al. [2]. This has helped us to tackle the problem by figuring out how given documents and queries were collected and which were the main goals of the task. In the paper are also given some *baseline* performances that we have used to benchmark our system during its development.

We decided to exploit query expansion techniques basing our knowledge of the theme on the works from Carpineto and Romano [3] and Azad and Deepak [4]. We have decided to use various techniques, such as word N-grams and synonyms, which we have then tried during the development of our system.

Moreover, we choose to approach the problem by also using *reranking*. To do so we have used the work from Alkhalifa et al. [5], which explains the problem of using a language model to address the longitudinal evaluation task. We build our reranking approach based on the work from Yilmaz et al. [6] who developed the Birch system.

3. Methodology

In this section, we describe the methodology we have adopted in order to develop an IR system for the task.

3.1. Parser

We manually inspected the documents provided by CLEF in order to understand their *structure* and be able to extract the body and ID of each document from them.

In order to do that, we created a tool called *parser* that has been essential for extracting information from documents in the specified format used by *Text REtrieval Conference (TREC)*. The parser helps us create structured objects that are used for analysis and indexing within the IR system.

Here are the key components we implemented:

- **ParsedDocument**: represents a parsed document to be indexed. It has two attributes: ID for the unique identifier of the document and body for the document's content. This class provides functionalities to set and retrieve documents' attributes.
- **DocumentParser**: represents an abstract class providing basic functionalities to iterate over the elements of a ParsedDocument, reading and parsing its content.
- **LongEvalDocumentParser**: specific DocumentParser for the LongEval corpus. It provides an implementation of a parser for the documents in the TREC format. The parser reads a document and returns a ParsedDocument that contains the ID and the body of the document.

We used the LongEvalDocumentParser and the ParsedDocument in the indexer to represent the content of the documents that are in the directory specified by docsDir. The first one has been used to iterate over the content of the specified document, while the second one has been used to represent a document to be indexed.

3.2. Analyzer

In order to *process texts* from documents and queries, we have implemented custom analyzers: since the collections are provided in both French and English language, two of them have been implemented.

3.2.1. French analyzer

The `FrenchLongEvalAnalyzer` component is in charge of processing French language texts, it is composed by:

- **Tokenizing:** the `StandardTokenizer` is used, which exploits the Word Break rules from Unicode Text Annex #29 [7];
- **Character folding:** the `ICUFoldingFilter` is used, which applies the foldings from Unicode Technical Report #30 [8]. This is useful to fold upper cases, accents and other kinds of complex characters;
- **Elision removal:** the `ElisionFilter` is used, which removes the elision from words;
- **Stopword removal:** the `StopFilter` is used, which removes the given stopwords from the tokens. In this system we have tried using the default Lucene¹ stoplist and custom one, generated by picking the 50 most frequent terms in the documents;
- **Position filtering:** a custom `TokenFilter` has been implemented to set the `positionIncrement` attribute of all tokens to a specific value. This will be useful to ignore the `positionIncrement` due to removed stopwords in the search phase when we will use the proximity between tokens, as explained in Section 3.4.3;
- **Stemming:** this process is useful to reduce words to their base form, in this system we have tried using the Snowball French [9] and Light [10] stemmers.

In Table 1 we show an example of the analyzing process for the French language.

3.2.2. English analyzer

The `EnglishLongEvalAnalyzer` component is in charge of processing English language texts, it is composed by:

- **Tokenizing:** the `StandardTokenizer` is used, which exploits the Word Break rules from Unicode Text Annex #29 [7];
- **Character folding:** the `ICUFoldingFilter` is used, which applies the foldings from Unicode Technical Report #30 [8]. This is useful to fold upper cases, accents and other kinds of complex characters;
- **Possessive removal:** the `EnglishPossessiveFilter` is used, which removes the very frequent possessives ('s) from words;
- **Stopword removal:** the `StopFilter` is used, which removes the given stopwords from the tokens. In this system we have tried using the default Lucene stoplist and custom one, generated by picking the 50 most frequent terms in the documents;

¹https://lucene.apache.org/core/9_5_0/index.html

Table 1
French analyzer process

Step	Tokens
	La méthode d'analyse de texte est essentielle pour l'extraction d'informations.
Tokenizing	[La, méthode, d'analyse, de, texte, est, essentielle, pour, l'extraction, d'informations]
Character folding	[la, methode, d'analyse, de, texte, est, essentielle, pour, l'extraction, d'informations]
Stopword removal (50 most freq.)	[methode, analyse, texte, essentielle, extraction, informations]
Stemming (Ligth)	[method, analys, text, esentiel, extraction, inform]

Table 2
English analyzer process

Step	Tokens
	The text analysis method's importance lies in its role in information extraction.
Tokenizing	[the, text, analysis, method, importance, lies, in, its, role, in, information, extraction]
Character folding	[the, text, analysis, method, importance, lies, in, its, role, in, information, extraction]
Stopword removal (50 most freq.)	[text, analysis, method, importance, lies, role, information, extraction]
Stemming (Krovetz)	[text, analysis, method, importance, lie, role, information, extraction]

- **Position filtering:** a custom `TokenFilter` has been implemented to set the `positionIncrement` attribute of all tokens to a specific value. This will be useful to ignore the `positionIncrement` due to removed stopwords in the search phase when we will use the proximity between tokens, as explained in Section 3.4.3;
- **Stemming:** this process is useful to reduce words to their base form, in this system we have tried using the Snowball English (Porter2) [11] and the Krovetz [12] stemmers.

In Table 2 we show an example of the analyzing process for the English language.

3.3. Indexer

Indexing is a crucial step where we create a searchable database, called *index*, for the parsed documents. The index contains important information about the documents, such as the words

Table 3
Indexing performances

Collection	Docs size (GB)	Stoplist	Stemmer	Body terms	Index size (GB)	Time (s)
French	7.99	Default	Snowball	7,497,875	6.98	1224
		50 most freq.	Light	7,459,058	6.95	842
English	7.27	Default	Snowball	7,253,947	6.49	1041
		50 most freq.	Krovetz	7,451,647	6.43	848

and phrases they contain, their frequency, and their location within the document. Indexing allows us to store the documents in a *structured* manner, which greatly speeds up the retrieval process by enabling users to search for documents based on keywords or phrases. To achieve this, we developed the following components:

- **DirectoryIndexer**: indexes the documents located in a specified directory and its sub-directories. It accepts various parameters, including the directory containing the documents to be indexed, the `DocumentParser`, the `Analyzer`, the `Similarity` to be used for indexing, the expected number of documents and the location where the index will be stored. Our code ensures that the document directory is readable and the index directory is writable before initiating the indexing process. Additionally, it keeps track of statistics, such as the number of indexed files and documents.
The main component of the class is the `index` method, which is in charge of performing the actual indexing of the documents. This method iterates through the documents in the directory, extracting their content and adding it to the index. During all the iterations, some statistics indexing is given, such as the time taken every 10 thousand documents. Finally, the index is closed.
- **BodyField**: represents the body field of a document. This field has the following properties:
 - it is *tokenized*, meaning that the body is broken into words, or tokens, to make the search more accurate and flexible;
 - *frequencies* and also the *positions* of the tokens are stored, in order to allow for phrase queries with proximity, as explained in Section 3.4.3;
 - the body content is *stored*, even if this had an impact on the index size, this was needed in the search phase in order to pass documents bodies to the reranker, as explained in Section 3.4.5.

In Table 3 are reported the index performances obtained using the analyzers described in Section 3.2 and the setup described in Section 4.

3.4. Searcher

In the searcher we have used a *boolean query* approach, in this way it was possible to create complex queries by combining, using the *boolean operators*, different components to be matched.

Table 4
NDCG results with different BM25 parameters

Run	FADERIC_French-Stop50-Stem-Shingle-Fuzzy								
Measure	nDCG								
	k1								
	0.6	0.8	1.0	1.2	1.4	1.6	1.8	2.0	
b	0.30	0.3941	0.3952	0.3954	0.3957	0.3963	0.3936	0.3948	0.3934
	0.40	0.3967	0.398	0.3984	0.3992	0.3992	0.3992	0.3981	0.3975
	0.50	0.3999	0.4004	0.4008	0.4013	0.4014	0.4014	0.4014	0.4011
	0.60	0.3999	0.4013	0.4017	0.4025	0.4026	0.4024	0.4019	0.4008
	0.70	0.4021	0.4029	0.4034	0.4038	0.4043	0.4047	0.4046	0.4039
	0.75	0.4018	0.4028	0.4035	0.4039	0.4038	0.4043	0.4037	0.4035
	0.80	0.401	0.4025	0.4033	0.404	0.4041	0.4043	0.4047	0.4039
	0.90	0.3985	0.3998	0.4009	0.4014	0.4018	0.4021	0.4015	0.4011

The components we have used in our queries are explained from Section 3.4.1 to Section 3.4.5. Note that in the following, the word N-grams component explained in Section 3.4.3 will be referenced also with the Lucene's jargon *shingles*.

3.4.1. BM25

Since the document collection has a significant size, as a base of our system we opted for a classic BM25 [13] approach due to the *efficiency* of it and higher *effectiveness* compared to other methods.

