



INNOVATION. AUTOMATION. ANALYTICS

PROJECT ON

Enhancing Search Engine Relevance for Video Subtitles

By

Team ID

SAIPRANEETH S – IN1240164

HEMA B - IN1240288

Abstract:

We're using advanced techniques like understanding the meaning of words and how they fit together. Instead of just looking for specific words, we're trying to understand what the user wants and find subtitles that best match that. We have three main steps: getting the subtitles ready for analysis, turning the words into numbers so we can compare them, and then figuring out how similar the subtitles are to what the user asked for. By doing this, we hope to provide more helpful search results and make subtitles more accessible for everyone.

1. Reading the given data:

- The dataset comprises 82,498 subtitle files obtained from opensubtitles.org.
- These subtitles predominantly cover movies and TV series released between 1990 and 2024.
- The database file is named eng_subtitles_database.db.
- Within the database, there exists a table named 'zipfiles' containing three columns: 'num' (unique subtitle ID on opensubtitles.org), 'name' (subtitle file name), and 'content' (compressed subtitle data stored as binary using 'latin-1' encoding).
- Additional details about each subtitle can be accessed using the 'num' column in the URL: <https://www.opensubtitles.org/en/subtitles/{num}>, with {num} representing the unique subtitle ID.

```
In [1]: import sqlite3
import pandas as pd
```

Step 1 - Reading the Tables from Database file

```
In [2]: cd C:\Users\DELL\OneDrive\Desktop\data-20240407T052706Z-001\data
C:\Users\DELL\OneDrive\Desktop\data-20240407T052706Z-001\data
```

```
In [3]: # Read the code below and write your observation in the next cell

conn = sqlite3.connect("eng_subtitles_database.db")
cursor = conn.cursor()
cursor.execute("SELECT name FROM sqlite_master WHERE type='table'")
print(cursor.fetchall())

[('zipfiles',)]
```

1.1 Loading the database

This Python script demonstrates the process of fetching subtitle data from a SQLite database, decoding it, cleaning it, and storing it in a Pandas DataFrame.

- The script establishes a connection to the SQLite database containing subtitle data using the `sqlite3.connect()` function. The database path is specified as `database_path`.

- A SQL query is executed to select data from the "zipfiles" table within the SQLite database. The query retrieves columns named "num", "name", and "content", representing subtitle ID, name, and content respectively.

Step 3 - Loading the Database Table inside a Pandas DataFrame

```
In [5]: df = pd.read_sql_query("""SELECT * FROM zipfiles""", conn)
df.head()
```

```
Out[5]:
```

	num	name	content
0	9180533	the.message.(1976).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x1c\xa9\x...
1	9180583	here.comes.the.grump.s01.e09.john.jack.in.bo...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x17\xb9\x...
2	9180592	yumis.cells.s02.e13.episode.2.13.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00L\xb9\x99V...
3	9180594	yumis.cells.s02.e14.episode.2.14.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00U\xa9\x99V...
4	9180600	broker.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x001\xa9\x99V...

1.2 Unzipping the content and decoding using latin-1

The `extract_content` function processes binary content representing a zip archive. It employs an in-memory binary stream and iterates through the archive's files. Each file's content is decoded under the assumption that it's text data, using the "latin-1" encoding. The function returns the content of the first file encountered. This process is then applied to the 'subtitle_content' column of DataFrame `data`, enabling extraction of zip archive content stored within this column.

Applying for entire data

```
In [10]: import zipfile
import io

count = 0

def decode_method(binary_data):
    global count
    # Decompress the binary data using the zipfile module
    # print(count, end=" ")
    count += 1
    with io.BytesIO(binary_data) as f:
        with zipfile.ZipFile(f, 'r') as zip_file:
            # Assuming there's only one file in the ZIP archive
            subtitle_content = zip_file.read(zip_file.namelist()[0])

    # Now 'subtitle_content' should contain the extracted subtitle content
    return subtitle_content.decode('latin-1') # Assuming the content is UTF-8 encoded text
```

```
In [11]: df['file_content'] = df['content'].apply(decode_method)
df.head()
```

```
Out[11]:
```

	num	name	content	file_content
0	9180533	the.message.(1976).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x1c\xa9\x...	Trin00.00.06.000 --> 00.00.12.074r'nWatch an...
1	9180583	here.comes.the.grump.s01.e09.john.jack.in.bo...	b'PK\x03\x04\x14\x00\x00\x00\x08\x00\x17\xb9\x...	Trin00.00.29.358 --> 00.00.12.048r'nAht Ther...
2	9180592	yumis.cells.s02.e13.episode.2.13.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00L\xb9\x99V...	Trin00.00.53.200 --> 00.00.56.030r'nYumf...
3	9180594	yumis.cells.s02.e14.episode.2.14.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x00U\xa9\x99V...	Trin00.00.06.000 --> 00.00.12.074r'nWatch an...
4	9180600	broker.(2022).eng.1cd	b'PK\x03\x04\x14\x00\x00\x00\x08\x001\xa9\x99V...	Trin00.00.06.000 --> 00.00.12.074r'nWatch...

```
In [20]: data.head()
```

```
Out[20]:
```

	num	name	file_content
0	9180533	the.message.(1976)	1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch an...
1	9180583	here.comes.the.grump.s01.e09.joltin.jack.in.bo...	1\r\n00:00:29,359 --> 00:00:32,048\r\nAh! Ther...
2	9180592	yumis.cells.s02.e13.episode.2.13.(2022)	1\r\n00:00:53,200 --> 00:00:56,030\r\n<i>Yumi'...
3	9180594	yumis.cells.s02.e14.episode.2.14.(2022)	1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch an...
4	9180600	broker.(2022)	i»¿1\r\n00:00:06,000 --> 00:00:12,074\r\nWatch...

