

## Documentation

Text slicing project base of the context limit size.

### Steps of the project:

- Measure length of documents[\*]
- pass to chat gpt if its length was below the limit size[\*]
- slice the document if its size is bigger than the limit size with the aims of the `split_into_slices` function in `nlp_utils.py` file[\*]
  - tokenize the input text including:
    - \* Tokenization
    - \* stopword removal;
    - \* lemmatization
  - calculate the length of tokenized text
  - calculate the number of slices and size of each slice based on the input text size and context limit size
  - calculate the start index and end index of sliced text and slice it
  - check the text is overlap with previous sliced text
  - check the similarity of text with previous sliced text
  - if it does not have overlap and was similar changing the start index
- checking the overlap
  - make the text lower case
  - find all overlapping sequences of a certain length
- checking similarity
  - vectorization and calculation of cosine distance with two kinds of approach (default one is `TfidfVectorizer` which considers the count of occurrences weighted on the length of the document):
    - \* `CountVectorizer`
    - \* `TfidfVectorizer`
    - \* Also another approach is to check the similarity meaning of the word in the text
  - calculate cosine distance and compare it to the threshold

### Project Structure:

- config file: include the config data like context limit size
- llm file: contains the code which is responsible for connecting to the chat gpt api.
- main file: the start point of application
- nlp utils file: contains all the methods which are related to the text processing like:
  - word tokenizing
  - checking similarity
  - checking overlapping text
  - splitting the text

## Running Project

In windows run these commands in rout of project: \* `pip install -r requirements.txt` \* `py ./main.py`

In Linux base systems: \* `pip3 install -r requirements.txt` \* `python3 ./main.py`