



AI-Enhanced Information Retrieval System with Big Data Analytics and NLP Transformers

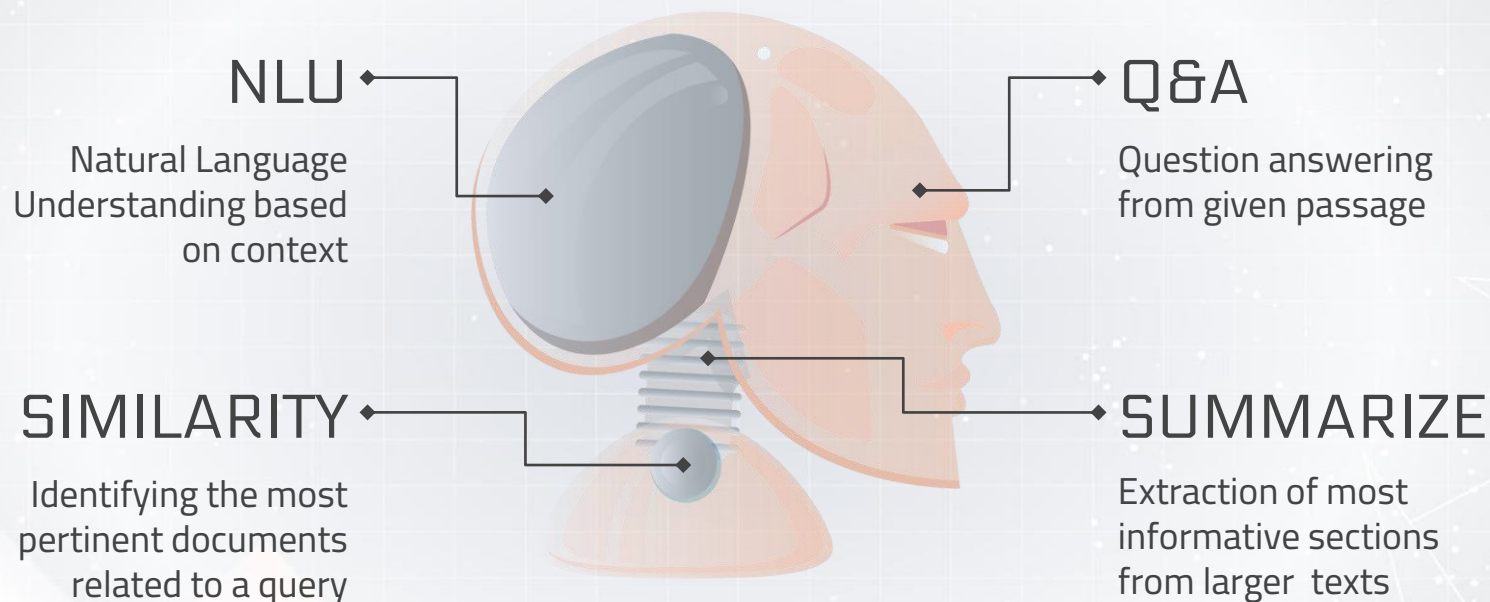
YUNUS EMRE ISIKDEMIR

MOTIVATION

As the volume of data continues to grow, obtaining relevant information from it becomes more challenging, particularly since not all parts of the data are equally important. Therefore, it is necessary to handle, filter, and process the data before modeling it to obtain better inferences. In this regard, Big Data Tools are essential, as they can efficiently process vast amounts of data and retrieve relevant information in a matter of seconds. To this end, various AI-based technologies such as ChatGPT and BARD have emerged, which can extract abstractive information from this massive amount of data.

In response to the challenge of extracting relevant information from large volumes of data, I developed an AI-based information retrieval system in this study. The system enables rapid access to pertinent data and performs specific tasks based on user specifications.