2.Data Preprocessing

The subtitle content within the 'File_content' column typically contains time stamps, HTML tags, and various forms of noise. It's imperative to clean this text data as part of the preprocessing for natural language processing (NLP) tasks. Cleaning involves eliminating irrelevant information, noise, and inconsistencies to ensure the text is better suited for subsequent processing steps like vectorization.

```
In [34]: import string
from nltk.corpus import stopwords
from nltk.tokenize import word_tokenize
from bs4 import BeautifulSoup
import unicodedata
from nltk.stem import WordNetLemmatizer
import re

def clean_text(sentence):
    # Remove timestamps
    clean_sentence = re.sub(r'\d+:\d+:\d+,?\d* --> \d+:\d+:\d+,?\d*', '', sentence)

    # Remove special characters and extra spaces
    clean_sentence = re.sub(r'[\u00A0-\u2019\u201C\u201D\u201E\u201F\u2020\u2021\u2022\u2023\u2024\u2025\u2026\u2027\u2028\u2029\u2030\u2031\u2032\u2033\u2034\u2035\u2036\u2037\u2038\u2039\u2040\u2041\u2042\u2043\u2044\u2045\u2046\u2047\u2048\u2049\u2050\u2051\u2052\u2053\u2054\u2055\u2056\u2057\u2058\u2059\u2060\u2061\u2062\u2063\u2064\u2065\u2066\u2067\u2068\u2069\u2070\u2071\u2072\u2073\u2074\u2075\u2076\u2077\u2078\u2079\u2080\u2081\u2082\u2083\u2084\u2085\u2086\u2087\u2088\u2089\u2090\u2091\u2092\u2093\u2094\u2095\u2096\u2097\u2098\u2099\u20A0\u20A1\u20A2\u20A3\u20A4\u20A5\u20A6\u20A7\u20A8\u20A9\u20B0\u20B1\u20B2\u20B3\u20B4\u20B5\u20B6\u20B7\u20B8\u20B9\u20C0\u20C1\u20C2\u20C3\u20C4\u20C5\u20C6\u20C7\u20C8\u20C9\u20D0\u20D1\u20D2\u20D3\u20D4\u20D5\u20D6\u20D7\u20D8\u20D9\u20E0\u20E1\u20E2\u20E3\u20E4\u20E5\u20E6\u20E7\u20E8\u20E9\u20F0\u20F1\u20F2\u20F3\u20F4\u20F5\u20F6\u20F7\u20F8\u20F9\u2100\u2101\u2102\u2103\u2104\u2105\u2106\u2107\u2108\u2109\u2110\u2111\u2112\u2113\u2114\u2115\u2116\u2117\u2118\u2119\u2120\u2121\u2122\u2123\u2124\u2125\u2126\u2127\u2128\u2129\u2130\u2131\u2132\u2133\u2134\u2135\u2136\u2137\u2138\u2139\u2140\u2141\u2142\u2143\u2144\u2145\u2146\u2147\u2148\u2149\u2150\u2151\u2152\u2153\u2154\u2155\u2156\u2157\u2158\u2159\u2160\u2161\u2162\u2163\u2164\u2165\u2166\u2167\u2168\u2169\u2170\u2171\u2172\u2173\u2174\u2175\u2176\u2177\u2178\u2179\u2180\u2181\u2182\u2183\u2184\u2185\u2186\u2187\u2188\u2189\u2190\u2191\u2192\u2193\u2194\u2195\u2196\u2197\u2198\u2199\u21A0\u21A1\u21A2\u21A3\u21A4\u21A5\u21A6\u21A7\u21A8\u21A9\u21B0\u21B1\u21B2\u21B3\u21B4\u21B5\u21B6\u21B7\u21B8\u21B9\u21C0\u21C1\u21C2\u21C3\u21C4\u21C5\u21C6\u21C7\u21C8\u21C9\u21D0\u21D1\u21D2\u21D3\u21D4\u21D5\u21D6\u21D7\u21D8\u21D9\u21E0\u21E1\u21E2\u21E3\u21E4\u21E5\u21E6\u21E7\u21E8\u21E9\u21F0\u21F1\u21F2\u21F3\u21F4\u21F5\u21F6\u21F7\u21F8\u21F9\u2200\u2201\u2202\u2203\u2204\u2205\u2206\u2207\u2208\u2209\u2210\u2211\u2212\u2213\u2214\u2215\u2216\u2217\u2218\u2219\u2220\u2221\u2222\u2223\u2224\u2225\u2226\u2227\u2228\u2229\u2230\u2231\u2232\u2233\u2234\u2235\u2236\u2237\u2238\u2239\u2240\u2241\u2242\u2243\u2244\u2245\u2246\u2247\u2248\u2249\u2250\u2251\u2252\u2253\u2254\u2255\u2256\u2257\u2258\u2259\u2260\u2261\u2262\u2263\u2264\u2265\u2266\u2267\u2268\u2269\u2270\u2271\u2272\u2273\u2274\u2275\u2276\u2277\u2278\u2279\u2280\u2281\u2282\u2283\u2284\u2285\u2286\u2287\u2288\u2289\u2290\u2291\u2292\u2293\u2294\u2295\u2296\u2297\u2298\u2299\u22A0\u22A1\u22A2\u22A3\u22A4\u22A5\u22A6\u22A7\u22A8\u22A9\u22B0\u22B1\u22B2\u22B3\u22B4\u22B5\u22B6\u22B7\u22B8\u22B9\u22C0\u22C1\u22C2\u22C3\u22C4\u22C5\u22C6\u22C7\u22C8\u22C9\u22D0\u22D1\u22D2\u22D3\u22D4\u22D5\u22D6\u22D7\u22D8\u22D9\u22E0\u22E1\u22E2\u22E3\u22E4\u22E5\u22E6\u22E7\u22E8\u22E9\u22F0\u22F1\u22F2\u22F3\u22F4\u22F5\u22F6\u22F7\u22F8\u22F9\u2300\u2301\u2302\u2303\u2304\u2305\u2306\u2307\u2308\u2309\u2310\u2311\u2312\u2313\u2314\u2315\u2316\u2317\u2318\u2319\u2320\u2321\u2322\u2323\u2324\u2325\u2326\u2327\u2328\u2329\u2330\u2331\u2332\u2333\u2334\u2335\u2336\u2337\u2338\u2339\u2340\u2341\u2342\u2343\u2344\u2345\u2346\u2347\u2348\u2349\u2350\u2351\u2352\u2353\u2354\u2355\u2356\u2357\u2358\u2359\u2360\u2361\u2362\u2363\u2364\u2365\u2366\u2367\u2368\u2369\u2370\u2371\u2372\u2373\u2374\u2375\u2376\u2377\u2378\u2379\u2380\u2381\u2382\u2383\u2384\u2385\u2386\u2387\u2388\u2389\u2390\u2391\u2392\u2393\u2394\u2395\u2396\u2397\u2398\u2399\u23A0\u23A1\u23A2\u23A3\u23A4\u23A5\u23A6\u23A7\u23A8\u23A9\u23B0\u23B1\u23B2\u23B3\u23B4\u23B5\u23B6\u23B7\u23B8\u23B9\u23C0\u23C1\u23C2\u23C3\u23C4\u23C5\u23C6\u23C7\u23C8\u23C9\u23D0\u23D1\u23D2\u23D3\u23D4\u23D5\u23D6\u23D7\u23D8\u23D9\u23E0\u23E1\u23E2\u23E3\u23E4\u23E5\u23E6\u23E7\u23E8\u23E9\u23F0\u23F1\u23F2\u23F3\u23F4\u23F5\u23F6\u23F7\u23F8\u23F9\u2400\u2401\u2402\u2403\u2404\u2405\u2406\u2407\u2408\u2409\u2410\u2411\u2412\u2413\u2414\u2415\u2416\u2417\u2418\u2419\u2420\u2421\u2422\u2423\u2424\u2425\u2426\u2427\u2428\u2429\u2430\u2431\u2432\u2433\u2434\u2435\u2436\u2437\u2438\u2439\u2440\u2441\u2442\u2443\u2444\u2445\u2446\u2447\u2448\u2449\u2450\u2451\u2452\u2453\u2454\u2455\u2456\u2457\u2458\u2459\u2460\u2461\u2462\u2463\u2464\u2465\u2466\u2467\u2468\u2469