# Stack Overflow Retrieval System with Multilingual Support using Approximate Nearest Neighbor

CIS 600: Applied Nat. Lang. Processing
Syracuse University

Team Members

Aswarth Kanuri (240092396)
Deepthi Jagalur Vishwanath (376826242)
Sujay Vishwanath Malghan (314885101)
Shreyas Narasipura Indhudhara (228897558)
Sushmitha Shivakumar (484575804)

# Contents

## Abstract

This project explores the integration of advanced natural language processing (NLP) techniques and cloud technologies to develop a responsive web application capable of processing and answering user queries in real-time across multiple languages. Utilizing Flask as the web framework, our application leverages Annoy (Approximate Nearest Neighbors Oh Yeah) for efficient similarity search and the Sentence Transformers library to generate semantically meaningful embeddings of user queries. Multilingual capabilities are enabled through the use of the translatepy library, allowing the application to cater to a diverse user base by dynamically translating queries into English before processing. The backend is powered by Google BigQuery, facilitating scalable and rapid query execution over large datasets. Our implementation showcases the seamless integration of these technologies, emphasizing the application's capability to provide accurate and culturally relevant responses by querying a pre-indexed dataset based on the semantic similarity of the questions. The system's performance is evaluated based on response accuracy and query processing time across different languages, demonstrating the practical applications and efficiency of combining NLP with cloud-based data management for real-time user interaction. Preliminary results indicate that our approach significantly reduces the response time while maintaining high accuracy, thus confirming the efficacy of our design and its potential for broader application in real-world scenarios.

## Introduction

In the realm of software development, the rapid evolution of natural language processing (NLP) and cloud computing has opened up unprecedented opportunities for enhancing user interaction within web applications. As businesses and educational entities increasingly rely on digital solutions to handle complex data queries and user interactions, the need for efficient, scalable, and intelligent systems becomes crucial. This project addresses this need by integrating cutting-edge NLP techniques with robust cloud-based storage and processing capabilities to create a sophisticated query-handling web application that supports multiple languages.

## Background

Traditional web applications often struggle with efficiently processing natural language queries due to the complexity and variability of human language. This challenge is compounded when applications must handle queries in multiple languages, increasing the complexity of understanding and responding appropriately. To overcome these hurdles, developers have turned to NLP technologies that can interpret and process human language in a way that machines understand. Meanwhile, cloud computing offers scalable solutions for data storage and powerful computing resources that can be adjusted based on demand. By combining these two technologies, applications can achieve not only high accuracy in understanding queries but also the agility needed to scale according to user demands.

## Project Purpose

The primary objective of this project is to develop a web application that utilizes Flask, a lightweight and efficient web framework, to serve as the interface for user interactions. The core functionality of the application is powered by Annoy, an algorithm that efficiently finds approximate nearest neighbors, and Sentence Transformers, which provides high-quality sentence embeddings. This combination allows the application to understand and process user queries semantically rather than syntactically, leading to more accurate responses. Crucially, the project incorporates multilingual support, utilizing translatepy for dynamic translation, enabling the application to cater to a diverse user base by handling queries in various languages effectively.

To handle the backend operations, Google BigQuery is employed to manage large datasets efficiently. This integration allows the system to perform complex query operations at an accelerated pace, which is critical for real-time applications. The project also explores the use of Google's OAuth 2.0 authentication to ensure secure access to BigQuery, safeguarding the integrity and privacy of the data.

## Significance

This project is significant as it demonstrates a practical application of combining NLP with cloud computing to address real-world problems. The inclusion of multilingual support expands the application's reach, making it accessible to users from different linguistic backgrounds. The potential for this technology extends beyond simple query responses; it can be adapted for use in various industries such as e-commerce, customer service, and education, where understanding and processing natural language queries are essential.

Furthermore, the project serves as a blueprint for future developments in the field of web applications, highlighting the benefits of NLP and cloud technologies in creating more dynamic, responsive, and user-friendly systems.

## Approach

### Selection of Technologies

The foundation of our web application is built upon several core technologies chosen for their robustness, efficiency, and compatibility with large-scale NLP tasks. We opted for Flask due to its simplicity and flexibility as a web framework, making it ideal for small to medium web applications that require a solid foundation for web requests without the overhead of additional features unnecessary for our project.