The Lucene implementation has default values $k1 = 1.2$ and $b = 0.75$ which worked fine, but we decided to *fine-tune* the parameters to improve measures, in particular, *Normalized Discounted Cumulated Gain (nDCG)* since it is the relevant measure for the LongEval campaign. The results of our experiments are reported in Table 4.

The best result was given by two pairs: $k1 = 1.8, b = 0.8$ and $k1 = 1.6, b = 0.7$; we chose pair $k1 = 1.8, b = 0.8$ due to having the highest *Mean Average Precision (MAP)*. Note that the performance gain was minimal, just 0.006 from the score with default parameters, which was expected.

3.4.2. Fuzzy

A *fuzzy search*, or approximate search, is a technique used to search for approximate or *partial* matches between a search term and documents in a collection of texts. Unlike exact search, in which the match must be exact and precise, fuzzy search allows you to find results even when the words you search for do not exactly match those in the documents. Fuzzy search is particularly useful when you want to get results even if there are *misspellings*, *language variants*, *abbreviations* or other forms of variations in the search terms or texts of the documents. For example, if you search for the term "roam" with a fuzzy search, the document containing the term "foam" might also be returned.

Fuzzy search techniques are based on the use of algorithms that evaluate the similarity between text strings. One of the most common algorithms used for fuzzy search is the Levenshtein algorithm [14], which calculates the edit distance between two strings, that is, the minimum number of operations (character insertions, deletions and substitutions) required to transform one string into the other. Lucene, for example, uses a variant of the algorithm just described, the Damerau-Levenshtein algorithm, named after the Damerau algorithm [15], which also considers character transpositions among its allowable operations.

Lucene also allows you to add an additional (optional) parameter with which to specify the maximum number of changes allowed. In our case we decided to set a manual *threshold* to choose the value of the parameter; if the word length is greater than or equal to the threshold then the fuzzy parameter is set to 2 otherwise 1 is used. The threshold is called "fuzzyThreshold" and can be set in the configuration file; we decided to set it to 10. Finally, to avoid performance degradation, in our IR system, fuzzy search is applied only if the query contains a single term.

3.4.3. Word N-grams

Word N-grams are a sentence analysis technique that consists in dividing the words of a sentence into sequences of n consecutive words. For example, the sentence "the dog barks," can generate the N-grams "the dog" and "dog barks". Word N-grams are useful because they capture *local relationships* between words within a text, so this approach helps identify similar, though not identical, phrases and can *improve search relevance*. The maximum number of words within an N-gram can be set in the configuration file in `maxShingleSize`; in our case, we decided to use a maximum of 3 words. Also, we avoided generating unigrams as they do not capture any local relationships within the query.

We then decided to set up a proximity search within each N-gram. The proximity of terms can be used to identify documents in which the search terms occur in a certain *spatial relationship*. For example, if we are searching for the terms "dog" and "brown" with the proximity of 3 words, we want to find documents in which these two terms appear within a maximum distance of three words from each other. Thus, if a document contains sentences such as "I saw a brown dog in the park" or "The brown dog was running fast", these documents would be considered relevant because the terms "dog" and "brown" are close to each other. In our case, the proximity parameter is set to 5.

Finally, we applied a boost to all word N-grams based on the number of generated N-grams for a certain query.

3.4.4. Synonyms

Managing synonyms in a retrieval system is not straightforward. The use of synonyms may not necessarily improve results, and this depends on how they are used.

Synonyms may be useful in the context of IR to broaden the search to include more related terms. However, the introduction of synonyms can also create problems such as noise in the information retrieval process. Here are some reasons why the use of synonyms may not improve results:

- **Polysemy:** words may have *more than one meaning*. Introducing synonyms might lead to an increase in irrelevant results if a synonym is associated with another meaning of the searched word. For example, if one searches for "bank" in the context of a financial bank, introducing the synonym "riverbank" could generate irrelevant results.
- **Irrelevant synonyms:** not all synonyms are equally *relevant in the context* of a given query. Some synonyms might be too general or too specific concerning the user's intent, leading to inappropriate or insufficient results.
- **Redundancy:** adding synonyms can lead to redundancy in the answers provided. If multiple synonymous words or phrases are used in the query, there may be significant overlap between the results, reducing the overall effectiveness of the IR system.

However, it should be kept in mind that the effectiveness of using synonyms in improving the results of the IR system also depends on the specific implementation and the characteristics of the domain or context in which the system is used. In some cases, the use or *expansion* of synonyms could actually improve the accuracy of information retrieval.

Furthermore, the use of synonyms can increase computational complexity in the information retrieval process. Because synonyms require accurate correspondence with indexed documents, additional computations must be performed to identify and compare matching synonyms in indexed texts.

We made several attempts to implement synonyms in our system; a summary description of what we did is given below.

Firstly, since handling synonyms in the index takes too much computation time, we decided to handle them directly in the search. We added synonyms in the queries with a query expansion approach.

In addition, we decided to use the WordNet² dictionary. WordNet is a lexical database that groups English words into sets of synonyms called synsets, providing semantic relationships and definitions. It offers a comprehensive resource for natural language processing tasks, such as word sense disambiguation, information retrieval, and sentiment analysis. Since WordNet is written in C, it was necessary to use an additional API in order to use the dictionary on our system. That API is called extJWNL³ (Extended Java WordNet Library) and does not need the WordNet database installed locally.

Also, to improve dictionary lookup, we decided to use the OpenNLP⁴ library for natural language processing and limit polysemy. Each word in the original query was processed by an OpenNLP *Part of Speech (PoS)* Tagger in order to obtain the tag associated with the word. That function analyzes the context in which the word is used and returns the associated tag. The model used for the pos tagger was `en-pos-maxent.bin` and the tags are associated with WordNet section as shown in Table 5.

Knowing the tag of each word in a query made it possible to look up the word in the corresponding dictionary section. For example, for the query "free software," free was searched in the adjective section and software in the noun section. This strategy improved the metrics very little probably because the queries provided by Long Eval are very short, averaging 2/3 words.

²<https://wordnet.princeton.edu/>

³<https://github.com/extjwnl/extjwnl>

⁴<https://opennlp.apache.org/>

Table 5
OpenNLP Tags compared with WordNet Sections

OpenNLP Tag	WordNet Section
JJ	Adjectives
VB	Verbs
RB	Adverbs
NN	Nouns
Others	No synonyms retrieved

In addition, OpenNLP works well with *properly formulated sentences*, including consideration of capitalization and punctuation. In this case, queries are very crudely formulated, for example, some begin with a capital letter and some do not, as a result, OpenNLP does not always provide the correct tag. Then, this strategy might be more useful in the case of more complex queries, such as those characterizing a conversational IR system.

Subsequently, we tried to give a *different boost* to each synonym. As a first approach, we decided to provide a boost based on the amount of synonyms returned. In this case, the boost was calculated in this way:

$$boost = \frac{BoostBase}{SynonymListLength} \quad (1)$$

This approach was used to limit the importance associated with each synonym if the returned synonym list is long, being more likely to get *irrelevant synonyms*.

Finally, we moved synonym management, creating two new Analyzers: `SynonymAnalyzer` and `SynonymPOSAnalyzer`, which are applied only in the search part:

- **SynonymAnalyzer:** uses as input a query already previously analyzed with the standard Analyzer, i.e. `EnglishLongEvalAnalyzer` or `FrenchLongEvalAnalyzer`, after applies: `SynonymGraphFilter`, `FlattenGraphFilter` and `StopFilter`. `SynonymGraphFilter` represents a filter that can be directly applied to a `TokenStream` within an Analyzer. The filter creates a synonym graph based on specified configurations and expands the terms found in the analyzed text by adding the corresponding synonyms to the token graph. `FlattenGraphFilter` converts an incoming graph token stream, such as one from `SynonymGraphFilter`, into a flat form so that all nodes form a single linear chain with no side paths. Every path through the graph touches every node. This is necessary when indexing a graph token stream because the index does not save `PositionLengthAttribute` and so it cannot preserve the graph structure. However, at search time, query parsers can correctly handle the graph and this token filter should not be used. This Analyzer uses a list of synonyms in .txt format, available in two versions: standard and custom. Before being processed by the Analyzer, the synonym list is transformed into a `SynonymMap` object via the `AnalyzerUtil`'s `loadSynonymList` function.
- **SynonymPOSAnalyzer:** takes as input `EnglishLongEvalAnalyzer`, then applies an `OpenNLPPosFilter`, so that each word is assigned the associated tag. Then, it applies

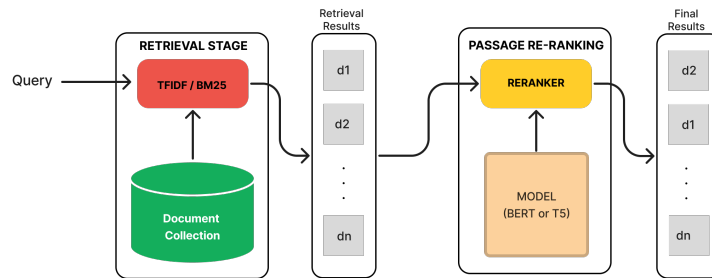


Figure 1: Retrieve-then-rerank framework

a custom filter called `SynonymPOSFilter` to manage the tags and look up words in the WordNet dictionary. Creating a custom filter was not trivial as the information about it in the documentation and online is very limited. The filter allows us to:

1. fetch the input tokens,
2. retrieve their associated tag,
3. search the synonyms in the dictionary based on their tag,
4. process the synonyms with the standard analyzer i.e. `EnglishLongEvalAnalyzer`
5. and return as output a `TokenStream` containing all the synonyms found.