\u2470\u2471\u2472\u2473\u2474\u2475\u2476\u2477\u2478\u2479\u2480\u2481\u2482\u2483\u2484\u2485\u2486\u2487\u2488\u2489\u2490\u2491\u2492\u2493\u2494\u2495\u2496\u2497\u2498\u2499\u24A0\u24A1\u24A2\u24A3\u24A4\u24A5\u24A6\u24A7\u24A8\u24A9\u24B0\u24B1\u24B2\u24B3\u24B4\u24B5\u24B6\u24B7\u24B8\u24B9\u24C0\u24C1\u24C2\u24C3\u24C4\u24C5\u24C6\u24C7\u24C8\u24C9\u24D0\u24D1\u24D2\u24D3\u24D4\u24D5\u24D6\u24D7\u24D8\u24D9\u24E0\u24E1\u24E2\u24E3\u24E4\u24E5\u24E6\u24E7\u24E8\u24E9\u24F0\u24F1\u24F2\u24F3\u24F4\u24F5\u24F6\u24F7\u24F8\u24F9\u2500\u2501\u2502\u2503\u2504\u2505\u2506\u2507\u2508\u2509\u2510\u2511\u2512\u2513\u2514\u2515\u2516\u2517\u2518\u2519\u2520\u2521\u2522\u2523\u2524\u2525\u2526\u2527\u2528\u2529\u2530\u2531\u2532\u2533\u2534\u2535\u2536\u2537\u2538\u2539\u2540\u2541\u2542\u2543\u2544\u2545\u2546\u2547\u2548\u2549\u2550\u2551\u2552\u2553\u2554\u2555\u2556\u2557\u2558\u2559\u2560\u2561\u2562\u2563\u2564\u2565\u2566\u2567\u2568\u2569\u2570\u2571\u2572\u2573\u2574\u2575\u2576\u2577\u2578\u2579\u2580\u2581\u2582\u2583\u2584\u2585\u2586\u2587\u2588\u2589\u2590\u2591\u2592\u2593\u2594\u2595\u2596\u2597\u2598\u2599\u25A0\u25A1\u25A2\u25A3\u25A4\u25A5\u25A6\u25A7\u25A8\u25A9\u25B0\u25B1\u25B2\u25B3\u25B4\u25B5\u25B6\u25B7\u25B8\u25B9\u25C0\u25C1\u25C2\u25C3\u25C4\u25C5\u25C6\u25C7\u25C8\u25C9\u25D0\u25D1\u25D2\u25D3\u25D4\u25D5\u25D6\u25D7\u25D8\u25D9\u25E0\u25E1\u25E2\u25E3\u25E4\u25E5\u25E6\u25E7\u25E8\u25E9\u25F0\u25F1\u25F2\u25F3\u25F4\u25F5\u25F6\u25F7\u25F8\u25F9\u2600\u2601\u2602\u2603\u2604\u2605\u2606\u2607\u2608\u2609\u2610\u2611\u2612\u2613\u2614\u2615\u2616\u2617\u2618\u2619\u2620\u2621\u2622\u2623\u2624\u2625\u2626\u2627\u2628\u2629\u2630\u2631\u2632\u2633\u2634\u2635\u2636\u2637\u2638\u2639\u2640\u2641\u2642\u2643\u2644\u2645\u2646\u2647\u2648\u2649\u2650\u2651\u2652\u2653\u2654\u2655\u2656\u2657\u2658\u2659\u2660\u2661\u2662\u2663\u2664\u2665\u2666\u2667\u2668\u2669\u2670\u2671\u2672\u2673\u2674\u2675\u2676\u2677\u2678\u2679\u2680\u2681\u2682\u2683\u2684\u2685\u2686\u2687\u2688\u2689\u2690\u2691\u2692\u2693\u2694\u2695\u2696\u2697\u2698\u2699\u26A0\u26A1\u26A2\u26A3\u26A4\u26A5\u26A6\u26A7\u26A8\u26A9\u26B0\u26B1\u26B2\u26B3\u26B4\u26B5\u26B6\u26B7\u26B8\u26B9\u26C0\u26C1\u26C2\u26C3\u26C4\u26C5\u26C6\u26C7\u26C8\u26C9\u26D0\u26D1\u26D2\u26D3\u26D4\u26D5\u26D6\u26D7\u26D8\u26D9\u26E0\u26E1\u26E2\u26E3\u26E4\u26E5\u26E6\u26E7\u26E8\u26E9\u26F0\u26F1\u26F2\u26F3\u26F4\u26F5\u26F6\u26F7\u26F8\u26F9\u2700\u2701\u2702\u2703\u2704\u2705\u2706\u2707\u2708\u2709\u2710\u2711\u2712\u2713\u2714\u2715\u2716\u2717\u2718\u2719\u2720\u2721\u2722\u2723\u2724\u2725\u2726\u2727\u2728\u2729\u2730\u2731\u2732\u2733\u2734\u2735\u2736\u2737\u2738\u2739\u2740\u2741\u2742\u2743\u2744\u2745\u2746\u2747\u2748\u2749\u2750\u2751\u2752\u2753\u2754\u2755\u2756\u2757\u2758\u2759\u2760\u2761\u2762\u2763\u2764\u2765\u2766\u2767\u2768\u2769\u2770\u2771\u2772\u2773\u2774\u2775\u2776\u2777\u2778\u2779\u2780\u2781\u2782\u2783\u2784\u2785\u2786\u2787\u2788\u2789\u2790\u2791\u2792\u2793\u2794\u2795\u2796\u2797\u2798\u2799\u27A0\u27A1\u27A2\u27A3\u27A4\u27A5\u27A6\u27A7\u27A8\u27A9\u27B0\u27B1\u27B2\u27B3\u27B4\u27B5\u27B6\u27B7\u27B8\u27B9\u27C0\u27C1\u27C2\u27C3\u27C4\u27C5\u27C6\u27C7\u27C8\u27C9\u27D0\u27D1\u27D2\u27D3\u27D4\u27D5\u27D6\u27D7\u27D8\u27D9\u27E0\u27E1\u27E2\u27E3\u27E4\u27E5\u27E6\u27E7\u27E8\u27E9\u27F0\u27F1\u27F2\u27F3\u27F4\u27F5\u27F6\u27F7\u27F8\u27F9\u2800\u2801\u2802\u2803\u2804\u2805\u2806\u2807\u2808\u2809\u2810\u2811\u2812\u2813\u2814\u2815\u2816\u2817\u2818\u2819\u2820\u2821\u2822\u2823\u2824\u2825\u2826\u2827\u2828\u2829\u2830\u2831\u2832\u2833\u2834\u2835\u2836\u2837\u2838\u2839\u2840\u2841\u2842\u2843\u2844\u2845\u2846\u2847\u2848\u2849\u2850\u2851\u2852\u2853\u2854\u2855\u2856\u2857\u2858\u2859\u2860\u2861\u2862\u2863\u2864\u2865\u2866\u2867\u2868\u2869\u2870\u2871\u2872\u2873\u2874\u2875\u2876\u2877\u2878\u2879\u2880\u2881\u2882\u2883\u2884\u2885\u2886\u2887\u2888\u2889\u2890\u2891\u2892\u2893\u2894\u2895\u2896\u2897\u2898\u2899\u28A0\u28A1\u28A2\u28A3\u28A4\u28A5\u28A6\u28A7\u28A8\u28A9\u28B0\u28B1\u28B2\u28B3\u28B4\u28B5\u28B6\u28B7\u28B8\u28B9\u28C0\u28C1\u28C2\u28C3\u28C4\u28C5\u28C6\u28C7\u28C8\u28C9\u28D0\u28D1\u28D2\u28D3\u28D4\u28D5\u28D6\u28D7\u28D8\u28D9\u28E0\u28E1\u28E2\u28E3\u28E4\u28E5\u28E6\u28E7\u28E8\u28E9\u28F0\u28F1\u28F2\u28F3\u28F4\u28F5\u28F6\u28F7\u28F8\u28F9\u2900\u2901\u2902\u2903\u2904\u2905\u2906\u2907\u2908\u2909\u2910\u2911\u2912\u2913\u2914\u2915\u2916\u2917\u2918\u2919\u2920\u2921\u2922\u2923\u2924\u2925\u2926\u2927\u2928\u2929\u2930\u2931\u2932\u2933\u2934\u2935\u2936\u2937\u2938\u2939\u2940\u2941\u2942\u2943\u2944\u2945\u2946\u2947\u2948\u2949\u2950\u2951\u2952\u2953\u2954\u2955\u2956\u2957\u2958\u2959\u2960\u2961\u2962\u2963\u2964\u2965\u2966\u2967\u2968\u2969\u2970\u2971\u2972\u2973\u2974\u2975\u2976\u2977\u2978\u2979\u2980\u2981\u2982\u2983\u2984\u2985\u2986\u2987\u2988\u2989\u2990\u2991\u2992\u2993\u2994\u2995\u2996\u2997\u2998\u2999\u29A0\u29A1\u29A2\u29A3\u29A4\u29A5\u29A6\