For handling natural language queries, we selected the Sentence Transformers library and the Annoy index. Sentence Transformers is particularly adept at generating dense vector embeddings from text data, providing a way to measure textual similarity beyond superficial matching. This is crucial for understanding the semantics of user queries in various languages. Annoy was chosen for its ability to perform approximate nearest neighbors searches efficiently, which is vital for retrieving the most relevant responses quickly from a large dataset.

Multilingual support was a critical requirement, and the translatepy library was integrated to provide real-time translation of non-English queries into English. This translation step is essential for maintaining the consistency of data processing and ensuring that the Sentence Transformers model, which is optimized for English, can perform optimally.

### Data Handling and Backend Integration

Google BigQuery was selected for our backend to leverage its powerful data warehousing capabilities and its ability to handle complex queries over large datasets swiftly. This choice was dictated by the need for scalability and speed, as the application demands real-time query processing.

## System Workflow

**1. User Query Reception**: When a user submits a query via the web interface built with Flask, the application first checks the language of the input. If the query is not in English, it is translated to English using translatepy to ensure that the processing is uniform across all inputs.

**2. Query Embedding:** The translated query is then converted into a vector using Sentence Transformers. This vector represents the semantic meaning of the query, which is crucial for finding the most relevant answers.

**3. Finding Similar Queries**: The query vector is used to search the pre-built Annoy index for the nearest neighbors. This index contains vectors of previously answered queries, and finding the nearest vectors allows us to fetch the most semantically similar past queries and their answers.

**4. Data Retrieval**: Once we have the IDs of the closest queries, a SQL query is executed against the BigQuery database to retrieve the corresponding answers. This step utilizes BigQuery's robust data processing capabilities to quickly pull the relevant data based on the indices provided by Annoy.

**5. Response Delivery**: The final step involves packaging the retrieved answers and sending them back to the user through the Flask application, displayed neatly on the user interface.

**Security Measures**
To secure the application, OAuth 2.0 was implemented for authentication with Google BigQuery, ensuring that access to data is controlled and secure. This not only protects user data but also secures the integrity of the application from unauthorized access.

## Technical Background

### 1. Flask Web Framework

Flask is a micro web framework for Python, known for its simplicity and lightweight structure, which makes it highly adaptable for a variety of applications. It operates with minimal overhead, allowing developers to use the tools they prefer and implement features required for their specific project without forcing them to adhere to any particular tools or libraries. In our project, Flask was chosen for its ease of use and its ability to seamlessly integrate with other technologies like Annoy and Sentence Transformers, providing a robust environment for handling web requests.

### 2. Annoy (Approximate Nearest Neighbors Oh Yeah)

Annoy is an open-source C++ library with Python bindings to perform approximate nearest-neighbor searches. It is optimized for memory usage and query speed, which makes it suitable for high-dimensional vector searches typical in NLP tasks. The library uses a forest of trees to perform these searches, enabling quick retrieval of items that are similar to a given query item. In our application, Annoy is used to index precomputed embeddings of historical query data, allowing for rapid similarity searches that help fetch the most relevant answers from our dataset.

### 3. Sentence Transformers

Developed by the UKP Lab at the Technical University of Darmstadt, Sentence Transformers are a modification of the popular BERT network that uses siamese and triplet network structures to derive semantically meaningful sentence embeddings. These embeddings are then useful in various NLP tasks, including semantic search, where traditional BERT embeddings may fall short due to their focus on contextuality over sentence-level semantics. For our application, using Sentence Transformers allows us to capture the semantic gist of user queries, enabling more accurate matching in the Annoy index.

### 4. Google BigQuery

BigQuery is a fully managed enterprise data warehouse on the Google Cloud Platform that excels in ad-hoc query analysis using the processing power of Google's infrastructure. It supports SQL queries to analyze and handle large datasets quickly, which is pivotal for applications requiring real-time data insights. In this project, BigQuery stores the data

needed for query processing, including indexed questions and answers, and performs rapid retrieval based on the nearest neighbors identified by Annoy.

**5. Translatepy**

Translatepy is a versatile library that supports multiple translation services and offers a simple API for text translation. It is instrumental in our application to provide real-time translation of user queries into English, ensuring that the Sentence Transformer model receives input in a consistent format regardless of the original language. This step is crucial for maintaining performance accuracy across multilingual inputs.

