<div align="center">

**Wikipedia Discovery**
-using different text summarization-
Extractive Summarization

</div>

Iliescu Mihail Dorin

## Problem statement

This project is aiming to create summary text based on a wikipedia article.

## Motivations

Our motivations for this project were to implement and compare text summarization through different algorithms and state what could be the best ways to achieve text summarization.

## Pipeline

Our program is defined and is excuted in this order:

1. Web crawler/scraper

2. Pre-processsing, tokenization

3. Algorithms computation

4. Results formationg
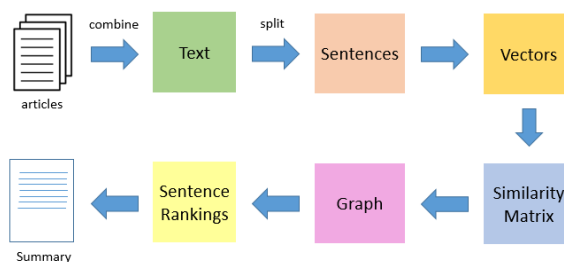
5. Results evaluation

6. Display output

## Proposed solution & Implementation

This section discuss about what methods we decided to use for our project and why.

- Cosine similarity

Cosine similarity through vector space model can be a way to find similarity in dataset. It means that if we can get the most similar sentences to the wole text, we can simply select the n most relevant sentences in the result set.

- TextRank



Derived from Google Page Rank algorithm [1], TextRank algorithm ranks parts of text. This ranking is defined by the number of relation between sentences. We base our evaluation on word frequency such as TF-IDF.

- K-Mean clustering

**Extractive Summarization**

K-Mean clustering algorithm [2] is an unsupervised classification algorithm frequently used in the world of Machine Learning and Data Science. Its main purpose is to, given a n dimension data set, be able to classify this data in categories (called clusters) according to them features. So using this algorithm, we can classify any type of data (images, text...).

Even K-Mean clustering has a data classification purpose, we decided to try to adapt its feature to text summarization. In our case, if we give as input our sentences from original text to K-Mean, it will classify our sentences by topics. Then we can pick the n most relevant sentences from each cluster/topic to form our summary.

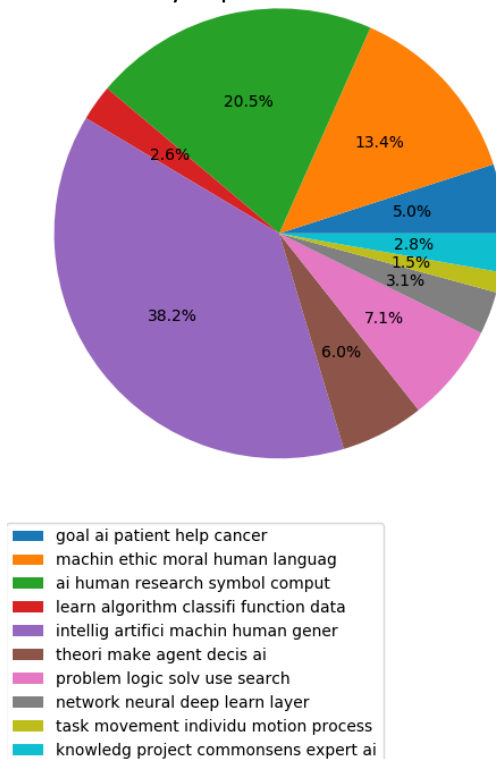I based my implementation on these two papers [3] and [4].



Chart showing data repartition of topic "Artificial Intelligence" over 10 clusters. Each category shows the top words in the cluster (note that words are stemmed).

- goal ai patient help cancer
- machin ethic moral human languag
- ai human research symbol comput
- learn algorithm classifi function data
- intellig artifici machin human gener
- theori make agent decis ai
- problem logic solv use search
- network neural deep learn layer
- task movement individu motion process
- knowledg project commonsens expert ai

To know what are the most important sentences of each cluster, we proceed as follow. After classifying input text in n clusters, we then pick one sentence from each cluster based on minimum distance from sentences of cluster. We apply this until we reach the number of sentences present in the reference summary

# Wikipedia Discovery
## -using different text summarization-
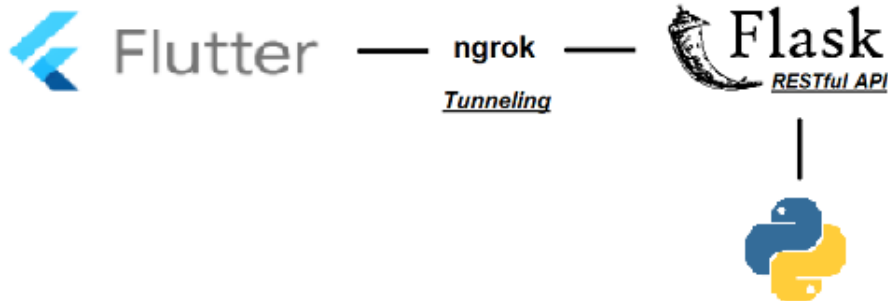### Extractive Summarization