u29A7\u29A8\u29A9\u29B0\u29B1\u29B2\u29B3\u29B4\u29B5\u29B6\u29B7\u29B8\u29B9\u29C0\u29C1\u29C2\u29C3\u29C4\u29C5\u29C6\u29C7\u29C8\u29C9\u29D0\u29D1\u29D2\u29D3\u29D4\u29D5\u29D6\u29D7\u29D8\u29D9\u29E0\u29E1\u29E2\u29E3\u29E4\u29E5\u29E6\u29E7\u29E8\u29E9\u29F0\u29F1\u29F2\u29F3\u29F4\u29F5\u29F6\u29F7\u29F8\u29F9\u2A00\u2A01\u2A02\u2A03\u2A04\u2A05\u2A06\u2A07\u2A08\u2A09\u2A10\u2A11\u2A12\u2A13\u2A14\u2A15\u2A16\u2A17\u2A18\u2A19\u2A20\u2A21\u2A22\u2A23\u2A24\u2A25\u2A26\u2A27\u2A28\u2A29\u2A30\u2A31\u2A32\u2A33\u2A34\u2A35\u2A36\u2A37\u2A38\u2A39\u2A40\u2A41\u2A42\u2A43\u2A44\u2A45\u2A46\u2A47\u2A48\u2A49\u2A50\u2A51\u2A52\u2A53\u2A54\u2A55\u2A56\u2A57\u2A58\u2A59\u2A60\u2A61\u2A62\u2A63\u2A64\u2A65\u2A66\u2A67\u2A68\u2A69\u2A70\u2A71\u2A72\u2A73\u2A74\u2A75\u2A76\u2A77\u2A78\u2A79\u2A80\u2A81\u2A82\u2A83\u2A84\u2A85\u2A86\u2A87\u2A88\u2A89\u2A90\u2A91\u2A92\u2A93\u2A94\u2A95\u2A96\u2A97\u2A98\u2A99\u2AA0\u2AA1\u2AA2\u2AA3\u2AA4\u2AA5\u2AA6\u2AA7\u2AA8\u2AA9\u2AB0\u2AB1\u2AB2\u2AB3\u2AB4\u2AB5\u2AB6\u2AB7\u2AB8\u2AB9\u2AC0\u2AC1\u2AC2\u2AC3\u2AC4\u2AC5\u2AC6\u2AC7\u2AC8\u2AC9\u2AD0\u2AD1\u2AD2\u2AD3\u2AD4\u2AD5\u2AD6\u2AD7\u2AD8\u2AD9\u2AE0\u2AE1\u2AE2\u2AE3\u2AE4\u2AE5\u2AE6\u2AE7\u2AE8\u2AE9\u2AF0\u2AF1\u2AF2\u2AF3\u2AF4\u2AF5\u2AF6\u2AF7\u2AF8\u2AF9\u2B00\u2B01\u2B02\u2B03\u2B04\u2B05\u2B06\u2B07\u2B08\u2B09\u2B10\u2B11\u2B12\u2B13\u2B14\u2B15\u2B16\u2B17\u2B18\u2B19\u2B20\u2B21\u2B22\u2B23\u2B24\u2B25\u2B26\u2B27\u2B28\u2B29\u2B30\u2B31\u2B32\u2B33\u2B34\u2B35\u2B36\u2B37\u2B38\u2B39\u2B40\u2B41\u2B42\u2B43\u2B44\u2B45\u2B46\u2B47\u2B48\u2B49\u2B50\u2B51\u2B52\u2B53\u2B54\u2B55\u2B56\u2B57\u2B58\u2B59\u2B60\u2B61\u2B62\u2B63\u2B64\u2B65\u2B66\u2B67\u2B68\u2B69\u2B70\u2B71\u2B72\u2B73\u2B74\u2B75\u2B76\u2B77\u2B78\u2B79\u2B80\u2B81\u2B82\u2B83\u2B84\u2B85\u2B86\u2B87\u2B88\u2B89\u2B90\u2B91\u2B92\u2B93\u2B94\u2B95\u2B96\u2B97\u2B98\u2B99\u2BA0\u2BA1\u2BA2\u2BA3\u2BA4\u2BA5\u2BA6\u2BA7\u2BA8\u2BA9\u2BB0\u2BB1\u2BB2\u2BB3\u2BB4\u2BB5\u2BB6\u2BB7\u2BB8\u2BB9\u2BC0\u2BC1\u2BC2\u2BC3\u2BC4\u2BC5\u2BC6\u2BC7\u2BC8\u2BC9\u2BD0\u2BD1\u2BD2\u2BD3\u2BD4\u2BD5\u2BD6\u2BD7\u2BD8\u2BD9\u2BE0\u2BE1\u2BE2\u2BE3\u2BE4\u2BE5\u2BE6\u2BE7\u2BE8\u2BE9\u2BF0\u2BF1\u2BF2\u2BF3\u2BF4\u2BF5\u2BF6\u2BF7\u2BF8\u2BF9\u2C00\u2C01\u2C02\u2C03\u2C04\u2C05\u2C06\u2C07\u2C08\u2C09\u2C10\u2C11\u2C12\u2C13\u2C14\u2C15\u2C16\u2C17\u2C18\u2C19\u2C20\u2C21\u2C22\u2C23\u2C24\u2C25\u2C26\u2C27\u2C28\u2C29\u2C30\u2C31\u2C32\u2C33\u2C34\u2C35\u2C36\u2C37\u2C38\u2C39\u2C40\u2C41\u2C42\u2C43\u2C44\u2C45\u2C46\u2C47\u2C48\u2C49\u2C50\u2C51\u2C52\u2C53\u2C54\u2C55\u2C56\u2C57\u2C58\u2C59\u2C60\u2C61\u2C62\u2C63\u2C64\u2C65\u2C66\u2C67\u2C68\u2C69\u2C70\u2C71\u2C72\u2C73\u2C74\u2C75\u2C76\u2C77\u2C78\u2C79\u2C80\u2C81\u2C82\u2C83\u2C84\u2C85\u2C86\u2C87\u2C88\u2C89\u2C90\u2C91\u2C92\u2C93\u2C94\u2C95\u2C96\u2C97\u2C98\u2C99\u2CA0\u2CA1\u2CA2\u2CA3\u2CA4\u2CA5\u2CA6\u2CA7\u2CA8\u2CA9\u2CB0\u2CB1\u2CB2\u2CB3\u2CB4\u2CB5\u2CB6\u2CB7\u2CB8\u2CB9\u2CC0\u2CC1\u2CC2\u2CC3\u2CC4\u2CC5\u2CC6\u2CC7\u2CC8\u2CC9\u2CD0\u2CD1\u2CD2\u2CD3\u2CD4\u2CD5\u2CD6\u2CD7\u2CD8\u2CD9\u2CE0\u2CE1\u2CE2\u2CE3\u2CE4\u2CE5\u2CE6\u2CE7\u2CE8\u2CE9\u2CF0\u2CF1\u2CF2\u2CF3\u2CF4\u2CF5\u2CF6\u2CF7\u2CF8\u2CF9\u2D00\u2D01\u2D02\u2D03\u2D04\u2D05\u2D06\u2D07\u2D08\u2D09\u2D10\u2D11\u2D12\u2D13\u2D14\u2D15\u2D16\u2D17\u2D18\u2D19\u2D20\u2D21\u2D22\u2D23\u2D24\u2D25\u2D26\u2D27\u2D28\u2D29\u2D30\u2D31\u2D32\u2D33\u2D34\u2D35\u2D36\u2D37\u2D38\u2D39\u2D40\u2D41\u2D42\u2D43\u2D44\u2D45\u2D46\u2D47\u2D48\u2D49\u2D50\u2D51\u2D52\u2D53\u2D54\u2D55\u2D56\u2D57\u2D58\u2D59\u2D60\u2D61\u2D62\u2D63\u2D64\u2D65\u2D66\u2D67\u2D68\u2D69\u2D70\u2D71\u2D72\u2D73\u2D74\u2D75\u2D76\u2D77\u2D78\u2D79\u2D80\u2D81\u2D82\u2D83\u2D84\u2D85\u2D86\u2D87\u2D88\u2D89\u2D90\u2D91\u2D92\u2D93\u2D94\u2D95\u2D96\u2D97\u2D98\u2D99\u2DA0\u2DA1\u2DA2\u2DA3\u2DA4\u2DA5\u2DA6\u2DA7\u2DA8\u2DA9\u2DB0\u2DB1\u2DB2\u2DB3\u2DB4\u2DB5\u2DB6\u2DB7\u2DB8\u2DB9\u2DC0\u2DC1\u2DC2\u2DC3\u2DC4\u2DC5\u2DC6\u2DC7\u2DC8\u2DC9\u2DD0\u2DD1\u2DD2\u2DD3\u2DD4\u2DD5\u2DD6\u2DD7\u2DD8\u2DD9\u2DE0\u2DE1\u2DE2\u2DE3\u2DE4\u2DE5\u2DE6\u2DE7\u2DE8\u2DE9\u2DF0\u2DF1\u2DF2\u2DF3\u2DF4\u2DF5\u2DF6\u2DF7\u2DF8\u2DF9\u2E00\u2E01\u2E02\u2E03\u2E04\u2E05\u2E06\u2E07\u2E08\u2E09\u2E10\u2E11\u2E12\u2E13\u2E14\u2E15\u2E16\u2E17\u2E18\u2E19\u2E20\u2E21\u2E22\u2E23\u2E24\u2E25\u2E26\u2E27\u2E28\u2E29\u2E30\u2E31\u2E32\u2E33\u2E34\u2E35\u2E36\u2E37\u2E38\u2E39\u2E40\u2E41\u2E42\u2E43\u2E44\u2E45\u2E46\u2E47\u2E48\u2E49\u2E50\u2E51\u2E52\u2E53\u2E54\u2E55\u2E56\u2E57\u2E58\u2E59\u2E60\u2E61\u2E62\u2E63\u2E64\u2E65\u2E66\u2E67\u2E68\u2E69\u2E70\u2E71\u2E72\u2E73\u2E74\u2E75\u2E76\u2E77\u2E78\u2E79\u2E80\u2E81\u2E82\u2E83\u2E84\u2E85\u2E86\u2E87\u2E88\u2E89\u2E90\u2E91\u2E92\u2E93\u2E94\u2E95\u2E96\u2E97\u2E98\u2E99\u2EA0\u2EA1\u2EA2\u2EA3\u2EA4\u2EA5\u2EA6\u2EA7\u2EA8\u2EA9\u2EB0\u2EB1\u2EB2\u2EB3\u2EB4\u2EB5\u2EB6\u2EB7\u2EB8\u2EB9\u2EC0\u2EC1\u2EC2\u2EC3\u2EC4\u2EC5\u2EC6\u2EC7\u2EC8\u2EC9\u2ED0\u2ED1\u2ED2\u2ED3\u2ED4\u2ED5\u2ED6\u2ED7\u2ED8\u2ED9\u2EE0\u2EE1\u2EE2\u
```