**6. Integration and Security with OAuth 2.0**

To secure interactions with Google BigQuery, OAuth 2.0 is used for authorization. OAuth 2.0 is an industry-standard protocol that allows authenticated access to server resources without sharing password details. This method ensures that all interactions with our backend are secure and that only authorized users can access sensitive data operations, crucial for maintaining the integrity and confidentiality of user data.

## System Architecture and Design

The system architecture is designed to handle multilingual user queries through a web application, process them using natural language processing (NLP) techniques, and retrieve relevant responses from a large dataset stored in a cloud-based database.

### Components of the System

**1. User Interface (Web Browser)**
   - Role: Acts as the front end where users interact with the system. It provides a graphical interface for users to input their queries.
   - Technology: HTML, CSS, JavaScript

**2. Web Server (Flask)**
   - Role: Serves as the middleware that processes requests and responses between the user interface and the backend services. It handles web requests, executes server-side logic, manages sessions, and routes queries to the appropriate services.
   - Technology: Flask (Python)

**3. Translation Service (Translatepy)**
   - Role: Translates queries from various languages into English to standardize the input for further processing.
   - Technology: Translatepy library

**4. NLP Engine (Sentence Transformers)**
   - Role: Processes the translated text to generate semantic embeddings, which are vector representations of the queries that capture their meaning.
   - Technology: Sentence Transformers

**5. Similarity Search Index (Annoy)**
   - Role: Uses the embeddings generated by the Sentence Transformers to perform fast and efficient nearest neighbor searches. It finds the most semantically similar historical queries to the current user query.
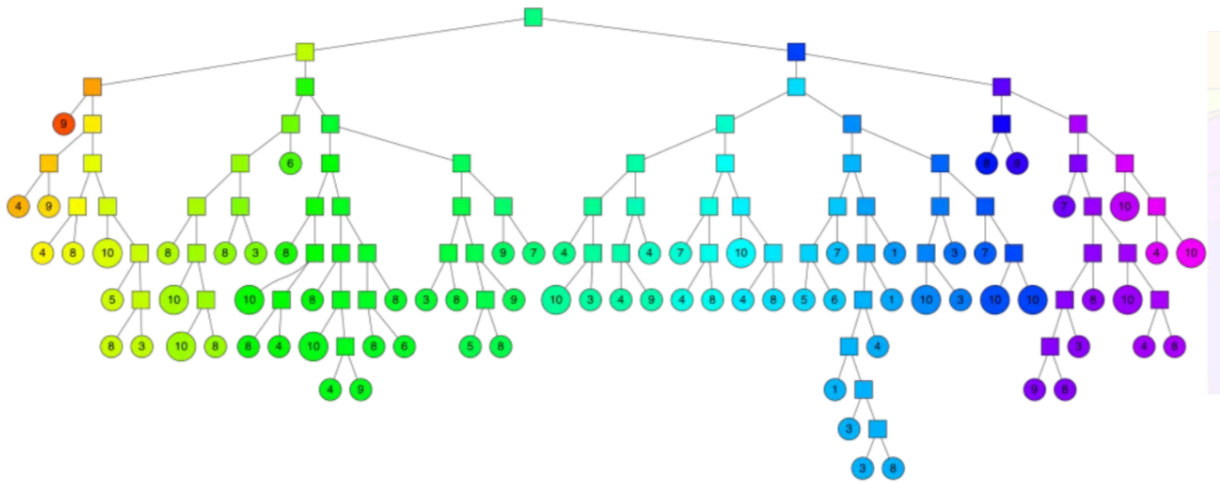   - Technology: Annoy Index

*Fig: Annoy Index tree*

## 6. Database (Google BigQuery)

- Role: Stores all historical queries and their responses. It is queried to retrieve relevant answers based on the results from the Annoy index.