PROBLEM DESCRIPTION

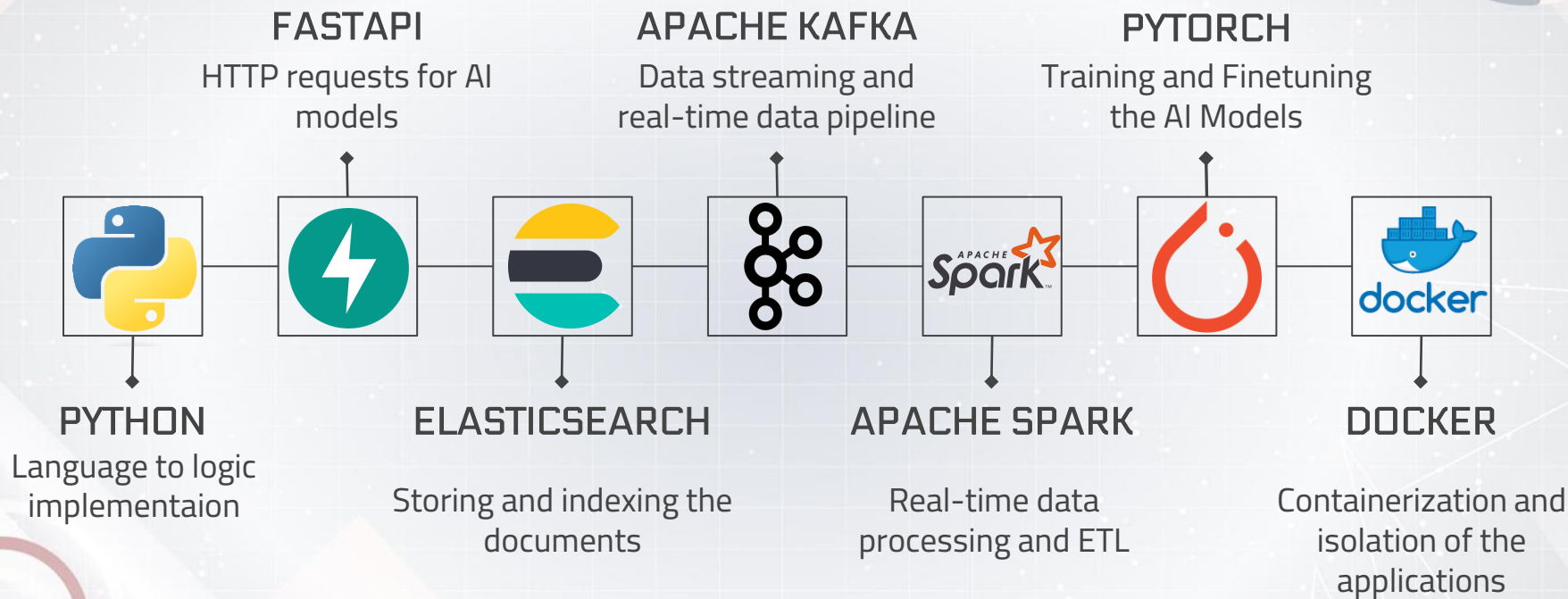




01

TECHNOLOGIES

Languages, Tools and Frameworks



The background is a light gray with a subtle grid pattern. There are several abstract geometric shapes, including triangles and circles, in shades of gray and blue. Some shapes have a slight 3D effect, appearing to be folded or layered.

02

DESIGN
SPECIFICATIONS

STAGES OF THE PIPELINE

01

ETL

The data is first preprocessed with Spark, then consumed from Kafka, before being utilized in Applications.

02

QUERY

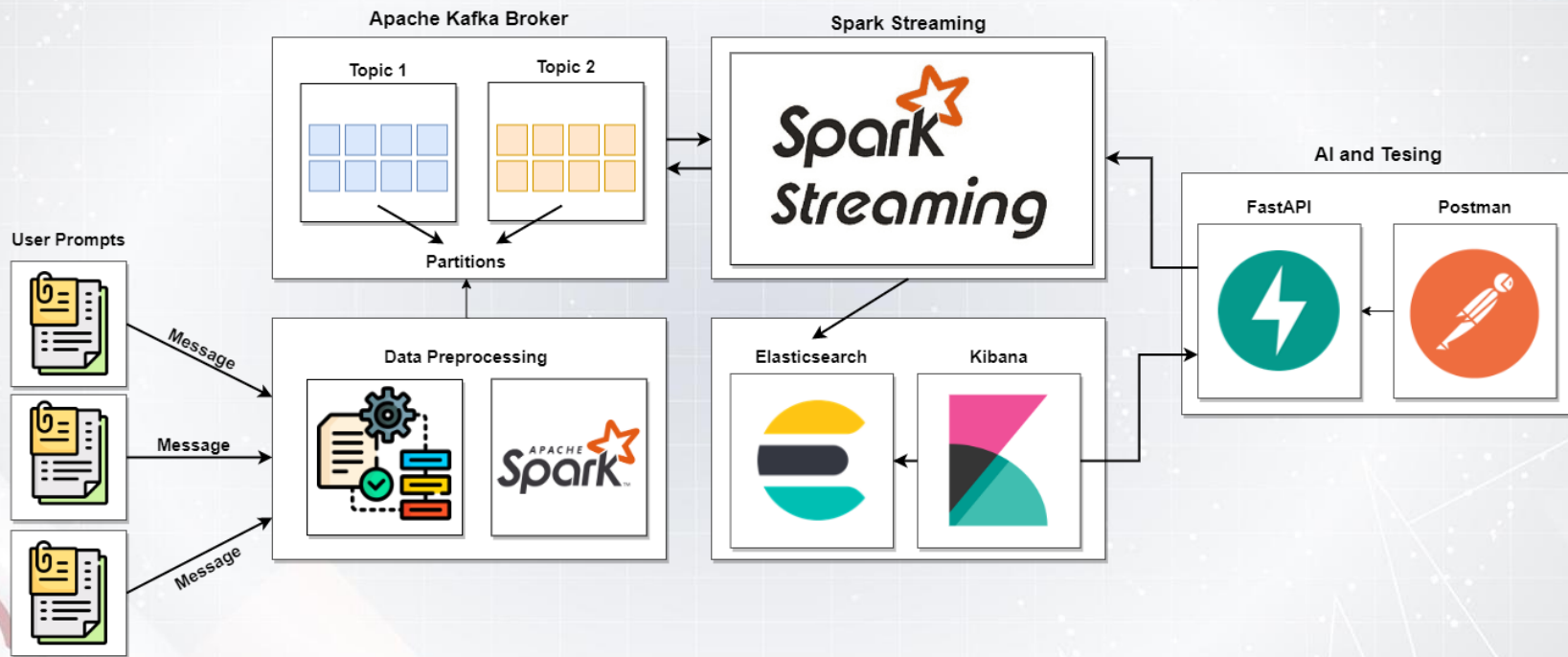
Prompt a query and send a HTTP request to trigger a particular operation