Remove Special Characters and Extra Spaces: It removes any characters that are not alphanumeric or whitespace using another regular expression.

Convert Text to Lowercase: All text is converted to lowercase to ensure consistency.

Remove Leading and Trailing Whitespace: It strips any leading or trailing whitespace.

Removing HTML Tags: BeautifulSoup is used to remove any HTML tags present in the text.

Removing URLs: It removes any words that start with 'http', assuming they are URLs.

Removing Punctuation: All punctuation marks are removed using a list comprehension.

Removing Numbers: Any digits are removed from the text.

Removing Stopwords: NLTK's stopwords list for English is used to remove common stopwords like 'the', 'and', 'is', etc.

Handling Special Characters: Special characters are normalized using Unicode normalization.

Lemmatization: Words are lemmatized using NLTK's WordNet lemmatizer, reducing them to their base or dictionary form.

Finally, the cleaned text is returned. This process helps in removing noise, irrelevant information, and inconsistencies from the text, making it more suitable for NLP tasks like text classification or sentiment analysis.

```
In [28]: data
```

```
Out[28]:
```

	num	name	file_content	clean_file_content
0	9180533	the message (1976)	1'in00:00:06,000 -> 00:00:12,074'inWatch an...	watch video online opensubtitles free browser ...
1	9180583	here comes the grump s01 e09 joltin jack in bo...	1'in00:00:29,368 -> 00:00:32,048'inAhl Ther...	ah there princess dawn terry blooney looney so...
2	9180592	yumis cells s02 e13 episode 2.13 (2022)	1'in00:00:53,200 -> 00:00:58,030'in-i>Yum?	yumis cell iepisode extremely polite yumy iy...
3	9180594	yumis cells s02 e14 episode 2.14 (2022)	1'in00:00:06,000 -> 00:00:12,074'inWatch an...	watch video online opensubtitles free browser ...
4	9180600	broker (2022)	ixg1'in00:00:06,000 -> 00:00:12,074'inWatch...	watch video online opensubtitles free browser ...
...
82493	9521935	the prophets game (2000)	ixg1'in00:01:16,284 -> 00:01:19,537'inGod,t...	god punishing red head us gun green chest us s...
82494	9521937	west beirut (1968)	1'in00:00:06,000 -> 00:00:12,074'inapi.Open...	apiopensubtitlesorg deprecated please implemen...
82495	9521938	frankenstein the true story (1973)	1'in00:00:01,001 -> 00:00:04,530'in(Dramati...	dramatic orchestral music advertise product br...
82496	9521940	frankenstein the true story (1973)	1'in00:00:06,000 -> 00:00:12,074'inAdvertis...	advertise product brand contact wwwopensubtitl...
82497	9521941	zombie island massacre (1984)	1'in00:00:01,919 -> 00:00:03,253'in(Sharp w...	sharp whistling man hey wait wait growling cla...

82498 rows x 4 columns

2.1 Converting to CSV:

Export DataFrame to CSV: It exports the DataFrame to a CSV file named "Data.csv", excluding the index column, and using '\' as the escape character.

Cleaning Text Data: It applies the clean_text function to the 'file_content' column of the DataFrame, storing the cleaned content in a new column named 'clean_file_content'.

Dropping Original Column: It drops the original 'file_content' column from the DataFrame.

```

In [32]: import pandas as pd
import os

# Example directory path
directory = 'D:\internship\Data science\TASKS\search engine\data'

# Check if the directory has write permission
if not os.access(directory, os.W_OK):
    print(f"No write permission in directory: {directory}")
    # You may choose to exit the script or handle this situation differently
else:
    # Assuming 'data' is your DataFrame
    try:
        data.to_csv(os.path.join(directory, 'eng_subtitles_database.csv'), index=False)
        print("Data successfully saved to CSV.")
    except PermissionError as e:
        print(f"PermissionError: {e}")

```

Data successfully saved to CSV.

3.Vectorize the given Subtitle Documents

Bag-of-Words (BOW) and Term Frequency-Inverse Document Frequency (TF-IDF) are traditional methods for generating sparse vector representations of text data.

```

In [13]: import pandas as pd
import numpy as np

In [14]: data=pd.read_csv("D:\internship\Data science\TASKS\search engine\data\eng_subtitles_database.csv")

In [15]: data.shape
Out[15]: (82498, 3)

In [16]: data.head()
Out[16]:

```

	num	name	clean_file_content
0	9180533	the.message (1976)	watch video online opensubtitles free browser ...
1	9180583	here.comes.the.grump.s01.e09.joltin.jack.in.bo...	ah there princess dawn terry blooney looney so...
2	9180592	yumis.cells.s02.e13.episode.2.13.(2022)	ryumis cell iepisode extremely polite yumii ry...
3	9180594	yumis.cells.s02.e14.episode.2.14.(2022)	watch video online opensubtitles free browser ...
4	9180600	broker.(2022)	watch video online opensubtitles free browser ...

3.1 Advantages:

Semantic Information: BERT-based Sentence Transformers capture semantic links between words and sentences, enabling a more nuanced understanding of the text.

Contextual Embeddings: Unlike BOW and TF-IDF, which treat words independently, Sentence Transformers generate dense embeddings that consider the context of words. This preserves contextual details and enhances representation quality.

Lower Dimensionality: Dense embeddings from Sentence Transformers typically have lower dimensionality than sparse representations, reducing computational overhead and memory usage.