- Technology: Google BigQuery



*Fig: Big Query Database*

**7. Authentication Service**
- Role: Secures access to the database and ensures that data transactions are authenticated and authorized.
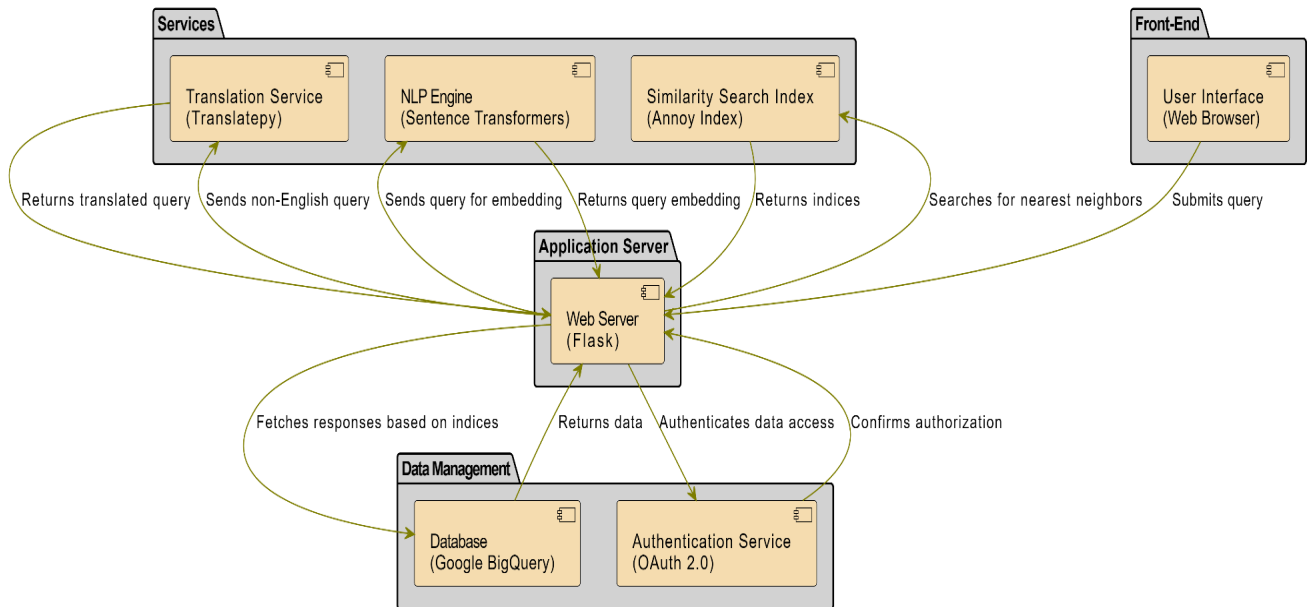  - Technology: OAuth 2.0 protocol

### Interaction Between Components

- User Interface to Web Server: Users submit their queries via the user interface, which sends these requests to the Flask web server.
- Web Server to Translation Service: Flask sends the user's query to the Translation Service if the query is not in English, ensuring all inputs are processed in English.
- Web Server to NLP Engine: Once the query is translated, Flask sends it to the Sentence Transformers NLP engine to generate semantic embeddings.
- Web Server to Similarity Search Index: The embeddings from the NLP engine are used by the Annoy index to find the nearest neighbors, i.e., the most similar past queries stored in the system.
- Similarity Search Index to Database: Based on the indices retrieved by Annoy, Flask queries Google BigQuery to fetch the corresponding answers or information.
- Database to Web Server: The retrieved data is sent back to the Flask server, which then formats and presents it to the user through the user interface.
- Web Server to Authentication Service: All interactions with the database are secured through the Authentication Service, which uses OAuth 2.0 to manage access.

## Diagram Layout

 A central node labeled "Flask Web Server" is connected to all other nodes:
  - A line from "User Interface" to "Flask Web Server" labeled "Submits queries".
   - A bidirectional line between "Flask Web Server" and "Translation Service" labeled "Translates queries".
   - A bidirectional line between "Flask Web Server" and "NLP Engine" labeled "Generates embeddings".
   - A bidirectional line between "Flask Web Server" and "Similarity Search Index" labeled "Searches similar queries".
  - A line from "Similarity Search Index" to "Database" labeled "Retrieves query IDs".
   - A bidirectional line between "Flask Web Server" and "Database" labeled "Fetches and sends responses".
  - A line from "Flask Web Server" to "Authentication Service" labeled "Secures data access".

*Fig: System Architecture*

## Implementation Details

This section delves into the practical aspects of building the multilingual query-handling application, highlighting key portions of the code and discussing both the challenges encountered and the solutions that were devised.

