

Project Title:

World's Top 50 Hospitals Data Collection and Private GPT Model Training

Table of Contents

- 1. Project Overview**
- 2. Data Collection**
 - 2.1 Data Sources**
 - 2.2 Web Scraping Methodology**
- 3. Data Cleaning**
 - 3.1 Cleaning Methodology**
- 4. Model Training**
 - 4.1 Model Repository**
 - 4.2 Model Training Process**
- 5. Documentation of Challenges**
- 6. Model Performance**
 - 6.1 Model Interaction Snippets**
- 7. Replication Guide**
 - 7.1 Installation**
 - 7.2 Data Scraping**
 - 7.3 Data Cleaning**
 - 7.4 Model Training**
- 8. Conclusion**

1. Project Overview

This project focuses on collecting data from the "World's Top 50 Hospitals 2023" ranking and training a Private GPT model for healthcare-related tasks. It includes data collection, cleaning, model training, and documentation of the process.

2. Data Collection

2.1 Data Sources

Data was obtained from the "World's Top 50 Hospitals 2023" ranking published by Newsweek <https://www.newsweek.com/rankings/worlds-best-hospitals-2023>.

2.2 Web Scraping Methodology

Python libraries, including BeautifulSoup and requests, were employed for web scraping.

Hospital names and website URLs were extracted from the Newsweek ranking.

Data from each hospital's website, including name, location, services, and awards, was collected.

3. Data Cleaning

3.1 Cleaning Methodology

Data cleaning steps included:

- Handling missing data through imputation or marking.
- Standardizing data to ensure consistent formats.
- Removing duplicate entries within the dataset.

4. Model Training

4.1 Model Repository

Model training utilized the Private GPT repository hosted on GitHub: Private GPT Repository
<https://github.com/imartinez/privateGPT>.

4.2 Model Training Process

The cleaned data was preprocessed for model training, including text encoding.

The Private GPT model was fine-tuned for healthcare-related tasks, such as generating descriptions or answering questions about hospitals.

5. Documentation of Challenges

Several challenges were encountered during the project:

- Rate Limiting: Website requests had to be managed carefully as frequent pinging was not allowed.
- Access Issues: Some websites were not accessible, resulting in permission denied errors.
- Data Structure: Data from various hospital websites often appeared cluttered.

6. Model Performance

Model performance was assessed through interactions with the Private GPT model. Sample interactions, questions, and responses were documented to gauge the model's effectiveness.

6.1 Model Interaction Snippets

snippets or examples of interactions with the model, demonstrating its capabilities and limitations.

7. Replication Guide

To replicate this project, follow these steps:

7.1 Clone the Private GPT model repository from [GitHub](<https://github.com/imartinez/privateGPT>).

7.2 Data Scraping

Code:

```
# Function to scrape, clean, and display hospital data
def scrape_clean_and_display_hospital_data():
    # Send an HTTP GET request to the Newsweek URL
    newsweek_response = requests.get(newsweek_url, headers=headers)

    # Check if the request to Newsweek was successful
    if newsweek_response.status_code == 200:
        # Parse the HTML content of the Newsweek page using BeautifulSoup
        newsweek_soup = BeautifulSoup(newsweek_response.text, 'html.parser')

        # Find the table element with the specified class
        table = newsweek_soup.find('table', class_='ranking-table')

        # Extract rows from the table
        rows = table.find_all('tr')

        # Initialize counters to skip the 9th and 27th links
        link_counter = 0

        # List to store scraped data
        scraped_data = []

        # Initialize variables to store the current hospital data
        current_hospital_data = {
            "p_tags": [],
```

```

    "h1_tags": [],
    "h2_tags": [],
    "h3_tags": [],
    "div_tags": [],
}
current_hospital_name = ""
current_hospital_rank = ""

# Loop through the rows and scrape content from the top 50 hospitals with links
for row in rows[1:51]: # Skip the header row and limit to the top 50 hospitals
    columns = row.find_all('td')
    rank = columns[0].get_text(strip=True)
    publication_name = columns[1].find('a').get_text(strip=True)
    country = columns[2].get_text(strip=True)
    city = columns[3].get_text(strip=True)
    state = columns[4].get_text(strip=True)
    link = columns[1].find('a')['href']

    # Increment the link counter
    link_counter += 1

    # Skip the 9th and 27th links and continue to the next link
    if link_counter == 9 or link_counter == 27:
        continue

    # Introduce a delay before sending the request to the hospital's link
    time.sleep(2) # Wait for 2 seconds (adjust the delay as needed)

    # Send an HTTP GET request to the hospital's link
    hospital_response = requests.get(link, headers=headers)

    # Check if the request to the hospital's link was successful
    if hospital_response.status_code == 200:
        # Parse the HTML content of the hospital's page using BeautifulSoup