Pre-trained Models: Sentence Transformers leverage pre-trained models like BERT, benefiting from extensive training on vast text corpora, thus capturing diverse linguistic patterns and semantic connections.

```
In [17]: data['clean_file_content'][0]
Out[17]: 'watch video online opensubtitles free browser extension osdblinkext name god gracious merciful muhammad messenger god heracl
ius emperor byzantium greeting follower righteous guidance bid hear divine call messenger god people accept islam salvation s
peaks new prophet arabia like john baptist came king herod desert cry salvation muqawqis patriarch alexandria kisra emperor p
ersia muhammad call call god accept islam salvation embrace islam come desert smelling camel goat tell persia kneel muhammad
messenger god gave authority god sent muhammad mercy mankind scholar historian islam university alazhar cairo high islamic co
ngress shiat lebanon maker film honour islamic tradition hold impersonation prophet offends spirituality message therefore pe
rson mohammad shown year christ died i when europe sunk dark agesi i and everywhere old civilization fall ing i muhammad born m
ecca arabia i mecca rich trading city ruled merchants i whose wealth multiplied unique privilege i they housed gods i every y
ear time great fair i the desert priest brought idols i and image god custody kaaba i once holy shrine abraham i the kaaba be
come house idolotry i hosting fewer different gods i mecca adi bilal today count unaya yet year god gold put god prophet toge
ther sit pretty hnn god place kaaba caravan syria hnn they must running theyll thirsty put five men north well many sheep shal
l kill give hundred mecca must keep name hospitality ten lamb leader bread water poet hakim house verse prose nightly put sla
ughter andand bread swear thinner water ohopen space open space lover poetry abu sufian willing patron art abu sofyam invite
poet joy kit love kin wine cake abound skill abu sofyam revel song begin abu sofyam invite poet silkworm china lady pleasure
limb lady see ravish eye yes length dinar abu sofyans wife oh gold god kaaba need upkeep man stood looked soul carry away mus
t muhammad come dont stop nephew maybe change change year old unnatural rich wife could afford best mecca yet chooses sit shi
vering cave unnatural man dare risk anger aluzza keep health manat god prosperity allat god family tribe hubal hubal start ca
ravan predicts fate challenge god within earshot god dangerous unreasonable rebellious blasphemous yes im afraid muhammad har
m always sad great fair might see next one abu talib abu talib catch breath zaid muhammad come mount hira yet he three day ha
```

BOW and TF-IDF have drawbacks:

Lack of Semantic Information: They fail to capture semantic connections among words or documents. Each word is treated in isolation, disregarding its context. Consequently, they may not fully grasp the intended meaning of the text.

High Dimensionality: BOW and TF-IDF yield high-dimensional sparse vectors, particularly with extensive vocabularies or datasets containing numerous unique tokens. This can lead to computational inefficiencies and increased memory usage.

Given these drawbacks, utilizing BOW or TF-IDF may not be optimal for tasks requiring semantic understanding or contextual comprehension. For endeavors like constructing a Semantic Search Engine, employing methods that encode semantic knowledge and contextual associations between words and documents is more advantageous.

4. Vectorizer

```
In [23]: count_vectorizer = CountVectorizer()
         tf_matrix = count_vectorizer.fit_transform(data['clean_file_content'])

In [24]: count_vectorizer
Out[24]: CountVectorizer()
In a Jupyter environment, please rerun this cell to show the HTML representation or trust the notebook.
On GitHub, the HTML representation is unable to render, please try loading this page with nbviewer.org.

In [25]: # TF-IDF transformer and transform the TF matrix
```

We use 'CountVectorizer' to convert text data into a matrix of token counts. This process transforms text into a numerical format suitable for machine learning algorithms, capturing the frequency of each word in the text. It's valuable for tasks like text classification or clustering, providing a foundational representation of text data for further analysis.

TF-IDF transformer

The 'TfidfTransformer' is used to convert a term frequency (TF) matrix into a TF-IDF matrix. This transformation normalizes TF values and applies inverse document frequency (IDF) weighting to emphasize rare terms and downweight common ones. It improves document representations by considering both local and global term characteristics, leading to more informative and discriminative document representations.

Calculating similarity using cosine similarity

To calculate similarity using cosine similarity, we compute the cosine of the angle between two vectors in a multidimensional space. In the context of text data, each document or text snippet is represented as a vector, and the similarity between two documents is determined by the cosine of the angle between their respective vectors. A cosine similarity of 1 indicates perfect similarity, while 0 indicates no similarity.

This measure is commonly used in information retrieval, recommendation systems, and clustering to assess the similarity between documents or text snippets based on their content. It's particularly useful when dealing with high-dimensional data, such as text data represented by TF-IDF or word embeddings.

```
In [29]: # Calculating similarity using cosine similarity
```

```
In [30]: query = input()
query_vector = count_vectorizer.transform([query])
query_tfidf = tfidf_transformer.transform(query_vector)
```

avatar

```
In [32]: similarity_scores = cosine_similarity(query_tfidf, tfidf_matrix)
```


5.Retrieving Documents

Retrieving Similar Documents: The code sorts the similarity scores in descending order to identify the top matching documents. It retrieves the content, subtitle IDs, and names of these documents.

Summarizing Documents: It iterates over the retrieved documents and generates summaries based on a given query. The function `generate_summarized_documents()` is called to create summaries for each document.

Displaying Results: For each retrieved document, it prints the summary, subtitle ID, and subtitle name (if available).

```
In [43]: top_indices = similarity_scores.argsort()[0][::-1]
top_n = 5
retrieved_documents = [data['clean_file_content'][idx] for idx in top_indices[:top_n]]
retrieved_subtitle_ids = [data['num'][idx] for idx in top_indices[:top_n]] # Assuming you h
retrieved_subtitle_names = [data['name'][idx] for idx in top_indices[:top_n]] # Assuming yo
4

In [46]: def generate_summarized_documents(query, document):
# Your implementation for summarizing the document based on the query
pass

In [47]: summarized_docs = []

# Iterate over each retrieved document and generate summaries
for doc in retrieved_documents:
    summarized_doc = generate_summarized_documents(query, doc)
    summarized_docs.append(summarized_doc)
```

This process allows users to input a query, find documents similar to the query, and display summarized information about those documents. It's useful for tasks such as document retrieval and summarization, aiding in information organization and understanding.

```

In [51]: # Assuming you have a list containing subtitle names called 'retrieved_subtitle_names'
for i, (summary, subtitle_id) in enumerate(zip(summarized_docs, retrieved_subtitle_ids), 1):
    print("Document", i, ":")
    print("Summary:", summary)
    print("Subtitle ID:", subtitle_id)

    if i <= len(retrieved_subtitle_names):
        print("Subtitle Name:", retrieved_subtitle_names[i - 1])
    else:
        print("Subtitle Name: Not available") # handle case where subtitle name list is shorter t

    print()

Document 1 :
Summary: None
Subtitle ID: 9233592
Subtitle Name: archer.s13.e03.saturday.(2022)

Document 2 :
Summary: None
Subtitle ID: 9233396
Subtitle Name: archer.s13.e03.saturday.(2022)

Document 3 :
Summary: None
Subtitle ID: 9233399
Subtitle Name: archer.s13.e03.saturday.(2022)

Document 4 :
Summary: None
Subtitle ID: 9233593
Subtitle Name: archer.s13.e03.saturday.(2022)

Document 5 :
Summary: None

```

This organizes the summarized documents and their associated subtitle IDs into a structured format for display. It ensures clarity and readability by presenting each document's summary alongside its unique identifier and, if applicable, its name. This approach facilitates efficient understanding and interpretation of the retrieved documents and their metadata.

```

In [52]: import joblib

# CountVectorizer
joblib.dump(count_vectorizer, 'count_vectorizer.joblib')

```

```

Out[52]: ['count_vectorizer.joblib']

```

```

In [53]: # TfidfTransformer
joblib.dump(tfidf_transformer, 'tfidf_transformer.joblib')

```

```

Out[53]: ['tfidf_transformer.joblib']

```

```

In [54]: joblib.dump(tfidf_matrix, 'tfidf_matrix.joblib')

```

```

Out[54]: ['tfidf_matrix.joblib']

```

```

In [55]: # cosine similarity model
joblib.dump(similarity_scores, 'cosine_similarity_scores.joblib')

```

```

Out[55]: ['cosine_similarity_scores.joblib']

```

We use this to save various components of our text processing and similarity modeling pipeline:

CountVectorizer: It's saved to preserve the mapping between words and their indices, allowing consistent vectorization of new data.