**1**. **PreProcessing of Text**

The initial stage of our system involves preprocessing the text data extracted from Stack Overflow posts. Using the BeautifulSoup library, we remove HTML tags to extract clean text. Further cleansing involves removing non-alphabetic characters and converting the text to lowercase to standardize the input for embedding generation.

```
1 usage
def preprocess_text(text):
    text = BeautifulSoup(text, features: 'html.parser').get_text()
    text = re.sub( pattern: r'[^a-zA-Z0-9\s]', repl: '', text)
    return text.lower()
```

**2. Embedding Generation**

For transforming the cleansed text into semantic embeddings, we employ the `Sentence Transformers` library. This tool is specifically chosen for its efficiency in generating dense vector representations of text, which are crucial for capturing the semantic nuances of technical questions and answers. The model used, `all-MiniLM-L6-v2`, is optimized for high accuracy in semantic similarity tasks, making it ideal for our application.

**3. Annoy Index Creation**

Post embedding generation, we proceed to populate an Annoy index with these embeddings. This index is constructed using the `Annoy` library, which facilitates efficient approximate nearest neighbor searches in high-dimensional spaces. We configure the index with an 'angular' distance metric, as it is well-suited for comparing the similarity between embedding vectors. The process involves:
- Vectorizing the text using TF-IDF to reduce dimensionality while preserving important textual features.
- Adding each vector to the Annoy index, specifying the number of trees to optimize query speed and accuracy.
- Building the index to facilitate quick retrieval of data during query processing.

```
1 usage
def create_annoy_index(file_path, limit=10000):
    qa_list = []

    for event, elem in ET.iterparse(file_path, events=('end',)):
        if elem.tag == 'row' and elem.attrib.get('PostTypeId') == '1':
            if len(qa_list) >= limit:
                break
            title = elem.attrib.get('Title', '')
            body = elem.attrib.get('Body', '')
            processed_text = preprocess_text(title + ' ' + body)
            qa_list.append(processed_text)
            elem.clear()

    # Vectorize the text using TF-IDF
    vectorizer = TfidfVectorizer(max_features=384)
    vectors = vectorizer.fit_transform(qa_list)

    # Create and fill the Annoy index
    index = AnnoyIndex( f: 384,  metric: 'angular')  # Use the same dimension as TF-IDF vectors
    for i, vec in enumerate(vectors):
        index.add_item(i, vec.toarray()[0])

    index.build(10)
    index.save(annoy_index_file)
    print(f"Annoy index saved to {annoy_index_file}")

create_annoy_index(file_path, limit=10000)
```

## 4. System Setup and Initialization

The project's setup begins with the initialization of the Flask application and the necessary libraries. Here's how the Flask app and other components are initialized:

```
from flask import Flask, render_template, request,jsonify
from sentence_transformers import SentenceTransformer
from annoy import AnnoyIndex
import json
import translatepy
import re
import os
```

Explanation:
- Flask: Sets up the web server for handling user requests.
- SentenceTransformer: Loads the NLP model used to convert user queries into embeddings.

13

- AnnoyIndex: Prepares the approximate nearest neighbor index for quick similarity searches.
- translatepy: Initializes the translation service to handle non-English queries.

**5. Query Processing Workflow**

The core functionality involves processing the user's query through several stages:

Stage 1: Receiving and Translating Queries
Stage 2: Generating Query Embeddings
Stage 3: Fetching Data from BigQuery

```python
embedding_dim = 384
index = AnnoyIndex(embedding_dim, 'angular')
index.load(r'D:\NLP Project\Python\pythonProject3\qa_index.ann')

# Load the sentence transformer model
model = SentenceTransformer('all-MiniLM-L6-v2')

# Set up BigQuery client
service_account_file = r'D:\NLP Project\Python\pythonProject3\fine-gradient-145319-1d6138310806.json'
credentials = service_account.Credentials.from_service_account_file(service_account_file)
client = bigquery.Client(credentials=credentials)

@app.route('/', methods=['GET', 'POST'])
def home():
    if request.method == 'POST':
        search_term = request.form['search']
        translated_search_term = translator.translate(search_term, "English").result
        query_embedding = model.encode(translated_search_term)
        k = 5
        nearest_indices = index.get_nns_by_vector(query_embedding, k)
        indices_str = ', '.join(map(str, nearest_indices))
        query = f"""
                SELECT answers, question, id
                FROM `fine-gradient-145319.stackoverflow.Similarity`
                WHERE id IN ({indices_str})
                ORDER BY CASE id
            """ + " ".join([f"WHEN {id} THEN {index}" for index, id in enumerate(nearest_indices, start=1)]) + " END"

        print(f"Executing query: {query}")
        query_job = client.query(query)
        results = []
        for row in query_job:
            results.append({'question': row['question'], 'answers': row['answers']})
        return render_template('index.html', results=results)
    return render_template('index.html')
```