The `tokenStream` returned as output by both Analyzers is then transformed into a list of strings which is in turn processed by the `Searcher` to apply a boost. Finally, the synonyms are added to the `BooleanQuery` along with the original query.

3.4.5. Reranker

After all of the previous steps, we obtained a working system that achieved satisfactory results at a good speed, so we tried to improve it by introducing a second stage, called *passage re-ranking*, in which each of the documents returned by the first retrieval model would be scored and re-ranked by a more computationally-intensive method involving Machine Learning.

We achieved this by leveraging a library called `PyGaggle`⁵, which provides some deep neural architectures for text ranking and question answering, and two transformer-based models, T5 and BERT, with different checkpoints and even our own trained checkpoint using code from another library [16].

In our system, we can choose how many documents are to be reranked from the top for two reasons: the first is raw computing performance, reranking all 1000 retrieved documents takes a long time and we don't have machines powerful enough to handle it; the second is that we saw

⁵<https://github.com/castorini/pygaggle>

that reranking more than 50 documents lowers our measures, even below the baseline measure without reranking.

Our first approach to reranking was to consider only the scores returned by the models to rerank the documents but we then switched to an approach where we consider also the BM25 score as follows. Let $Score_{BM25}(i)$ the score given by BM25 for the document at rank i and $Score_{rr}(i)$ the score given by the reranker for the document at rank i , and let n be the total number of reranked documents, we define:

$$nScore_{rr}(i) = \left(Score_{rr}(i) + \min_{j \in [1, n]} Score(j) \right) \cdot \frac{Score_{BM25}(1)}{Score_{rr}(1)} \quad (2)$$

as the normalized score for the reranked documents, since the models returned a score in the range $[-10, +10]$, which was not suitable for our case.

In our first approach, we simply passed the score to Lucene's ScoreDoc object and we were done. But in this way, we would lose information about the ranking given by BM25, which is still relevant, so we defined a new score:

$$finalScore(i) = mntr + (1 - \alpha) \cdot Score_{BM25}(i) + \alpha \cdot nScore_{rr}(i) \quad (3)$$

where $mntr$ is the maximum score from docs which are not reranked, in this way, we preserve the order of this docs. With this approach, we can give a weight to the reranker to find the balance and we do not lose the information given to us by the first stage. Note that $\alpha = 1$ corresponds to considering only scores from the reranker.

Pretrained Models

We tried two transformer-based models, T5 and BERT since they are supported by PyGaggle. First, we tried T5 [17], but with BERT [18] we got better results. Starting from the same base model, t5-base⁶ for T5 bert-base-uncased⁷ for BERT, we tried different checkpoints⁸ fine-tuned specifically for reranking tasks and we even tried to train our checkpoint. The pre-trained checkpoints that we used are:

- *monot5-base-msmarco-10k*⁹
- *bert-base-mdoc-bm25*¹⁰

The results for the T5 model are reported in Table 6 and the results for the BERT model are reported in Table 7.

The BERT pretrained checkpoint with 20 reranked documents improved nDCG by 3.56% and MAP by 8.5%.

Training our own checkpoint

At this point, BERT gave us good results so we took it a step further and we tried to find ways to finetune it to our data. The training process is pretty straightforward:

⁶<https://huggingface.co/t5-base>

⁷<https://huggingface.co/bert-base-uncased>

⁸saving the model's parameters and optimizer state during the training process

⁹<https://huggingface.co/castorini/monot5-base-msmarco-10k>

¹⁰<https://huggingface.co/Luyu/bert-base-mdoc-bm25>

Table 6

monot5-base-msmarco-10k model with different number of documents to rerank

	nDCG	MAP
0	0.4075	0.2411
10	0.414	0.2502
20	0.4119	0.2477
50	0.4083	0.242
100	0.405	0.2376
250	0.3987	0.2301

Table 7

bert-base-mdoc-bm25 model with different number of documents to rerank

	nDCG	MAP
0	0.4075	0.2411
10	0.4207	0.2608
20	0.4222	0.2617
50	0.4212	0.2598
100	0.4184	0.2563
250	0.4104	0.2478

- **Data pre-processing.** Transformers expect batches of tensors as input, so we need to preprocess our data to the expected format. For processing textual data the tokenizer tool is used, which splits text into tokens; in our case, we exploit the pre-trained BERT tokenizer which returns tokens that are not necessarily words, but rather subwords: frequently used words are (or should) not split into smaller subwords, but rare words should be decomposed into meaningful subwords [19]. Further, the tokenizer adds, at the beginning and at the end, two special tokens, respectively [CLS] and [SEP]. The tokenizer returns a dictionary with three items:

- `input_ids`, indices corresponding to each token in the sentence
- `attention_mask`, indicates whether a token should be attended to or not
- `token_type_ids`, identifies which sequence a token belongs to when there is more than one sequence

An important note is that BERT accepts input sequences of up to 512 tokens, so the tokenizer truncates longer documents.

- **Training.** The easiest way to train is to use the Trainer¹¹ API from PyTorch.

To ease development, we used a package for training deep language model rerankers [16] and adapted the example code to our collection.

The `convert_to_training.py` takes care of converting data to the training file, given the ranking file, the qrels, the query collection and the docs collection. Then it is sufficient to use the `trainer.py` code to get the trained model. Unfortunately, we don't have access to

¹¹https://huggingface.co/docs/transformers/main_classes/trainer

Table 8

Our trained model with different number of documents to rerank

	nDCG	MAP
0	0.4075	0.2411
10	0.3910	0.2253
20	0.3741	0.1975
50	0.3405	0.1580

sufficiently powerful machines so we were forced to train on a Google Colab notebook. This came with a major drawback: the maximum runtime is 12 hours, so we couldn't train on more than 1 epoch since a 2 epoch model was estimated to take 15 hours to train. This obviously tanked our model performances, and, as we can see in Table 8.

The Colab notebook with all the hyperparameters can be found at <https://colab.research.google.com/drive/1oFeYSkR31A-MUibwWrfNcKLUyPxEYbOq?usp=sharing>. The trained model can be found at https://huggingface.co/enricobolzonello/clef_longeval.

Integrating with Lucene

One problem that emerged while working with the reranker was integrating it with Lucene since the reranker is written in Python and our main program is in Java. To solve the issue we came up with three approaches:

1. passing intermediate text files, with documents, query to the reranker and returning the ranking to Java. This solution was used for initial testing but was deemed too inefficient and prone to errors
2. using Python as the system's entry point and using the PyJNIus¹² library to access Java classes. The same approach was used by Birch [6], but we should have changed the classes too much to integrate tightly the reranker and more importantly we didn't want to change the entry point of our system
3. our final solution was to call in some way Python from Java

The library we used to achieve this is JEP¹³, which uses *Java Native Interface (JNI)* and the CPython API to start up the Python interpreter inside the *Java Virtual Machine (JVM)*. Thanks to the Python interpreter, we can call, at each query, the reranker which returns a list of generic Objects that is converted to a list of Float.

4. Experimental Setup

The experimental setup of this project consists of the following:

- The *project* is available at <https://bitbucket.org/upd-dei-stud-prj/seupd2223-faderic/src/master/>;
- The *collections* are taken from <https://clef-longeval.github.io/data/>;

¹²<https://github.com/kivy/pyjnius>

¹³<https://github.com/ninia/jep>

- The *evaluation tool* used is `trec_eval` v9.0.7, available https://trec.nist.gov/trec_eval/;
- To compute the runs we have used the following *hardware*: CPU AMD Ryzen 7 2700, GPU Zotac GeForce RTX 2060 6 GB, RAM 16 GB DDR4;

In order to reproduce the runs for this system, it is necessary to follow the instructions given in the project’s ReadMe file.

5. Results

In this Section we will show and analyze the results obtained by our system. In particular, Section 5.1 will be about the results of all the runs produced by our system on the training collection during its development, while Section 5.2 will consist of a statistical analysis of the runs submitted to CLEF on the test collections.

5.1. Training results

We have *combined* the different components explained in Section 3 and we conducted a thorough experimentation process to identify the best-performing system for our task. By tuning the parameters of each component and selecting the optimal ones, we were able to build a system that addresses the persistence issue of IR systems. In Table 9 are reported the MAP and nDCG scores obtained on those systems.