TfidfTransformer: Similarly, it's stored to retain information about the TF-IDF transformation applied to the data.

TF-IDF Matrix: The TF-IDF matrix itself is saved to avoid recomputation, ensuring consistency in future analysis.

Cosine Similarity Model: The cosine similarity scores calculated between documents are saved, enabling quick retrieval and comparison of document similarities without needing to recalculate them.

6.Streamlit Application Code:

```
1 import streamlit as st
2 import pandas as pd
3 from sklearn.feature_extraction.text import CountVectorizer, TfidfTransformer
4 from sklearn.metrics.pairwise import cosine_similarity
5 import joblib
6
7 # Define global variables to store history and chat data
8 history = []
9 chat_data = []
10
11 def load_data(csv_file):
12     data = pd.read_csv(csv_file)
13     return data
14
15 def load_models():
16     count_vectorizer = joblib.load(r"D:\Internship\Data science\TASKS\search engine\models\count_vectorizer.joblib")
17     tfidf_transformer = joblib.load(r"D:\Internship\Data science\TASKS\search engine\models\tfidf_transformer.joblib")
18     tfidf_matrix = joblib.load(r"D:\Internship\Data science\TASKS\search engine\models\tfidf_matrix.joblib")
19     return count_vectorizer, tfidf_transformer, tfidf_matrix
20
21 def retrieve_similar_documents(query, count_vectorizer, tfidf_transformer, tfidf_matrix, data, top_n=5):
22     query_vector = count_vectorizer.transform([query])
23     query_tfidf = tfidf_transformer.transform(query_vector)
24     similarity_scores = cosine_similarity(query_tfidf, tfidf_matrix)
25     top_indices = similarity_scores.argsort()[0][::-1]
26     retrieved_documents = [data['clean_file_content'][idx] for idx in top_indices[:top_n]]
27     retrieved_subtitle_names = [data['name'][idx] for idx in top_indices[:top_n]] # Assuming subtitle names are stored in 'name' column
28     retrieved_subtitle_nums = [data['num'][idx] for idx in top_indices[:top_n]] # Assuming subtitle numbers are stored in 'subtitle_num'
29
30     return retrieved_documents, retrieved_subtitle_names, retrieved_subtitle_nums
```

```
32 def main():
33     global history, chat_data
34
35     # Customizing title and header
36     st.title('🔍🎬 FilmFinder')
37     st.subheader('🎬 In the realm of movies, allow SeekSpot to guide you through a cinematic adventure, effortlessly finding the films')
38
39     # Sidebar navigation
40     st.sidebar.title('📌 Navigation')
41     if st.sidebar.button('🏠 Home'):
42         st.sidebar.text('Go to Home')
43         # Clear chat data and history
44         history = []
45         chat_data = []
46     if st.sidebar.button('📜 History'):
47         st.sidebar.text('View Search History')
48         # Display search history
49         st.sidebar.write(history)
50     if st.sidebar.button('📄 Export'):
51         st.sidebar.text('Export Data')
52         # Export chat data to a file
53         export_data(chat_data)
54     if st.sidebar.button('⚙️ Settings'):
55         st.sidebar.text('Change Settings')
56
57     # Search functionality
58     query = st.text_input('Enter your query:', '')
59     if st.button('🔍 Search'):
60         if query:
61             # Add query to history
62             history.append(query)
63             # Retrieve similar documents
64             retrieved_documents, retrieved_subtitle_names, retrieved_subtitle_nums = retrieve_similar_documents(query, count_vectorizer, tfidf_transformer, tfidf_matrix, data, top_n=5)
65             st.subheader('📄 Top 5 documents similar to the query:')
```

```

for i, (doc, subtitle_name, subtitle_num) in enumerate(zip(retrieved_documents, retrieved_subtitle_names, retrieved_subtitle_nums)):
    st.write(f"Document {i}")
    st.write(f"Subtitle Name: {subtitle_name}")
    st.write(f"Subtitle Number: {subtitle_num}")
    st.write(f"Summary:", doc)
    # Add search results to chat data
    chat_data.append((query, retrieved_documents, retrieved_subtitle_names, retrieved_subtitle_nums))

def export_data(data):
    # Export chat data to a file
    pass # Placeholder for actual export functionality

if __name__ == '__main__':
    data = load_data(r"b:\internship\Data science\TASKS\search engine\data\eng_subtitles_database.csv")
    count_vectorizer, tfidf_transformer, tfidf_matrix = load_models()
    main()

```

This Streamlit application, named "FilmFinder," serves as a movie search engine. Here's a summary of its features:

Imports: The code imports essential libraries like Streamlit for building the web app, pandas for data handling, and scikit-learn for text processing tasks such as vectorization and similarity computation.

Global Variables: Two lists are globally defined to store the search history (history) and chat data (chat_data).

Functions:

load_data(csv_file): Loads movie subtitle data from a CSV file into a pandas DataFrame.

load_models(): Loads pre-trained models for text processing tasks.

retrieve_similar_documents(query, count_vectorizer, tfidf_transformer, tfidf_matrix, data, top_n): Finds similar movie documents based on a user query using cosine similarity.

main(): The core function of the application. It constructs the layout, including the title, sidebar navigation, search input, and search button. Additionally, it manages the search functionality and displays results.

Main Features:

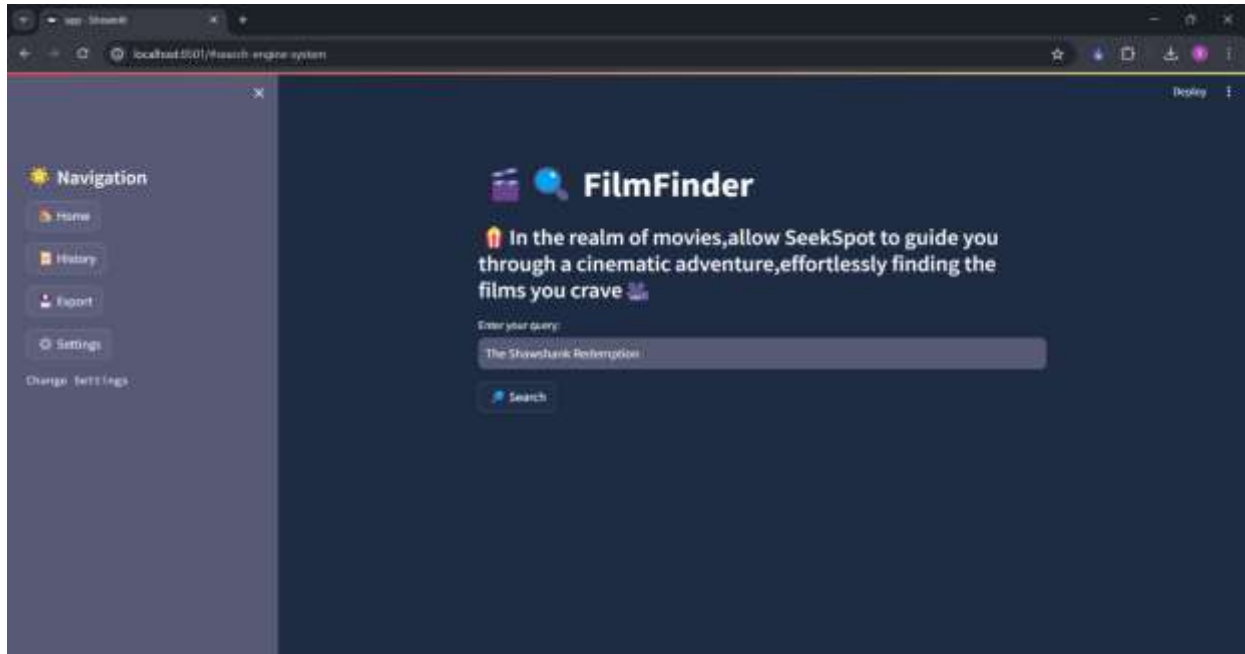
Title and Header: Establishes the title and a subheader for the application.