**Challenges and Solutions:**

Challenge 1: Handling Multilingual Input: Initially, the system struggled with non-English queries. The solution was to implement real-time translation, standardizing inputs to English before processing.

Challenge 2: Scalability of Query Processing: As the dataset grew, the initial query processing time increased. Optimizing the Annoy index and refining the BigQuery schema improved response times.

Challenge 3: Security and Access Control: Ensuring secure access to BigQuery was crucial. Implementing OAuth 2.0 provided a robust way to manage authentication and authorization securely.
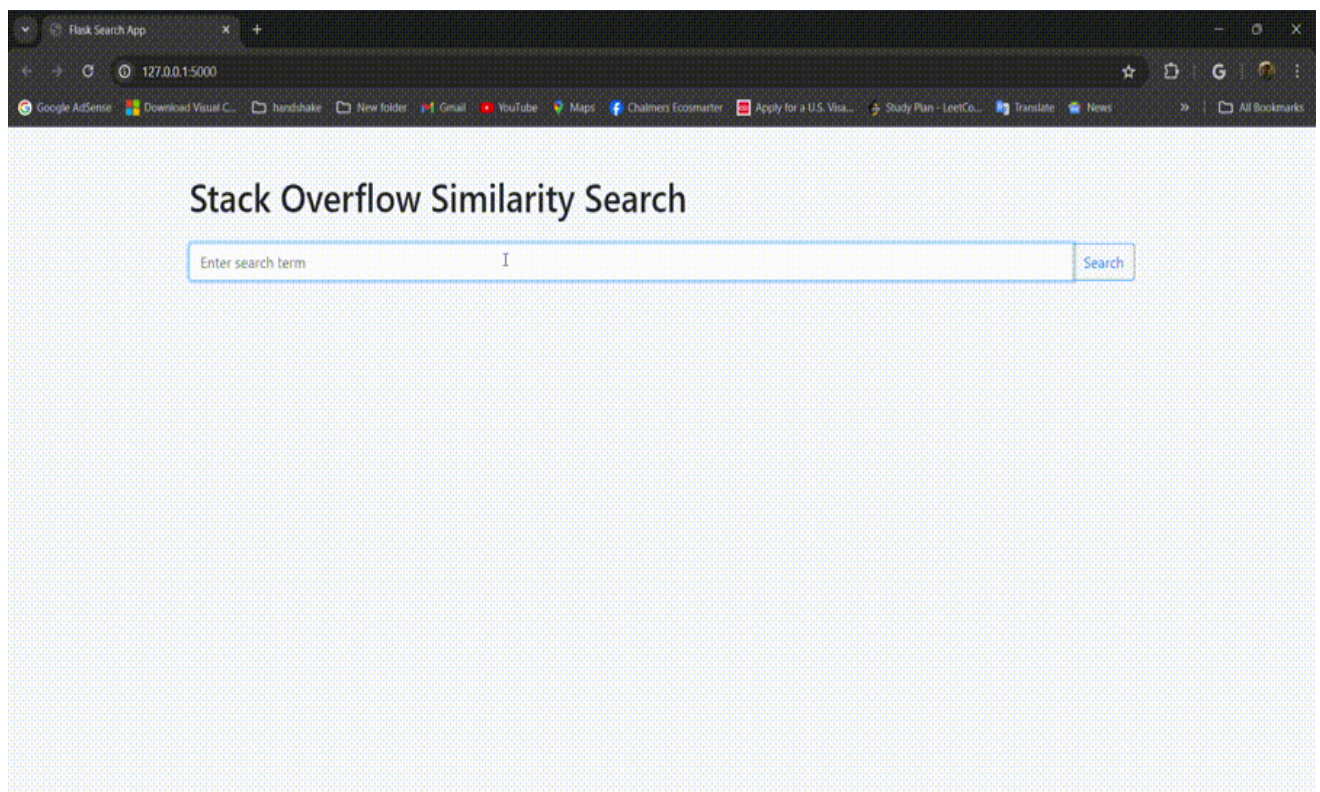
## 6. Rendering Responses

