# Exercise 1

## WS 2021 - 188.977 Grundlagen des Information Retrieval

## Organization

- Deadline: 9.12. 2021, 23:55 upload in TUWEL
- 3 persons per group
- Exercise will be evaluated based on your submitted report, results, and code
- In case of questions, use the TUWEL forum so that others also benefit from the question/answer. For specific questions, contact me markus.zlabinger@tuwien.ac.at directly.
- Use the programming language Python3

## Goal (Total of 100 points)

The main goal of this exercise is to create, evaluate, and describe your own search engine. To test it, we use a subset of Wikipedia articles. More details are in the notes on the next page. The individual goals of this exercise are the following:

- Implement efficient text processing code to create an inverted index from the given Wikipedia articles, including: tokenization, lowercasing, stemming, stop words. Make sure your preprocessing handles cases such as special characters, multi-blanks, punctuation, ... **(40 points)**

- Implement save/load to disk code for your inverted index **(5  points)**

- Implement query processing: (re-use tokenization, stemming) including: TF-IDF and BM25 scoring methods **(30 points)**
    - Exploration mode: allow free text input and return ranked documents in full text
    - Evaluation mode: read a topic file and output ranked document ids in the evaluation format, described in Section *Information on dataset & evaluation*

- Evaluate the given INEX Wiki test collection **(5 points)**
    - Report MAP, NDCG@10, Precision@10, Recall@10 for TF-IDF and BM25 scoring. Use the previously implemented *Evaluation mode* and *trec_eval* to generate these metrics.

- Write a report using the provided template: **(20 points)**
    - Describe your architecture briefly
    - Note 2 or more interesting details of your work (coming up with the interesting things is part of the exercise)
    - List your performance numbers (indexing time, query time, #indexed documents)
    - List the evaluation results + observations between the two scoring methods

# Bonus points

Do something extra - that we haven't thought of - but it was interesting to you **(up to 5 points)** Note: this extra thing can be some extraordinary efficient code (including, but not limited to: SIMD, GPU acceleration) or some visual analysis of the index, results, etc. or some inspection of the test collection. Describe your attempt to obtain the bonus points in the report.

# Notes

- **Hand-in** Upload your (i) code, (ii) report (as PDF), and (iii) results in TREC format in TUWEL. Use the provided template structure for the source code (folder *code*), results (folder *retrieval_results*) and report (folder *report*). You can find the template in TUWEL. You are allowed to add additional files to your code (e.g. helper functions, etc.).
- You are free to choose which architecture and methods you like best to achieve the goals (e.g. type of inverted index creation process, how you structure your data structures, etc.)
- **Libraries** The idea of this exercise is to write your own text-processing and indexing code. So that means in addition to low-level methods (such as encoding, file i/o, string methods, list, hash-tables, trees, numpy, memory management helpers, etc..) only a stemming library is allowed. Therefore, libraries for tokenization, automatic tf-idf generation, stop word remove methods and so on are not allowed. You can import an NLP or search engine library to get an array of stop words and the stemming class, but not more.
- **Memory management** The corpus contains slightly more than 2 GB of documents. If you index these documents using standard Python data types (e.g. int()), you will quickly drain your memory. We recommend that you store the posting using numpy arrays, where you can define data types for an efficient storage of your posing list. To illustrate the advantage in using numpy arrays, consider the following code example:

  ```
  $ var = np.array([1,2,3,4,5], dtype=np.uint32) # requires 20 bytes of memory
  ```

  ```
  $ var = [1,2,3,4,5] # requires 140 bytes of memory
  ```

- **Report** Describe your prototype using the provided report template. Describe briefly your index and what information is stored at indexing time. Finally, report your evaluation results calculated using *trec_eval*. Try to find a justification of why some of your retrieval models performed better than the others. The report must explain how to run the prototype. Maximum size: 4 pages or 2000 words.
- **Indexing Tips** Start indexing with only a few documents to see if your pipeline works correctly. Do some profiling and timing to identify code that runs slowly and try to fix these bottlenecks. Do sanity checks: e.g. use the *Exploration mode* to search for terms appearing in certain documents and see if your search engine returns the documents in the top ranks.
- You don't have to store the document text for the exploration mode, just save the original filename + start/end bytes and load the text. However, you might want to consider saving the preprocessed text (how it is added to the inverted index) per document to explore your tokenization etc. which gives you an accurate representation of your indexed data.
- You don't have to create a graphic user interface, just use a console based app, that can switch between the exploration and evaluation mode and between TFIDF and BM25.
- Feel free to look at code samples, search engine implementations, algorithm descriptions - but do not copy code. If you find some good references and ideas, add them to the report.

# Information on dataset & evaluation

## Dataset:

- Dataset consisting of Wikipedia articles
  - Documents of a total of 2.2 GB
  - At a minimum, use the *document-ID* and the *body* for indexing. However, you might consider using the other information (such as *category*), e.g., to improve your solution and get some extra points by describing your improvement in the report.
  - The ".xml" documents are not properly structured, and using an XML parser will raise Exceptions. Try an HTML parser instead. You can use an external library for parsing but can also write your own customer parser since the XML structure is very basic.
- "topics" file - with given evaluation queries and their description
- "qrel" file - which contains human relevance judgements for queries in the topics file = allows us to measure the performance of a search engine and compare different configurations (like TF-IDF and BM25)

## Evaluation:

1. Index the articles of the dataset
2. Parse the topics file and get query id and query string (you may choose the title only or also text from the description or narrative)
3. For every query, search in your index for relevant documents and output the top-100 documents per topic in the following format (per line):

```
{topic-id}  Q0  {document-id}  {rank}  {score}  {run-name}
```

- **topic-id** is the id per query from the topics file
- **document-id** is an identifier for the Wikipedia article
- **Q0** is just a legacy hardcoded string
- **rank** is an integer indicating the rank of the doc in the sorted list (normally, this should be ascending from the first line to the last line). Rank starts at 1 (not at 0)
- **score** the similarity score calculated by the BM25 or TF-IDF method
- **run-name** a name you give to your experiment (free for you to choose)

4. Use the trec_eval utility (https://github.com/usnistgov/trec_eval) to compute the performance metrics with the following command:

```
$ trec_eval -m map -m ndcg_cut.10 -m P.10 -m recall.10 path/to/qrel_file
path/to/output_from_3
```

The trec_eval utility has to be compiled with make (on windows: use the Linux subsystem (e.g. bash/ubuntu) for windows)

Put the result in a text file in the "retrieval_results/" folder in your code repository (the name should roughly describe the configuration like: tfidf_title_only.txt)

# References

Wikipedia evaluation collection from: https://www.mpi-inf.mpg.de/departments/databases-and-information-systems/software/inex/

Queries and Ratings (Adhoc Track 2010 Topics) from: https://inex.mmci.uni-saarland.de/data/documentcollection.html