```

```
hospital_soup = BeautifulSoup(hospital_response.text, 'html.parser')
```

```
# Extract and clean content from p, h1, h2, h3, and div tags
```

```
p_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('p')]
```

```
h1_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('h1')]
```

```
h2_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('h2')]
```

```
h3_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('h3')]
```

```
div_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('div')]
```

```
# Store the cleaned data in the current hospital data dictionary
```

```
current_hospital_data["p_tags"] = p_tags
```

```
current_hospital_data["h1_tags"] = h1_tags
```

```
current_hospital_data["h2_tags"] = h2_tags
```

```
current_hospital_data["h3_tags"] = h3_tags
```

```
current_hospital_data["div_tags"] = div_tags
```

```
# Store the current hospital name and rank
```

```
current_hospital_name = publication_name
```

```
current_hospital_rank = rank
```

```
# Append the current hospital data to the list
```

```
scraped_data.append({
```

```
    "Hospital Rank": current_hospital_rank,
```

```
    "Hospital Name": current_hospital_name,
```

```
    **current_hospital_data
```

```
})
```

```
return scraped_data
```

7.3 Data Cleaning

Function to scrape, clean, and display hospital data

def scrape_clean_and_display_hospital_data():

... (previous code)

Loop through the rows and scrape content from the top 50 hospitals with links

for row in rows[1:51]: # Skip the header row and limit to the top 50 hospitals

columns = row.find_all('td')

rank = columns[0].get_text(strip=True)

publication_name = columns[1].find('a').get_text(strip=True)

country = columns[2].get_text(strip=True)

city = columns[3].get_text(strip=True)

state = columns[4].get_text(strip=True)

link = columns[1].find('a')['href']

... (previous code)

Send an HTTP GET request to the hospital's link

hospital_response = requests.get(link, headers=headers)

Check if the request to the hospital's link was successful

if hospital_response.status_code == 200:

Parse the HTML content of the hospital's page using BeautifulSoup

hospital_soup = BeautifulSoup(hospital_response.text, 'html.parser')

Extract and clean content from p, h1, h2, h3, and div tags

p_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('p')]

h1_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('h1')]

h2_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('h2')]

h3_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('h3')]

div_tags = [clean_html_tags(tag.get_text()) for tag in hospital_soup.find_all('div')]

... (previous code)

```
# Store the cleaned data in the current hospital data dictionary
```

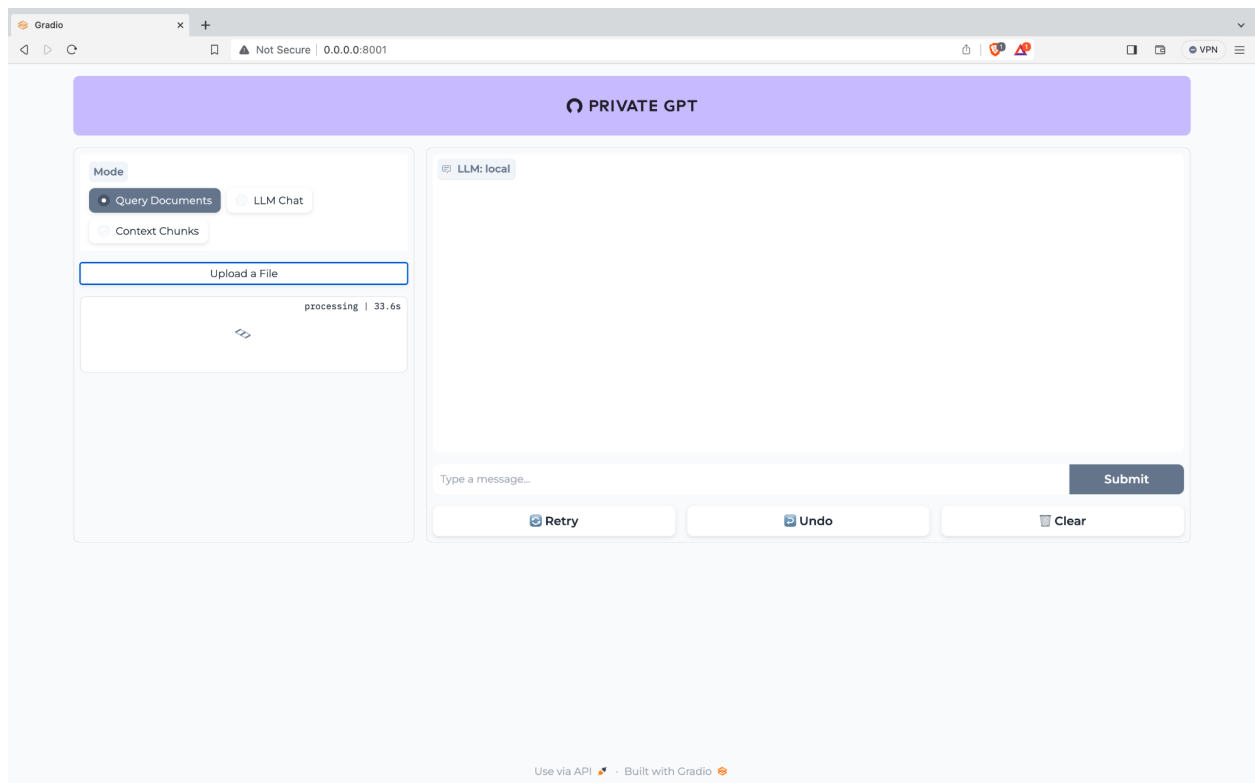
```
# Append the current hospital data to the list
```

```
scraped_data.append({  
    "Hospital Rank": current_hospital_rank,  
    "Hospital Name": current_hospital_name,  
    **current_hospital_data  
})
```

```
return scraped_data
```