The keywords reported in the names of the runs have the following meanings:

- French: used documents and queries from the French collection
- English: used documents and queries from the English collection

Table 9
nDCG and MAP values on train collection

Run name	nDCG	MAP
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom-Rerank20W6	0.4274	0.2671
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-Rerank30	0.4230	0.2632
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom	0.4079	0.2416
FADERIC_French-BM25Tuned-Stop50-LightStem-Shingle-Fuzzy	0.4047	0.2383
FADERIC_French-BM25-StopDefault-SnowStem	0.3786	0.2110
FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS-Rerank30	0.3271	0.1877
FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-Rerank30	0.3527	0.1873
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-TrainedRerank30	0.3599	0.1799
FADERIC_French-LMDirichlet-Stop50-LightStem	0.3398	0.1731
FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS	0.3081	0.1634
FADERIC_English-BM25-StopDefault-SnowStem	0.2927	0.1490
FADERIC_English-LMDirichlet-Stop50-KStem	0.2612	0.1228

- BM25: used Okapi BM25 similarity (with default parameters)
- BM25Tuned: used Okapi BM25 similarity, with parameters tuned on training collection: $k_1=1.6$ and $b=0.7$
- LMDirichlet: used Dirichlet smoothing (with default parameter)
- StopDefault: used Lucene's default stop words list
- Stop50: used stoplist built by picking the 50 most frequent terms in the documents indexed without stoplist and stemming
- LightStem: used the Light stemmer
- KStem: used Krovetz stemmer
- SnowStem: the Snowball stemmer
- Shingle: used word N-grams (max window size = 3) query expansion
- Fuzzy: used fuzzyness (threshold parameter = 10) query expansion
- SynCustom: used custom synonyms list
- SynPOS: used WordNet synonym list together with OpenNLP PoS tagging
- Rerank20W6: used reranker, reranking 20 documents with weight 0.6 given to the reranker scores
- Rerank30: used reranker, reranking 30 documents with weight 1 given to the reranker scores
- TrainedRerank30: used reranker, reranking 30 documents using our custom model

During the tuning of our system, we primarily focused on MAP and nDCG. However, in order to perform a more *comprehensive* analysis, we also considered additional measures such as Precision and Recall: the first one is the fraction of retrieved documents that are relevant to the user's query, indicating a measure of the accuracy of the system, while with the second is a measure of the completeness of the system in retrieving all relevant results, computed by the fraction of relevant documents retrieved.

In Figure 2, we show the *interpolated* Precision-Recall curve, which can be useful to show the inverse relationship between Precision and Recall, indicating the *trade-off* between these two measures.

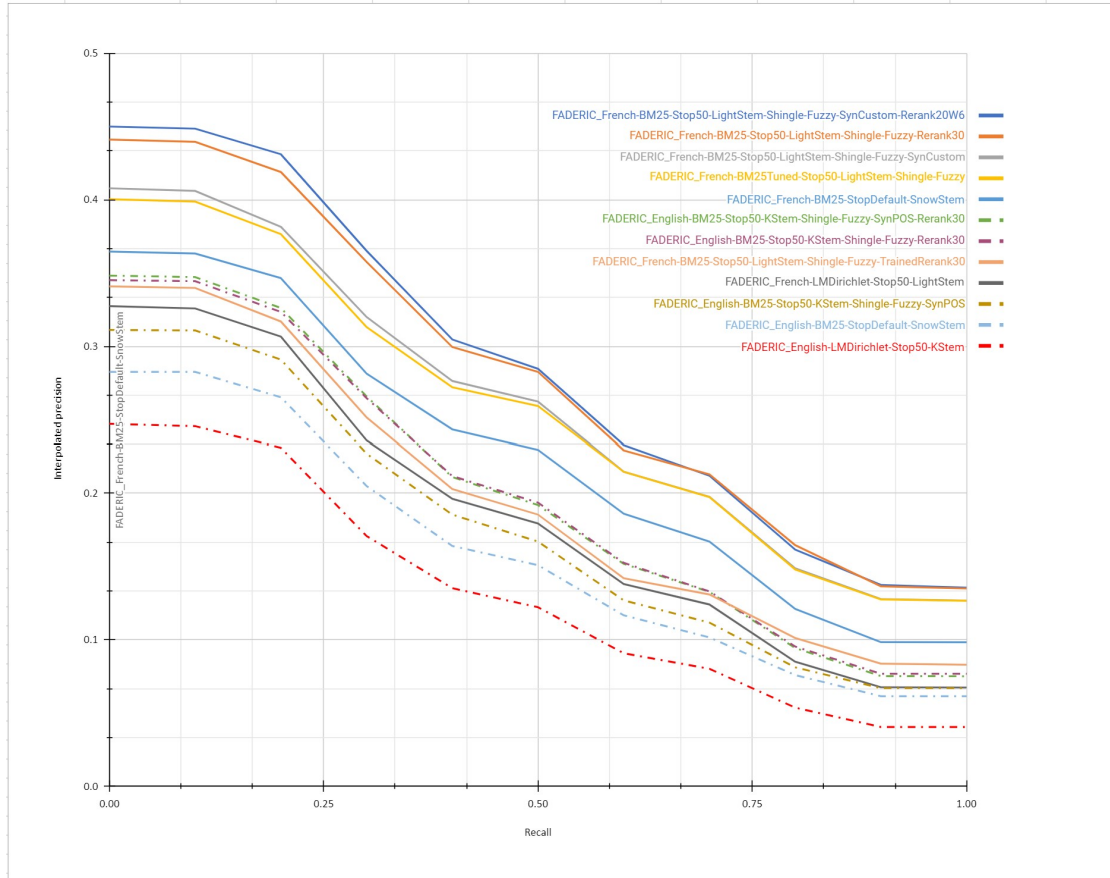


Figure 2: Interpolated Precision-Recall curve on train collection

In Table 10 and in Table 11 we have reported a more complete list of scores respectively for the French and English runs. For space reasons, we label the runs as:

- fr_1 = FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom-Rerank20W6
- fr_2 = FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-Rerank30
- fr_3 = FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom
- fr_4 = FADERIC_French-BM25Tuned-Stop50-LightStem-Shingle-Fuzzy
- fr_5 = FADERIC_French-BM25-StopDefault-SnowStem
- fr_6 = FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-TrainedRerank30
- fr_7 = FADERIC_French-LMDirichlet-Stop50-LightStem
- en_1 = FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS-Rerank30
- en_2 = FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-Rerank30
- en_3 = FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS
- en_4 = FADERIC_English-BM25-StopDefault-SnowStem
- en_5 = FADERIC_English-LMDirichlet-Stop50-KStem

Table 10

Measures for French runs on train collection

runid	all	fr_1	fr_2	fr_3	fr_4	fr_5	fr_6	fr_7
num_q	all	672	672	672	672	672	672	672
num_ret	all	660838	660838	660838	660838	658471	660838	658512
num_rel	all	2626	2626	2626	2626	2626	2626	2626
num_rel_ret	all	2316	2318	2316	2323	2271	2318	2156
ndcg	all	0.4274	0.4230	0.4079	0.4047	0.3786	0.3599	0.3398
map	all	0.2670	0.2632	0.2416	0.2383	0.2110	0.1799	0.1731
gm_map	all	0.0786	0.0765	0.0702	0.0696	0.0545	0.0552	0.0394
Rprec	all	0.2280	0.2278	0.1987	0.1946	0.1747	0.1365	0.1469
bpref	all	0.4128	0.4122	0.4085	0.4063	0.3753	0.3701	0.3561
recip_rank	all	0.4222	0.4140	0.3824	0.3734	0.3424	0.3170	0.3134
iprec_at_recall_0.00	all	0.4499	0.4410	0.4078	0.4003	0.3646	0.3409	0.3274
iprec_at_recall_0.10	all	0.4485	0.4395	0.4060	0.3987	0.3633	0.3398	0.3258
iprec_at_recall_0.20	all	0.4310	0.4189	0.3815	0.3765	0.3465	0.3169	0.3066
iprec_at_recall_0.30	all	0.3651	0.3575	0.3200	0.3131	0.2813	0.2515	0.2359
iprec_at_recall_0.40	all	0.3046	0.2995	0.2763	0.2720	0.2433	0.2030	0.1963
iprec_at_recall_0.50	all	0.2846	0.2825	0.2623	0.2592	0.2296	0.1855	0.1795
iprec_at_recall_0.60	all	0.2328	0.2293	0.2147	0.2147	0.1861	0.1421	0.1381
iprec_at_recall_0.70	all	0.2120	0.2130	0.1976	0.1974	0.1671	0.1310	0.1242
iprec_at_recall_0.80	all	0.1616	0.1647	0.1489	0.1481	0.1212	0.1013	0.0851
iprec_at_recall_0.90	all	0.1375	0.1366	0.1279	0.1276	0.0985	0.0838	0.0677
iprec_at_recall_1.00	all	0.1357	0.1351	0.1268	0.1267	0.0984	0.0831	0.0675
P_5	all	0.2074	0.2065	0.1875	0.1827	0.1637	0.1315	0.1375
P_10	all	0.1603	0.1560	0.1472	0.1469	0.1332	0.1007	0.1088
P_15	all	0.1233	0.1243	0.1183	0.1167	0.1090	0.0869	0.0889
P_20	all	0.0990	0.1022	0.0990	0.0990	0.0901	0.0799	0.0753
P_30	all	0.0748	0.0743	0.0748	0.0743	0.0683	0.0742	0.0580
P_100	all	0.0279	0.0279	0.0279	0.0280	0.0267	0.0279	0.0232
P_200	all	0.0150	0.0150	0.0150	0.0151	0.0146	0.0150	0.0133
P_500	all	0.0065	0.0065	0.0065	0.0065	0.0064	0.0065	0.0060
P_1000	all	0.0034	0.0034	0.0034	0.0035	0.0034	0.0034	0.0032
recall_5	all	0.2738	0.2732	0.2478	0.2417	0.2106	0.1749	0.1802
recall_10	all	0.4167	0.4042	0.3830	0.3828	0.3402	0.2606	0.2789
recall_15	all	0.4724	0.4757	0.4560	0.4499	0.4203	0.3366	0.3380
recall_20	all	0.5034	0.5193	0.5032	0.5051	0.4595	0.4119	0.3814
recall_30	all	0.5710	0.5657	0.5710	0.5677	0.5187	0.5651	0.4412
recall_100	all	0.7022	0.7019	0.7022	0.7028	0.6688	0.7019	0.5816
recall_200	all	0.7555	0.7560	0.7555	0.7587	0.7283	0.7560	0.6729
recall_500	all	0.8219	0.8232	0.8219	0.8235	0.7994	0.8232	0.7574
recall_1000	all	0.8663	0.8662	0.8663	0.8685	0.8485	0.8662	0.8072