03

RETRIEVE

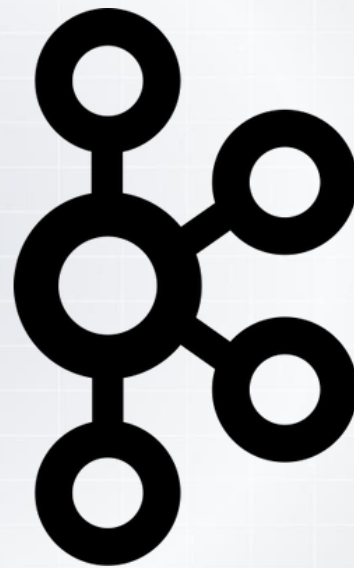
Retrieve relevant data from the database and perform specific actions based on the query.

DATA INGESTION PIPELINE



DATA STREAMING

Data is continuously streamed into Kafka in real-time, making it readily available for consumption by all applications.



DATA WRANGLING

Apache Spark is employed to preprocess text data prior to publishing it, as well as for performing read and write operations with Elasticsearch. By processing queries similarly to indexed documents, the data can be efficiently prepared for AI model execution.



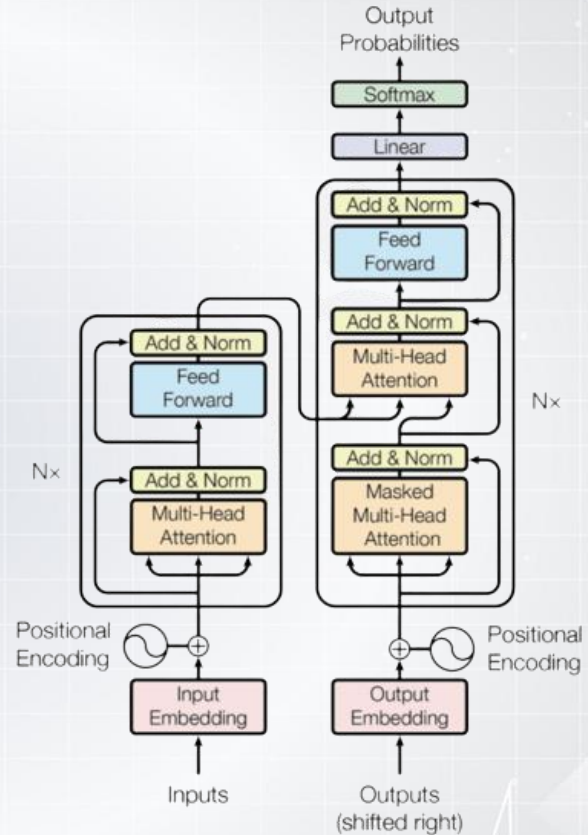
DOCUMENT INDEXING

Elasticsearch is used to store and index documents, allowing for efficient retrieval of information relevant to a given query.