7.4 Model Training

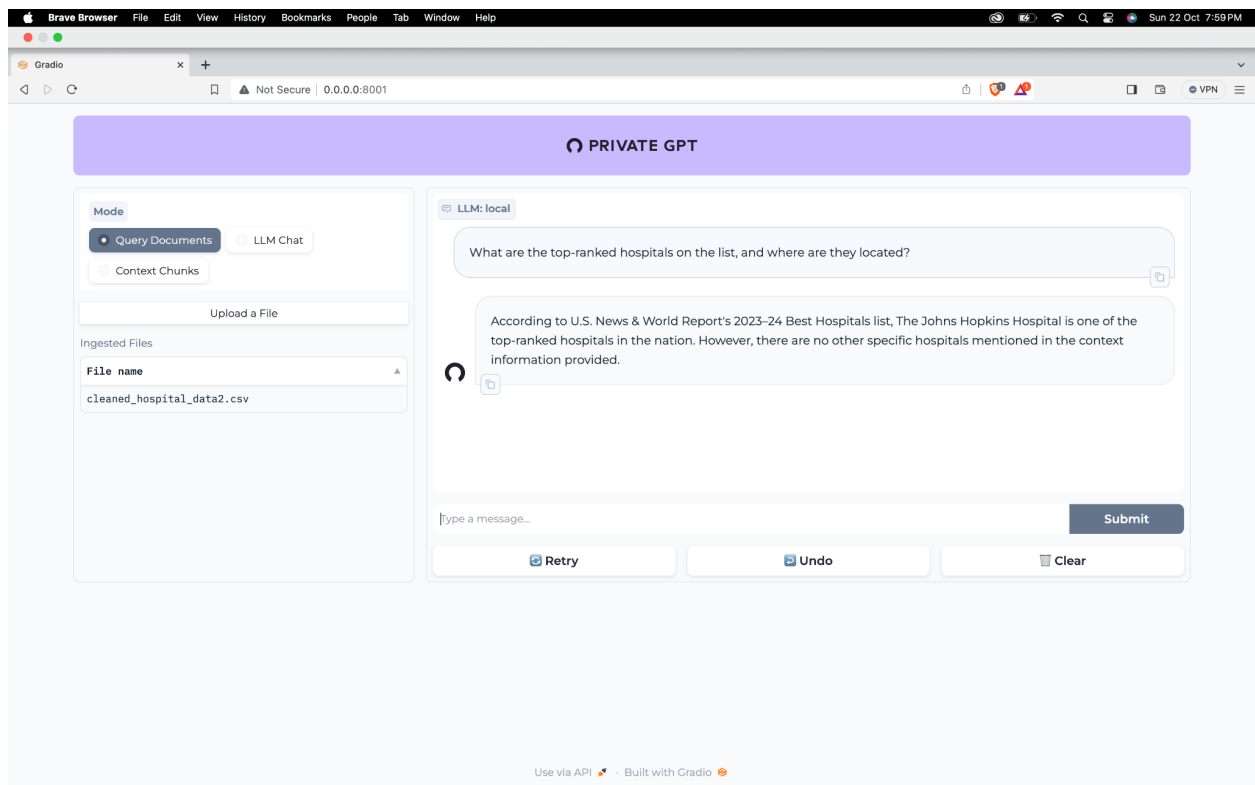
Img1:Uploding A Csv File.