Table 11

Measures for English runs on train collection

runid	all	en_1	en_2	en_3	en_4	en_5
num_q	all	672	672	672	671	671
num_ret	all	655602	655602	655602	654919	654764
num_rel	all	2626	2626	2626	2623	2623
num_rel_ret	all	1924	1915	1924	1878	1765
ndcg	all	0.3271	0.3257	0.3081	0.2927	0.2612
map	all	0.1877	0.1873	0.1634	0.1490	0.1228
gm_map	all	0.0221	0.0215	0.0195	0.0164	0.0122
Rprec	all	0.1706	0.1708	0.1385	0.1253	0.1092
bpref	all	0.3536	0.3524	0.3417	0.3263	0.3139
recip_rank	all	0.3322	0.3289	0.2967	0.2697	0.2364
iprec_at_recall_0.00	all	0.3482	0.3451	0.3111	0.2825	0.2471
iprec_at_recall_0.10	all	0.3472	0.3444	0.3108	0.2825	0.2455
iprec_at_recall_0.20	all	0.3261	0.3234	0.2909	0.2652	0.2310
iprec_at_recall_0.30	all	0.2658	0.2646	0.2270	0.2050	0.1709
iprec_at_recall_0.40	all	0.2111	0.2118	0.1855	0.1641	0.1353
iprec_at_recall_0.50	all	0.1920	0.1937	0.1671	0.1510	0.1223
iprec_at_recall_0.60	all	0.1519	0.1526	0.1270	0.1168	0.0909
iprec_at_recall_0.70	all	0.1328	0.1331	0.1118	0.1017	0.0803
iprec_at_recall_0.80	all	0.0945	0.0955	0.0813	0.0760	0.0538
iprec_at_recall_0.90	all	0.0753	0.0769	0.0671	0.0616	0.0406
iprec_at_recall_1.00	all	0.0752	0.0769	0.0671	0.0616	0.0406
P_5	all	0.1542	0.1542	0.1345	0.1195	0.1019
P_10	all	0.1137	0.1131	0.1042	0.0954	0.0793
P_15	all	0.0904	0.0898	0.0835	0.0784	0.0650
P_20	all	0.0745	0.0732	0.0705	0.0662	0.0543
P_30	all	0.0536	0.053	0.0537	0.0513	0.0433
P_100	all	0.0208	0.0209	0.0208	0.0201	0.0179
P_200	all	0.0116	0.0115	0.0116	0.0113	0.0103
P_500	all	0.0052	0.0052	0.0052	0.0051	0.0048
P_1000	all	0.0029	0.0028	0.0029	0.0028	0.0026
recall_5	all	0.2020	0.2017	0.1710	0.1512	0.1329
recall_10	all	0.2909	0.2887	0.2628	0.2437	0.1988
recall_15	all	0.3430	0.3404	0.3155	0.2970	0.2432
recall_20	all	0.3745	0.3683	0.3523	0.3318	0.2727
recall_30	all	0.4010	0.3970	0.4018	0.3861	0.3281
recall_100	all	0.5140	0.5167	0.5151	0.500	0.4486
recall_200	all	0.5769	0.5754	0.5774	0.5630	0.5180
recall_500	all	0.6607	0.6597	0.6612	0.6443	0.6087
recall_1000	all	0.7186	0.7150	0.7186	0.7027	0.6635

Based on these results, we can derive the following considerations:

- the runs on the French collections are significantly better than the ones on the English one, the reason for this is the *automatic translation* of the document collection, which has led to errors and inconsistencies in the English one;
- the system performs better when configured to use BM25 similarity instead of the Dirichlet smoothing;
- the *tuned* parameters on BM25 just gave slightly better results, probably *overfitting* the run we have tuned them on, for this reason we have chosen to use the default ones in most of the runs;
- the custom stoplist we have generated by picking the most frequent terms has outperformed Lucene's default ones because since it was based on the specific collection the *effectiveness* of using a stoplist has been maximized;
- in both French and English the Snowball stemmer performed worse than the Light and Krovetz stemmer, respectively;
- the use of word N-grams improved the performances, allowing to have more *contextual matches* by looking for group of words instead of single ones;
- synonyms have slightly improved the performances, this is due to the fact that sometimes they can be *misleading* and retrieve documents that are not contextual with the query;
- the reranking is a very *powerful* tool that has given a huge performance increase to our runs.

Based on the previous results and these considerations, we have decided to submit to CLEF the following systems:

- FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom-Rerank20W6, i.e., the best system overall;
- FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS-Rerank30, i.e., the best system overall on the English collection;
- FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom, i.e., the best system without the use of reranking;
- FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-Rerank30, i.e., the best system without the use of synonyms;
- FADERIC_French-BM25Tuned-Stop50-LightStem-Shingle-Fuzzy, i.e., the best system without the use of both synonyms and reranking;

5.2. Test results

In this Section, we will analyze the performance of the submitted runs on heldout, short term and long term collections. At first, we will tackle the performance changes and then we will perform a *statistical analysis*. In this last part, we will use ANOVA2 with a significance level $\alpha = 0.05$ in order to find out if we can reject the *null hypothesis*, i.e. there is no significant statistical difference between the results of the given runs.

Table 12

nDCG and MAP values on heldout collection

Run name	nDCG	MAP
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom-Rerank20W6	0.4169	0.2474
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-Rerank30	0.4147	0.2416
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom	0.4080	0.2376
FADERIC_French-BM25Tuned-Stop50-LightStem-Shingle-Fuzzy	0.4044	0.2324
FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS-Rerank30	0.3030	0.1626

Then we will use Tukey’s *Honestly Significant Difference (HSD)* test to perform *multiple* pairwise comparisons and determine which specific run means differ significantly from each other.

Since we have submitted four runs performed on the French collection and one on the English collection, in the statistical analysis we will compare to each other only the French runs.

5.2.1. Heldout

In Table 12 are reported the nDCG and MAP values obtained from the submitted runs on the heldout collection, while in Figure 3 is reported the interpolated Precision-Recall curve. Comparing these results with the ones obtained on training, shown in Table 9 and Figure 2 respectively, we can see that every run suffered a *performance drop* over all the measures. This worsening was *expected* and it can be due to the fact that the system has been tuned over a different set of queries. It should also be noticed that this set of queries is more than 6 times smaller compared to the training one, therefore the presence of some *outliers* could have caused the mean performances to drop and to not be a good descriptor of the system.

Observing the nDCG and *Average Precision (AP)* boxplots, shown in Figure 4, we can notice that the runs performed on the French collection have a similar structure in terms of median values and interquartile ranges. We can also notice that, in the AP boxplot, the reranked runs fr_1, fr_2 have longer whiskers, while the others show the presence of outliers.

From the ANOVA2 analysis, which results are reported in Table 13, we can see that $p\text{-value} > \alpha$, therefore we *cannot reject* the null hypothesis. Moreover, from Tukey’s HSD multiple comparison shown in Figure 5, we can derive that the French runs can be considered to be similar to each other.

5.2.2. Short term

In Table 14 are reported the nDCG and MAP values obtained from the submitted runs on the short term collection, while in Figure 6 is reported the interpolated Precision-Recall curve. Comparing these results with the ones obtained on heldout, shown in Table 12 and Figure 3 respectively, we can see that every run has *increased* its performances. This improvement was *not expected* since the performances should tend to drop over time. This can be due to the fact that this set of runs is almost 9 times bigger than the heldout, therefore we can consider the mean measures obtained to be more *reliable* than the ones on the heldout collection.