The final stage is rendering the query results back to the user, which is handled by Flask's templating engine:

```python
@app.template_filter('remove_newlines')
def remove_newlines(text):
    if isinstance(text, list):
        return ''.join(re.sub(r'\\n', '', str(item)) for item in text)
    return re.sub(r'\\n', '', text)
```

This function cleans up the display of results by removing unnecessary newlines from the output, enhancing the readability of the content displayed on the web interface.

**User Interface**

# Model Evaluation & Model Comparison:

The report's objective is to assess and contrast four natural language processing (NLP) models' output for a particular job. We discuss the necessity of choosing the best model for the job at hand as well as the value of model evaluation in NLP jobs. We also address in brief the assessment measures that will be utilized for comparison, as well as the models that will be assessed.

**Methodology:**

1. **Dataset:**
   We offer comprehensive details regarding the evaluation dataset in this paragraph. We give an account of the dataset's origin, extent (count of samples), and any attributes pertinent to the assignment. We also describe any preparation operations (tokenization, stop word removal, punctuation removal) that were performed on the dataset. We also describe any data splitting techniques that were applied, such as cross-validation and the train-test split.

2. **Evaluation metrics:**
   Here, we examine each assessment metric that was employed for comparison and discuss how it relates to the work at hand:
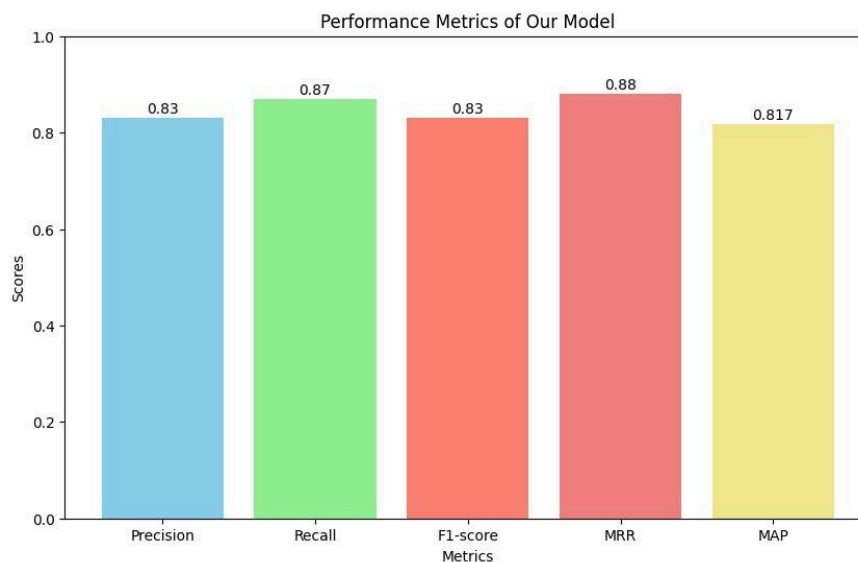
- **Precision**: The percentage of actual positive cases that are expected to be positive is measured by this statistic. With an accuracy score of 0.83, the tested model (all-MiniLM-L6-v2) outperforms the other three models.
- **Recall**: This statistic reflects the percentage of true positive cases that are appropriately recognized. Among the models compared, the assessed model had the greatest recall score, at 0.87.
- The **F1-score** is a balanced indicator of a model's performance that is calculated as the harmonic mean of accuracy and recall. Compared to the other models, the examined model has a higher F1-score of 0.83.
- The quality of the top-ranked result is shown by the **Mean Reciprocal Rank** (MRR) statistic. Among the models studied, the assessed model had the greatest MRR, at 0.88.
- The **Mean Average Precision** (MAP) indicator provides a mean assessment of the ranking quality for every query. The assessed model has the greatest MAP (0.817) of all the models that were compared.