**Application**



**NGROK**

**Tunneling**

**Backend (Python)**                    **Frontend (Flutter)**



**Frameworks** used:

- numpy~=1.20.1
- pandas~=1.2.4
- scikit-learn~=0.24.1
- scipy~=1.6.2
- **gensim**~=3.8.3 - an open-source library for unsupervised topic modeling and natural language processing
- **nltk**~=3.6.1 -a leading platform for building Python programs to work with human language data
- **textblob**~=0.15.3 -part-of-speech tagging, noun phrase extraction, sentiment analysis, classification, translation, and more.
- flask~=1.1.2
- **rouge**~=1.0.0
- Flutter (frontend)
- Python 3.9
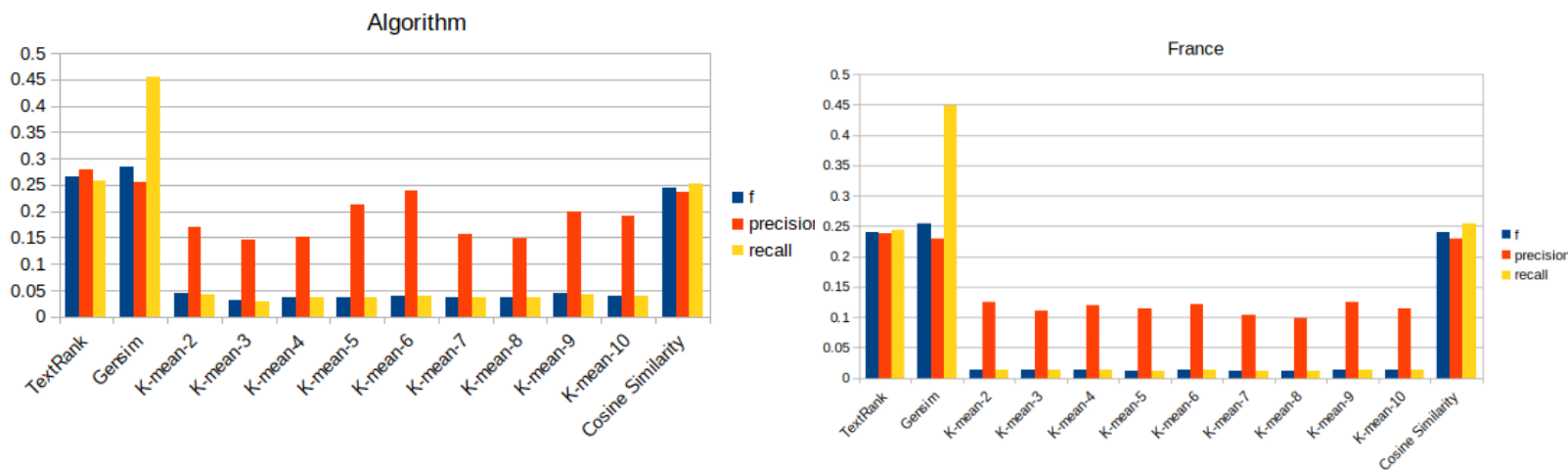- Ngrok.exe (tunneling on the web)

**Experiments and results**

In order to evalute generated summaries with reference summary, we need a relavant evaluation tool. How can we state in term of number how close is our summary to the reference one.

We decided to use ROUGE evaluation system [5]. ROUGE evaluation is a method to calculate the percentage of generate summary in reference summary and vice versa

$$\frac{\sum_r \sum_s \text{match}(\text{gram}_{s,c})}{\sum_r \sum_s \text{count}(\text{gram}_s)}$$

- ROUGE-n recall=40% means that 40% of the n-grams in the *reference* summary are also present in the *generated* summary.
- ROUGE-n precision=40% means that 40% of the n-grams in the *generated* summary are also present in the *reference* summary.
- ROUGE-n F1-score=40% is more difficult to interpret, like any F1-score.

Overall precision for all algorithms varies from 15% to 30%. It means that 15-30% of the n-grams in the generated summary are also present in the reference summary.

$$\text{RECALL} \qquad \frac{number\_of\_overlapping\_words}{total\_words\_in\_reference\_summary}$$

Overall recall gives us quite same values (about 30%) except Gensim which is way higher. It means that 30% of the n-grams in the reference summary are also present in the generated summary. However, K-mean implementation leads to a very low recall no matter how many clusters we choose.

$$\text{PRECISION} \qquad \frac{number\_of\_overlapping\_words}{total\_words\_in\_system\_summary}$$

Considering f score, it's way tougher to evaluate but its pretty well in agreement with precision and recall (except for K-Mean).

So why K-mean results are not relevant ?

We tried to figure out and tweak the algorithm but we guess that the approach we took was not optimal. It would have been way more relevant to combine K-Mean classification with other algorithm. We could have apply Text Rank Algorithm inside of each cluster.

**Conclusion**

Even K-mean can be in some way an approach for summarize text, our implementation didnt provide expected results. Although this basic implementation didnt do the trick, be believe that combined with other evaluation algorithms, it can be used in text summarization domain.

**References**

[1] Google Page Rank algorithm. https://en.wikipedia.org/wiki/PageRank

[2] K-means clustering, Wikipedia. https://en.wikipedia.org/wiki/K-means_clustering

[3] Automatic document summarization by sentence extraction. https://pdfs.semanticscholar.org/e7a4/8350000cec2025a212e7e3ca533b76351027.pdf

[4] Automatic extractive single document summarization, An unsupervised approach https://pdfs.semanticscholar.org/e7a4/8350000cec2025a212e7e3ca533b76351027.pdf

[5] ROUGE (Metric), Wikipedia. https://en.wikipedia.org/wiki/ROUGE_(metric)

[6] n-gram https://en.wikipedia.org/wiki/N-gram

[7] Lin, Chin-Yew. "Rouge: A package for automatic evaluation of summaries." http://anthology.aclweb.org/W/W04/W04-1013.pdf