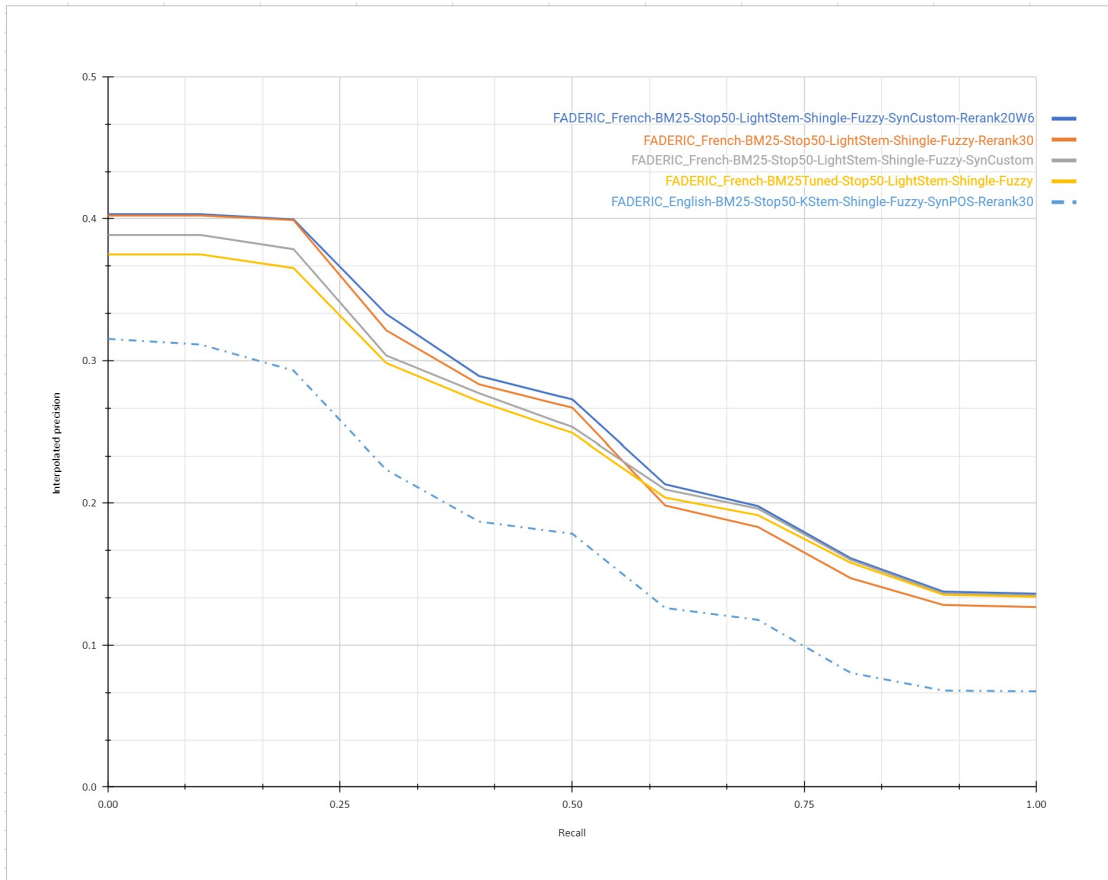
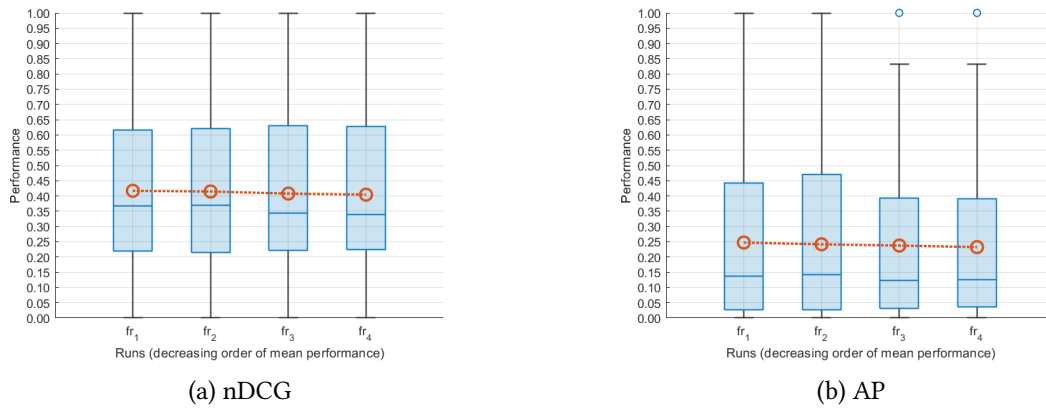


Figure 3: Interpolated Precision-Recall curve on heldout collection



(a) nDCG

(b) AP

Figure 4: Box plot on heldout collection, the mean values are shown in red

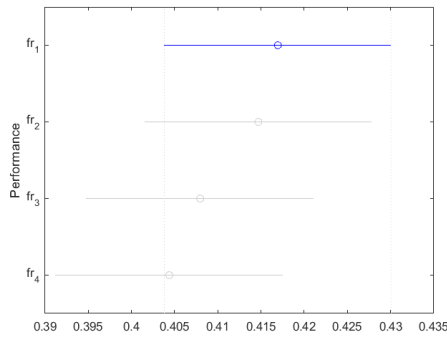
Table 13
ANOVA2 on heldout collection

(a) nDCG

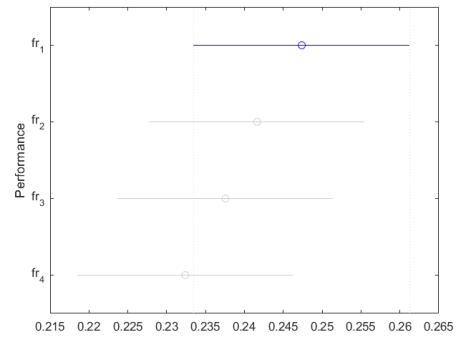
Source	SS	df	MS	F	Prob>F
Systems	0.01	3	0.003	0.64	0.58
Topics	23.54	97	0.242	47.20	1.97E-134
Error	1.49	291	0.005	-	-
Total	25.04	391	-	-	-

(b) AP

Source	SS	df	MS	F	Prob>F
Systems	0.01	3	0.003	0.68	0.55
Topics	23.35	97	0.240	42.10	8.26E-128
Error	1.66	291	0.057	-	-
Total	25.03	391	-	-	-



(a) nDCG



(b) AP

Figure 5: Tukey's HSD on heldout collection

Observing the nDCG and AP boxplots, shown in Figure 7, we can notice that the runs performed on the French collection have a similar structure in terms of median values and interquartile ranges. We can also notice that in the AP boxplot the fr_1 run has a longer whisker, while the others show the presence of outliers.

From the ANOVA2 analysis, which results are reported in Table 15, we can see that $p\text{-value} < \alpha$, therefore we can reject the null hypothesis. Moreover, from Tukey's HSD multiple comparison shown in Figure 8, we can derive that runs fr_3, fr_4 can be considered similar, while all the other runs differ from each other.