Sidebar Navigation: Offers buttons for Home, History, Export, and Settings. Each button triggers specific actions, such as clearing history, displaying search history, exporting data, or adjusting settings.

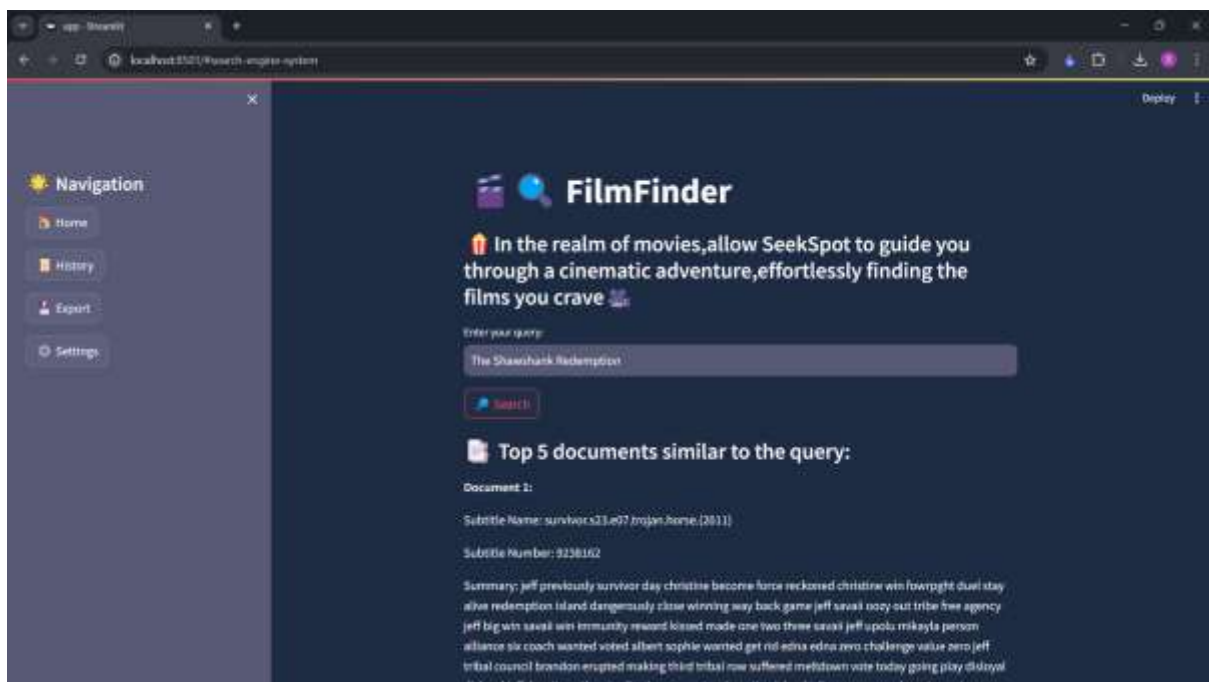
Search Feature: Enables users to input queries and initiate searches. Upon receiving a query, it retrieves similar movie documents using the cosine similarity method and presents the top 5 results.

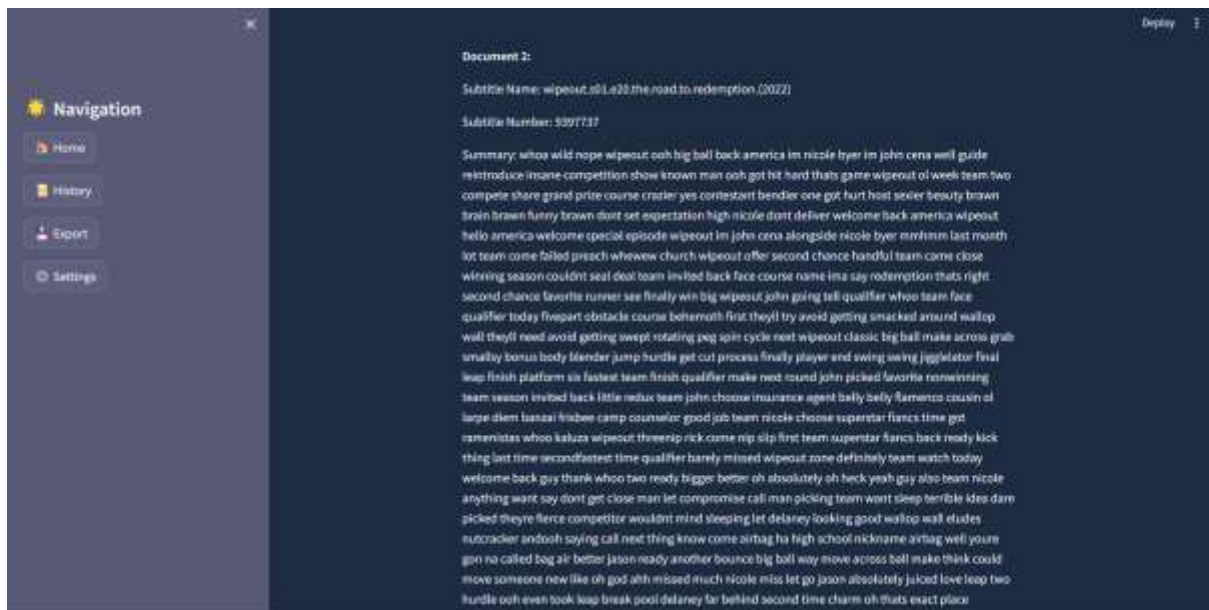
Export Functionality: The export_data function placeholder is set to save chat data to a file, though it currently lacks actual export capabilities.

7.Retrieving Documents using user input Application:

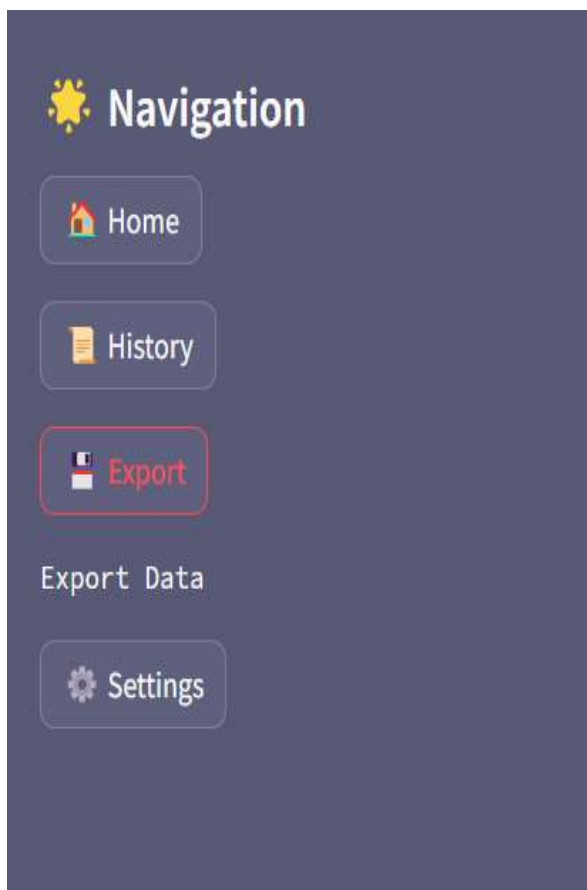


After the Search





Nagivations



Conclusion:

In conclusion, "FilmFinder" is a Streamlit web application designed to simplify the process of finding movies based on user queries. By leveraging machine learning techniques such as text vectorization and cosine similarity, users can input their queries and receive a curated list of similar movie documents. The application offers intuitive navigation through a sidebar menu, allowing users to explore search history, export data, and customize settings. While the current version lacks full export functionality, it provides a solid foundation for further development and enhancement. Overall, "FilmFinder" demonstrates the potential of leveraging modern technologies to streamline the movie search experience for users.