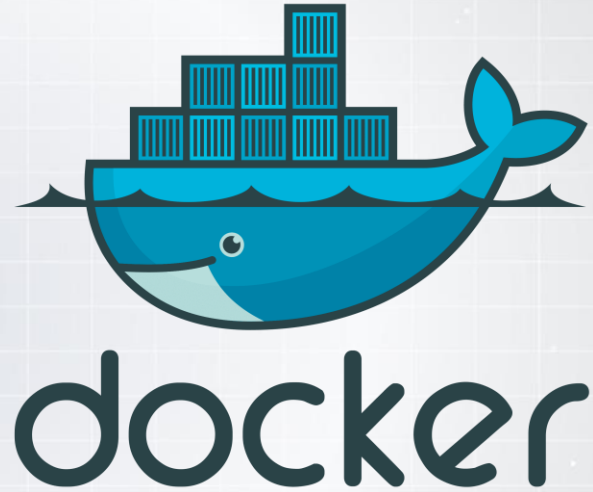
NLP AND NLU

Cutting-edge NLP architectures are utilized to comprehend the provided prompt and carry out a particular task.



CONTAINERIZATION

The services are containerized and isolated using Docker, a platform that allows for efficient and consistent deployment of applications across different environments.



API DEVELOPMENT

The provision of communication between the services is facilitated through HTTP requests via FastAPI.





03

EXPERIMENTAL
RESULTS

POST

{{URL}}/similar-documents?limit=10

Send

Params

Authorization

Headers (9)

Body

Pre-request Script

Tests

Settings

Cookies

Beautiful

none

form-data

x-www-form-urlencoded

raw

binary

GraphQL

JSON

```
1  [
2    {
3      "query": "fetch deforestation related documents"
```

Body

Cookies

Headers (4)

Test Results

Status: 200 OK

Time: 7 ms

Size: 9.08 KB

Save Response

Pretty

Raw

Preview

Visualize

JSON

```
1  [
2    {
3      "id": "1568e7b543748710d3cf3828e1d4a262",
4      "content": "The needs of soy farmers have been used to justify many of the controversial transportation projects that are currently developing in the Amazon. The first two highways
                    successfully opened up the rainforest and led to increased settlement and deforestation. The mean annual deforestation rate from 2000 to 2005 (22,392 km2 or 8,646 sq mi per year)
                    was 18% higher than in the previous five years (19,018 km2 or 7,343 sq mi per year). Although deforestation has declined significantly in the Brazilian Amazon between 2004 and 2014,
                    there has been an increase to the present day.",
5      "content_type": "text",
6      "meta": {},
7      "id_hash_keys": [
8        "content"
9      ],
10     "score": 0.756820706803314,
11     "embedding": null
12   },
13   {
14     "id": "e716f8261700f485df9072cc8b8aa422",
15     "content": "Deforestation is the conversion of forested areas to non-forested areas. The main sources of deforestation in the Amazon are human settlement and development of the land.
                    Prior to the early 1960s, access to the forest's interior was highly restricted, and the forest remained basically intact. Farms established during the 1960s were based on crop
                    cultivation and the slash and burn method. However, the colonists were unable to manage their fields and the crops because of the loss of soil fertility and weed invasion. The soils
                    in the Amazon are productive for just a short period of time, so farmers are constantly moving to new areas and clearing more land. These farming practices led to deforestation and
                    caused extensive environmental damage. Deforestation is considerable, and areas cleared of forest are visible to the naked eye from outer space.",
16     "content_type": "text",
```


POST

{{URL}}/question-answer

Send

Params

Authorization

Headers (9)

Body

Pre-request Script

Tests

Settings

Cookies

Beautify

none

form-data

x-www-form-urlencoded

raw

binary

GraphQL

JSON

```
1 {
2   ... "query": "What is the precise or exact rate of deforestation in numerical terms?"
3 }
```

Body

Cookies

Headers (4)

Test Results

Status: 200 OK

Time: 197 ms

Size: 18.54 KB

Save Response

Pretty

Raw

Preview

Visualize

JSON

```
1 {
2   "query": "What is the precise or exact rate of deforestation in numerical terms?",
3   "no_ans_gap": 5.262350082397461,
4   "answers": [
5     {
6       "answer": "22,392 km2 or 8,646 sq mi per year",
7       "type": "extractive",
8       "score": 0.5973910093307495,
9       "context": "The needs of soy farmers have been used to justify many of the controversial transportation projects that are currently developing in the Amazon. The first two highways
10        successfully opened up the rainforest and led to increased settlement and deforestation. The mean annual deforestation rate from 2000 to 2005 (22,392 km2 or 8,646 sq mi per
11        year) was 18% higher than in the previous five years (19,018 km2 or 7,343 sq mi per year). Although deforestation has declined significantly in the Brazilian Amazon between 2004
12        and 2014, there has been an increase to the present day.",
13       "offsets_in_document": [
14         {
15           "start": 312,
16           "end": 346
17         }
18       ],
19       "offsets_in_context": [
20         {
21           "start": 312,
22           "end": 346
23         }
24       ]
25     }
26   ]
27 }
```