VIET NAM NATIONAL UNIVERSITY HO CHI MINH CITY

UNIVERSITY OF SCIENCE

FACULTY OF INFORMATION TECHNOLOGY

---

# Assignment 3
# Text-based Search Engine

**Introduction to Information Retrieval - CS419**

---

21125148 - Mai Tan Dat

Ngày 20 tháng 11 năm 2024

# 1 Functionality

1. Text preprocessing using Underthesea library. Unlike English, Vietnamese is a monosyllabic language so stemming/lemmatization is unnecessary.

   - Text normalization: "oà"→ "òa", "uý"→ "úy", "đột quị"→ "đột quỵ".
   - Word tokenization: "Chàng trai 9X Quảng Trị khởi nghiệp từ nấm sò"→ "Chàng_trai 9X Quảng_Trị khởi_nghiệp từ nấm sò".
   - Stop words removal.

2. TF-IDF indexing using TfidfVectorizer from scikit-learn library.

3. Store the index in json file.

4. Graphical user interface using tkinter library, allow user to enter query string and search for top K relevant documents.

5. Rank documents by cosine similarity.
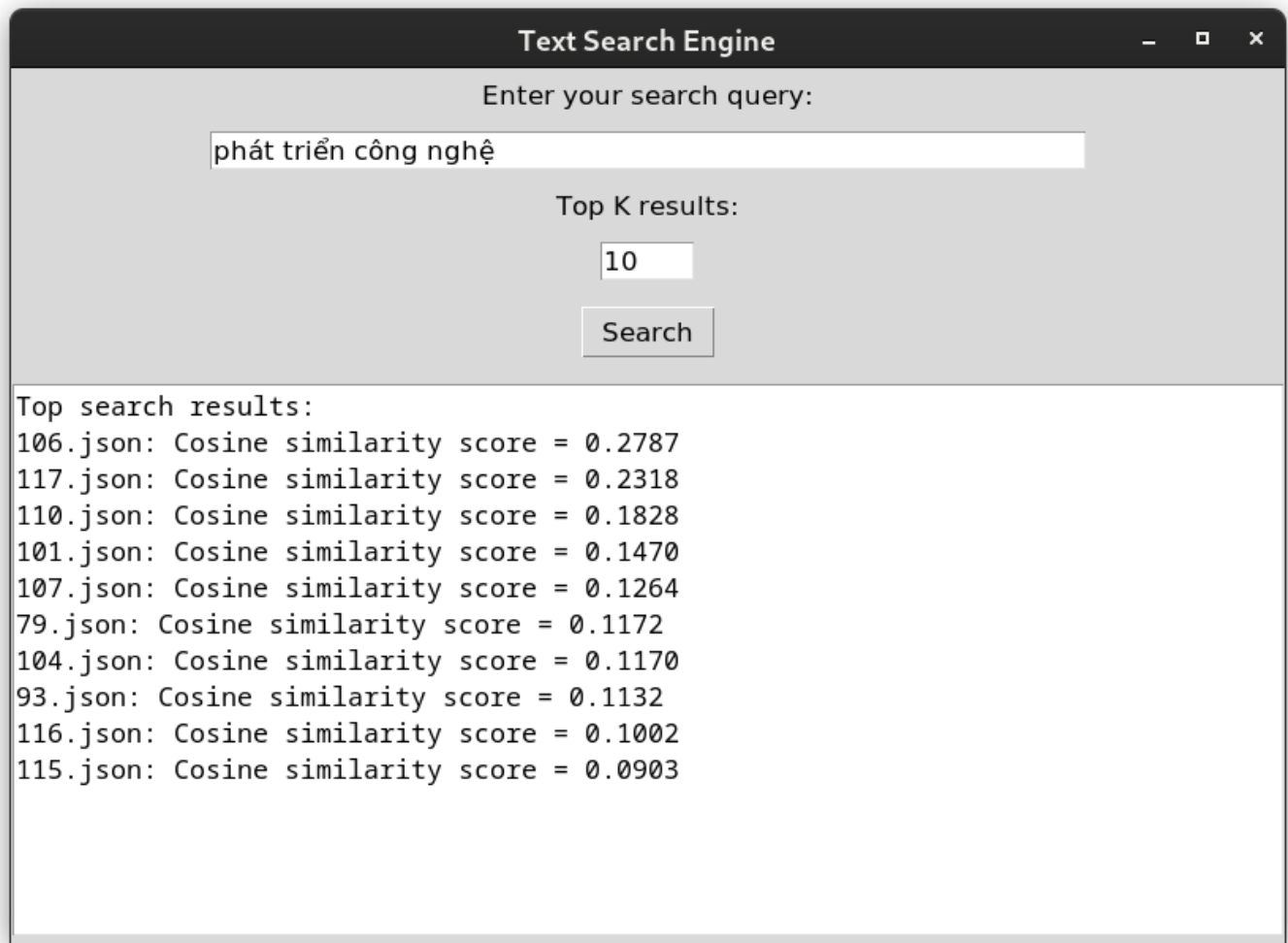
# 2 Usage instruction

Put document files in folder "./data".
A document file is a json file with format:

```
1  {
2      "Title": "This is sentence 1",
3      "Detail_sapo": "This is sentence 2",
4      "Content": [
5          "This is sentence 3",
6          "This is sentence 4",
7          "This is sentence 5",
8          ...
9      ],
10     ...
11 }
```

In the first run, the program load the data, create document index and save it to the folder "./vectorizer_data". Later, the program will load the saved index.
In Figure 1, when searching for "phát triển công nghệ", the result with the highest cosine similarity is ./data/106.json. In this document, the word "phát triển"appeared 11 times and the word "công nghệ"appeared 24 times.

Hình 1: Screenshot of the program