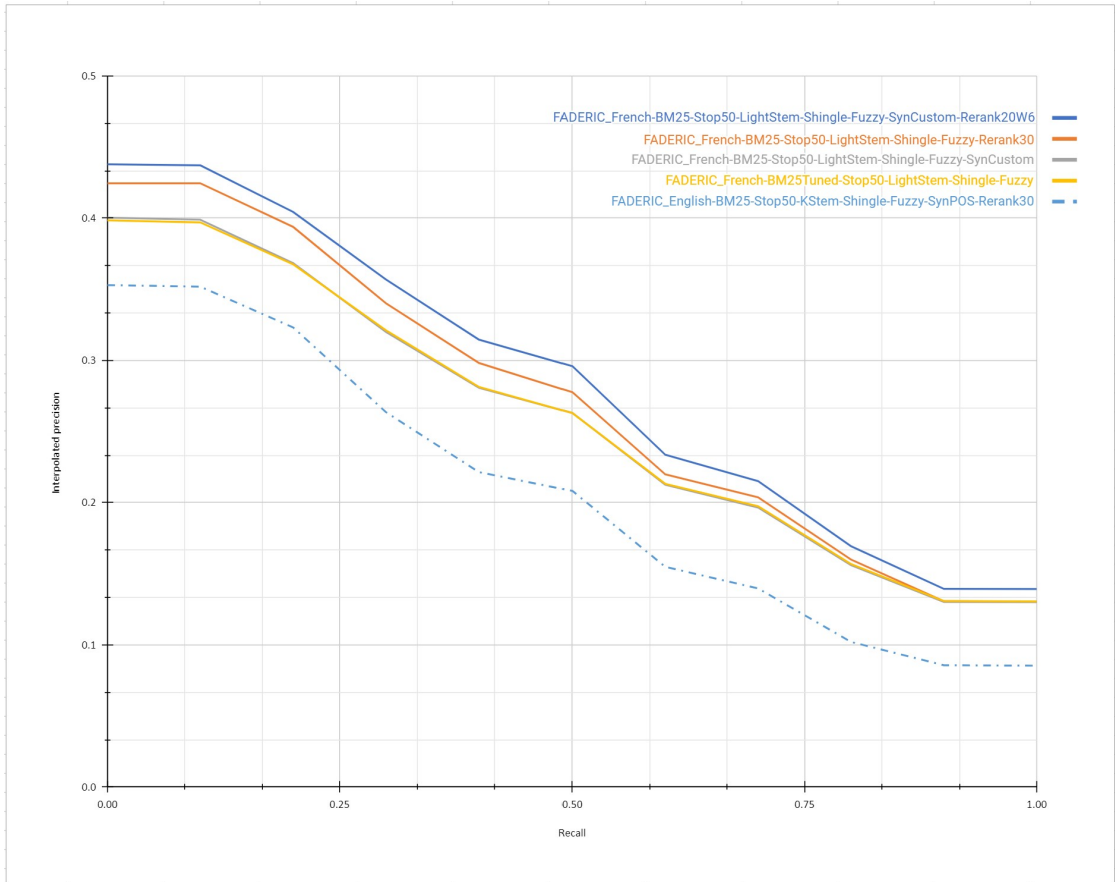


Figure 6: Interpolated Precision-Recall curve on short term collection

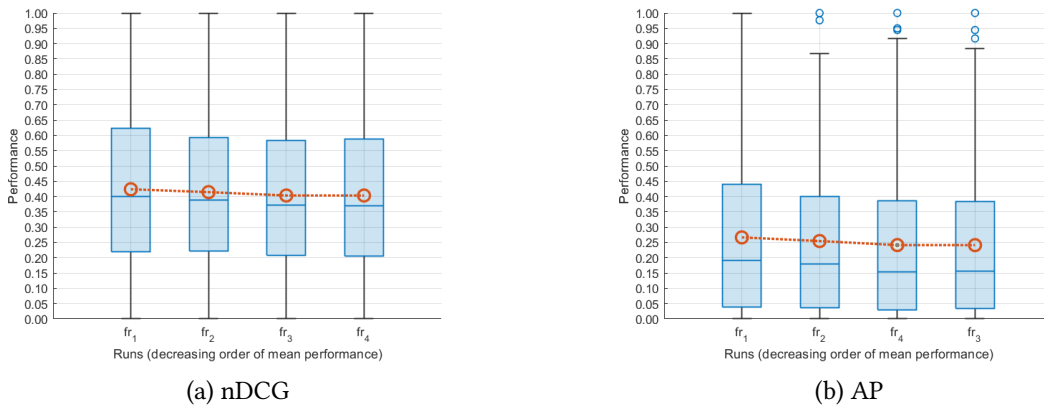


Figure 7: Box plot on short term collection, the mean values are shown in red

Table 14

nDCG and MAP values on short term collection

Run name	nDCG	MAP
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom-Rerank20W6	0.4239	0.2665
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-Rerank30	0.4145	0.2546
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom	0.4034	0.2412
FADERIC_French-BM25Tuned-Stop50-LightStem-Shingle-Fuzzy	0.4034	0.2414
FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS-Rerank30	0.3296	0.1931

Table 15

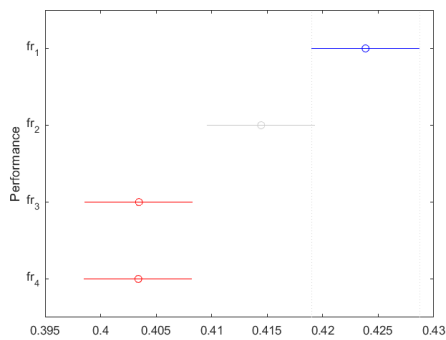
ANOVA2 on short term collection

(a) nDCG

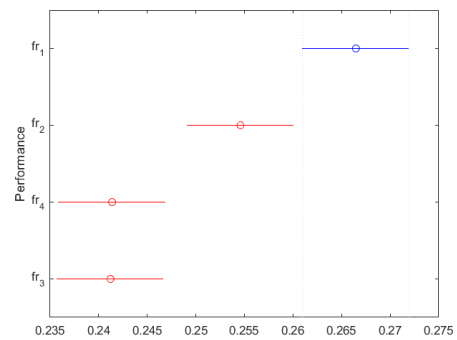
Source	SS	df	MS	F	Prob>F
Systems	0.25	3	0.085	13.58	8.51E-9
Topics	218.58	881	0.248	39.30	0
Error	16.68	2643	0.006	-	-
Total	235.52	3527	-	-	-

(b) AP

Source	SS	df	MS	F	Prob>F
Systems	0.38	3	0.129	16.05	2.43E-10
Topics	205.64	881	0.233	28.92	0
Error	21.32	2643	0.008	-	-
Total	227.36	3527	-	-	-



(a) nDCG



(b) AP

Figure 8: Tukey's HSD on short term collection

Table 16

nDCG and MAP values on long term collection

Run name	nDCG	MAP
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom-Rerank20W6	0.4153	0.2473
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-Rerank30	0.4146	0.2465
FADERIC_French-BM25-Stop50-LightStem-Shingle-Fuzzy-SynCustom	0.4091	0.2384
FADERIC_French-BM25Tuned-Stop50-LightStem-Shingle-Fuzzy	0.4071	0.2350
FADERIC_English-BM25-Stop50-KStem-Shingle-Fuzzy-SynPOS-Rerank30	0.3296	0.1809

5.2.3. Long term

Since some error occurred in the indexing phase when performing the submitted runs fr_1 and fr_2, in the following analysis we will use a *fixed* version of those runs.

In Table 16 are reported the nDCG and MAP values obtained from the submitted runs on the long term collection, while in Figure 9 is reported the interpolated Precision-Recall curve. Comparing these results with the ones obtained on the short term, shown in Table 14 and Figure 6 respectively, we can see that every run suffered a *performance drop* over all the measures. This worsening was expected and it can be due to the fact that the performances tend to drop over time, however, the decrease is not huge and we can consider the performances of the system to *remain satisfactory*.

Observing the nDCG and AP boxplots, shown in Figure 10, we can notice that the runs performed on the French collection have a similar structure in terms of median values and interquartile ranges. We can also notice that in the AP boxplot all the runs show the presence of outliers.

From the ANOVA2 analysis, which results are reported in Table 17, we can see that we obtained very different results for the two measures. In particular, on nDCG we have $p\text{-value} > \alpha$, which means we *cannot reject* the null hypothesis, while on AP we have $p\text{-value} < \alpha$, which means we *can reject* the null hypothesis. The same behavior is reflected in Tukey's HSD multiple comparison, shown in Figure 11, where on nDCG we can derive that the French runs can be considered to be similar to each other, while on AP we can see that runs fr_1, fr_2, fr_3 and runs fr_3, fr_4 can be considered similar to each other.