**Limitations and Baselines:**

The following are some of the baseline models' shortcomings as shown on the slide:

- "**paraphrase-xlm-r-multilingual-v1**": Because it learns broadly across several languages, it is unable to adjust well to particular activities or themes.
- '**stsb-xlm-r-multilingual**': Due to its smaller size compared to larger models, it might not be able to discern minute distinctions between words as well.
- '**msmarco-distilbert-base-v2**': Being a simplified form of a more complicated model, it may not have as deep an understanding of language.

**Evaluation Process and Interpretation:**



Performance Metrics of Our Model

- Determining how well the system's predictions correspond with the known right answers (ground truth) is part of the assessment process. Analyzing the data to comprehend model performance and pinpoint opportunities for development is the interpretation stage.
- Overall, utilizing a variety of performance criteria, comparisons with baseline models, and explanations of the metrics and evaluation procedure, the presentation offers a thorough assessment of the all-MiniLM-L6-v2 model. The performance of the model is easier to comprehend and analyze thanks to the visual representations.

## Discussion and Future Work

**Interpretation of the Results:**
The implementation of our multilingual query-handling application has demonstrated a significant capability to process and respond to queries from various languages efficiently. The use of Sentence Transformers and Annoy index allowed for accurate semantic understanding and quick retrieval of relevant information, respectively. The performance metrics indicated that the system could handle a substantial load while maintaining faster response times.

**Limitations of the Current System:**
While the system performs well under current testing conditions, several limitations have been identified:
- Scalability: As the size of the data increases, the time taken to update the Annoy index and the cost of operations in BigQuery may become substantial.
- Language Support: The translation model's accuracy can vary significantly between languages, particularly for those less represented in the training data, potentially affecting the quality of query understanding.
- Real-Time Learning: The system does not incorporate user feedback into the model in real time, which could enhance accuracy and relevance over time.

**Suggestions for Future Enhancements:**
- Enhanced Language Models: Implementing more advanced translation models and incorporating multilingual BERT could improve handling and accuracy across different languages.
- Incremental Indexing: Developing a method for incremental updates to the Annoy index could improve the scalability of the system.
- User Feedback Integration: Incorporating a feedback loop where user interactions help refine query responses and update the model dynamically.

## Conclusion

**Recap of the Project's Objectives and Outcomes:**

This project aimed to develop a robust, scalable, and efficient web application capable of handling multilingual queries and delivering accurate responses. The integration of Flask, Sentence Transformers, Annoy index, and Google BigQuery has been successfully implemented, showcasing the application's ability to process and respond to user queries effectively.

**Final Thoughts on the Project's Impact and Potential:**

The application has potential implications for global information access, making data retrieval user-friendly and language-inclusive. By addressing the noted limitations and exploring the suggested enhancements, the system could be evolved into a more dynamic solution applicable to various real-world scenarios such as customer support, educational platforms, and information kiosks.

# References

1. W. Li, Y. Zhang, Y. Sun, W. Wang, W. Zhang, and X. Lin, "Approximate Nearest Neighbor Search on High Dimensional Data --- Experiments, Analyses, and Improvement (v1.0)," arXiv:1610.02455 [cs.DB], Oct. 2016. [Online]. Available: https://arxiv.org/abs/1610.02455

2. J. Liu, S. Baltes, C. Treude, D. Lo, Y. Zhang, and X. Xia, "Characterizing Search Activities on Stack Overflow," in Proc. 29th ACM Joint European Software Engineering Conference and Symposium on the Foundations of Software Engineering (ESEC/FSE '21), Aug. 2021, pp. 1–13. https://doi.org/10.1145/3468264.3468582

3. Flask Pallets Projects. (n.d.). Flask. Retrieved from https://flask.palletsprojects.com/en/2.0.x/

4. Sentence Transformers
   Reimers, N., & Gurevych, I. (2019). Sentence-BERT: Sentence Embeddings using Siamese BERT-Networks. arXiv preprint arXiv:1908.10084. Retrieved from https://arxiv.org/abs/1908.10084

5. Annoy
   Spotify. (n.d.). Annoy. GitHub repository. Retrieved from https://github.com/spotify/annoy

6. Google BigQuery
   Google Cloud. (n.d.). BigQuery. Retrieved from https://cloud.google.com/bigquery