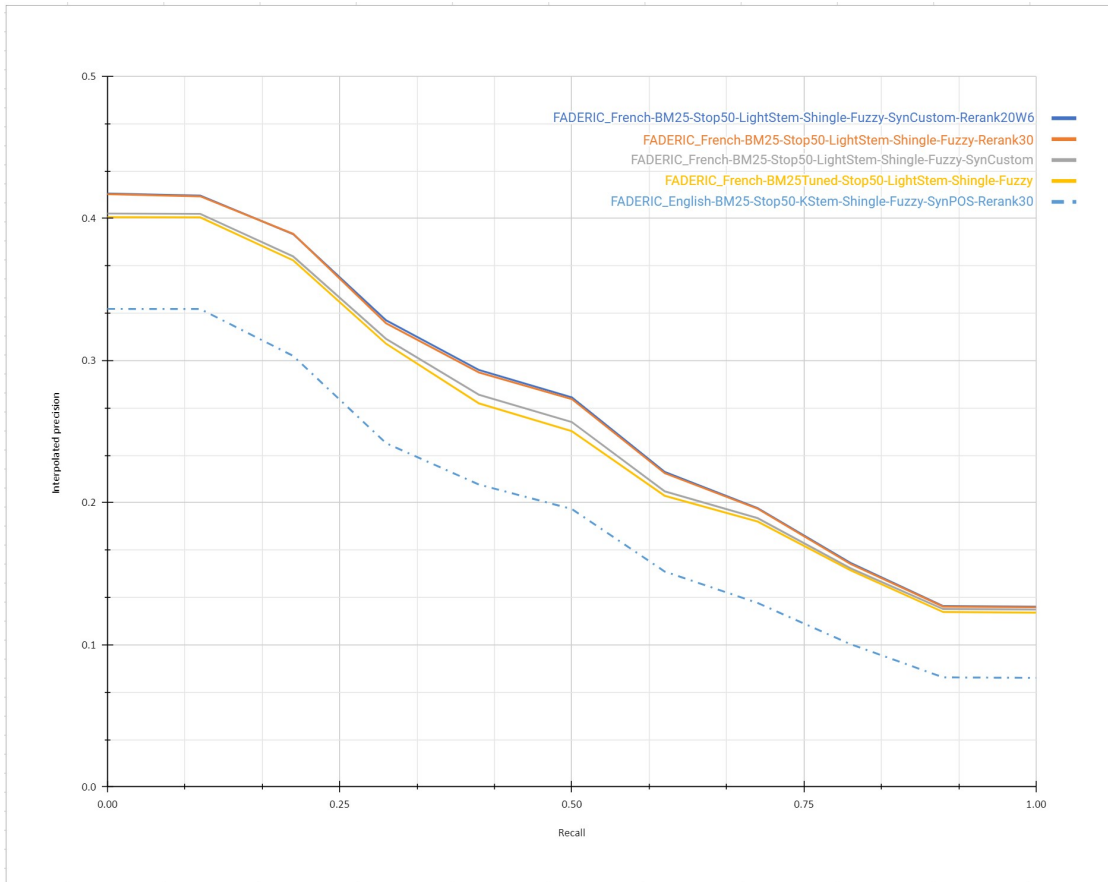


Figure 9: Interpolated Precision-Recall curve on long term collection

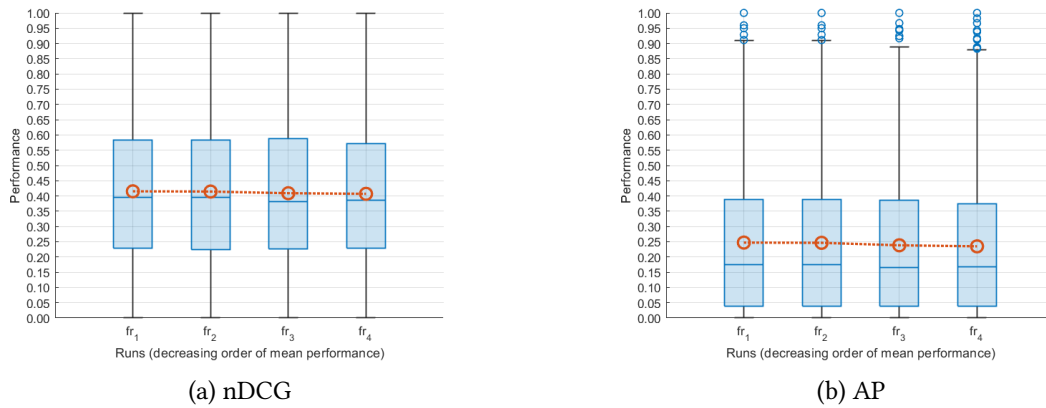


Figure 10: Box plot on long term collection, the mean values are shown in red

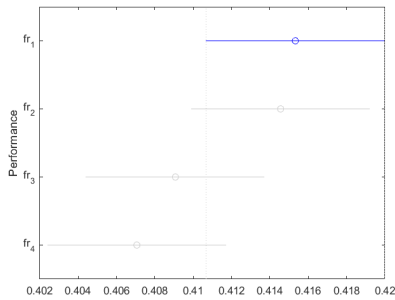
Table 17
ANOVA2 on long term collection

(a) nDCG

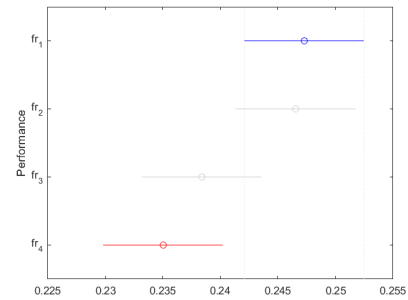
Source	SS	df	MS	F	Prob>F
Systems	0.04	3	0.015	1.51	0.056
Topics	202.35	922	0.219	16.17	0
Error	16.78	2766	0.006	-	-
Total	219.18	3691	-	-	-

(b) AP

Source	SS	df	MS	F	Prob>F
Systems	0.10	3	0.033	4.47	0.003
Topics	188.61	922	0.204	27.03	0
Error	20.92	2766	0.007	-	-
Total	209.64	3691	-	-	-



(a) nDCG



(b) AP

Figure 11: Tukey's HSD on long term collection

6. Conclusions and Future Work

From the obtained results, we can say that the system has kept *satisfactory performance* from the longitudinal evaluation point of view. In fact, on the short term, there has been no significant worsening, while on the long term, there has been just a moderate performance drop.

During the developing of the system we understood that *query expansion* and *reranking* play a major role in the IR systems. Those two features granted us the biggest performance improvements. Although the system has reached decent scores, our work could be *further developed* in different ways.

With respect to query expansion, we could improve the synonyms feature by using other dictionaries, since WordNet is quite outdated and a big portion of the queries were regarding very recent topics, and also by using better French dictionaries, since there are only a few available when working with languages other than English and we had to customize our own.

Another query expansion option could be switching from simple dictionaries to *neural models*, in order to expand a query with words related not just to the meaning of the single words but also to the *context* of the whole topic the user is looking for.

Regarding the reranker, there are some things we could improve, mainly focusing on the *trained model*. As discussed in Section 3.4.5, we did not have the required hardware to properly train a machine learning model, so we trained on Colab with a time limit of 8 hours. This forced us to train with only 1 epoch and we tested only one set of hyperparameters, so, assuming having the necessary hardware, future improvements could focus on training a better model, with more epochs and testing different hyperparameters. More affordable improvements could be done by dropping the libraries for training and inference and working directly with Torch for the interaction with the model, so we could be able to use the latest Torch version, which includes a new implementation of the Transformer API which speeds up training significantly.

References

- [1] R. Alkhalifa, I. Bilal, H. Borkakoty, J. Camacho-Collados, R. Deveaud, A. El-Ebshihy, L. Espinosa-Anke, G. Gonzalez-Saez, P. Galuscakova, L. Goeuriot, E. Kochkina, M. Liakata, D. Loureiro, H. T. Madabushi, P. Mulhem, F. Piroi, M. Popel, C. Servan, A. Zubiaga, Overview of the clef-2023 longeval lab on longitudinal evaluation of model performance, in: Experimental IR Meets Multilinguality, Multimodality, and Interaction. Proceedings of the Fourteenth International Conference of the CLEF Association (CLEF 2023), Lecture Notes in Computer Science (LNCS), Springer, Thessaliniki, Greece, 2023.
- [2] P. Galuščáková, R. Deveaud, G. Gonzalez-Saez, P. Mulhem, L. Goeuriot, F. Piroi, M. Popel, Longeval-retrieval: French-english dynamic test collection for continuous web search evaluation, arXiv preprint arXiv:2303.03229 (2023).
- [3] C. Carpineto, G. Romano, A survey of automatic query expansion in information retrieval, *Acm Computing Surveys (CSUR)* 44 (2012) 1–50.
- [4] H. K. Azad, A. Deepak, Query expansion techniques for information retrieval: a survey, *Information Processing & Management* 56 (2019) 1698–1735.
- [5] R. Alkhalifa, E. Kochkina, A. Zubiaga, Building for tomorrow: Assessing the temporal persistence of text classifiers, 2022. arXiv:2205.05435.
- [6] Z. A. Yilmaz, S. Wang, W. Yang, H. Zhang, J. Lin, Applying BERT to Document Retrieval with Birch, Technical Report, University of Waterloo, 2019.
- [7] C. Chapman, 2022, Unicode Text Segmentation, URL: <https://www.unicode.org/reports/tr29/>.
- [8] A. Freytag, 2004, Character Foldings, URL: <https://www.unicode.org/reports/tr30/tr30-4.html>.
- [9] M. Porter, R. Boulton, F. Brault, 2002, Snowball French stemmer, URL: <https://snowballstem.org/algorithms/french/stemmer.html>.
- [10] J. Savoy, Light stemming approaches for the french, portuguese, german and hungarian languages, in: Proceedings of the 2006 ACM Symposium on Applied Computing, Association for Computing Machinery, 2006, pp. 1031–1035.

- [11] M. Porter, R. Boulton, 2002, Snowball English stemmer, URL: <https://snowballstem.org/algorithms/english/stemmer.html>.
- [12] R. Krovetz, Viewing morphology as an inference process, *Artificial Intelligence* 118 (2000) 277–294.
- [13] S. E. Robertson, U. Zaragoza, The Probabilistic Relevance Framework: BM25 and Beyond, *Foundations and Trends in Information Retrieval (FnTIR)* 3 (2009) 333–389.
- [14] V. I. Levenshtein, et al., Binary codes capable of correcting deletions, insertions, and reversals, in: *Soviet physics doklady*, volume 10, Soviet Union, 1966, pp. 707–710.
- [15] F. J. Damerau, A technique for computer detection and correction of spelling errors, *Communications of the ACM* 7 (1964) 171–176.
- [16] L. Gao, Z. Dai, J. Callan, Rethink training of bert rerankers in multi-stage retrieval pipeline, in: *The 43rd European Conference On Information Retrieval (ECIR)*, 2021.
- [17] C. Raffel, N. Shazeer, A. Roberts, K. Lee, S. Narang, M. Matena, Y. Zhou, W. Li, P. J. Liu, Exploring the limits of transfer learning with a unified text-to-text transformer, *Journal of Machine Learning Research* 21 (2020).
- [18] J. Devlin, M.-W. Chang, K. Lee, K. Toutanova, Bert: Pre-training of deep bidirectional transformers for language understanding, 2019.
- [19] HuggingFace, 2023, Subword tokenization, URL: https://huggingface.co/docs/transformers/tokenizer_summary#subword-tokenization.