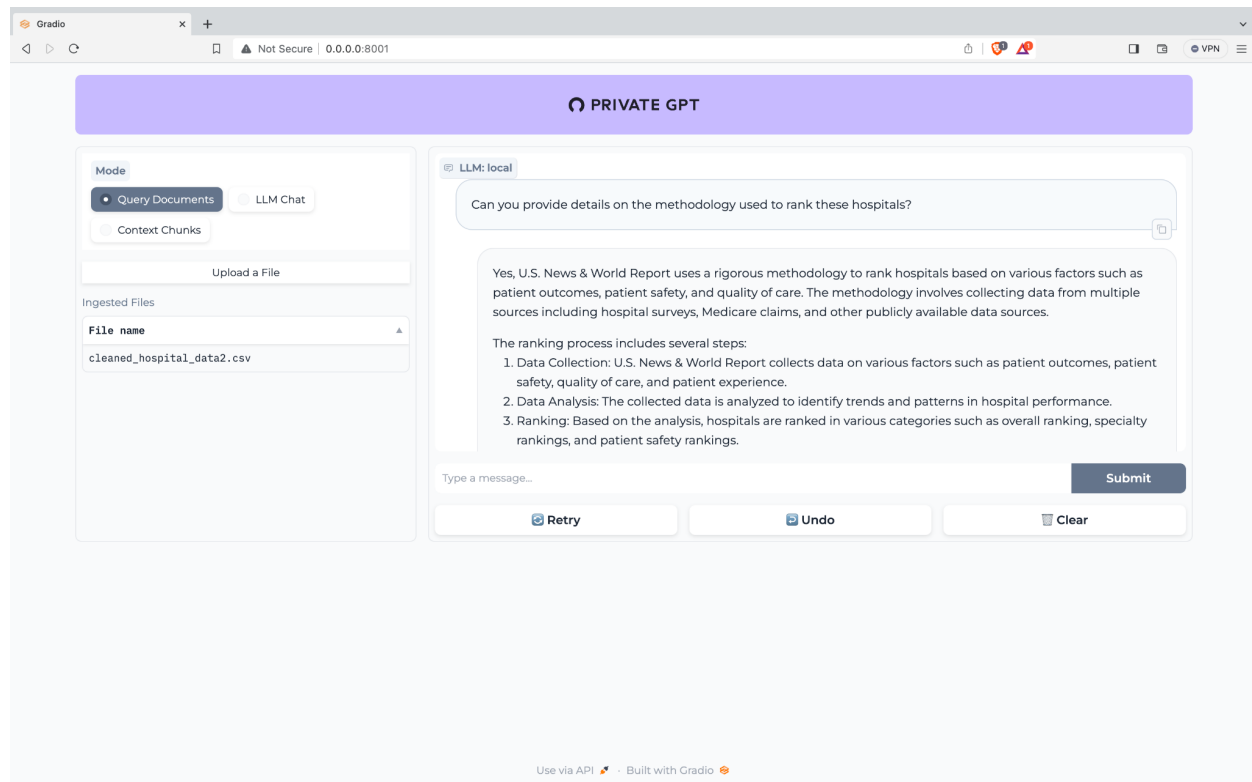
Img2: Executing a Csv File in Background.

```
~/Desktop/privateGPT -- Python - make run
llama_new_context_with_model: n_ctx = 5000
llama_new_context_with_model: freq_base = 48000.0
llama_new_context_with_model: freq_scale = 1
llama_new_context_with_model: kv self size = 487.68 MB
llama_new_context_with_model: compute buffer total size = 281.25 MB
llama_new_context_with_model: max tensor size = 182.64 MB
ggml_metal_add_buffer: allocated data buffer, size = 4896.00 MB, offs = 0
ggml_metal_add_buffer: allocated data buffer, size = 172.64 MB, offs = 4187422720, ( 4269.27 / 5461.34)
ggml_metal_add_buffer: allocated kv buffer, size = 489.68 MB, ( 4758.77 / 5461.34)
ggml_metal_add_buffer: allocated alloc buffer, size = 275.00 MB, ( 5034.14 / 5461.34)
AVX = 0 | AVX2 = 0 | AVX512 = 0 | AVX512_VBMI = 0 | AVX512_VNNI = 0 | FMA = 0 | NEON = 1 | ARM_FMA = 1 | F16C = 0 | FP16_VA = 1 | WASM_SIMD = 0 | BLAS = 1 | SSE3 = 0 | SSSE3 = 0 | VSX = 0 |
INFO: Started server process (18200)
INFO: Waiting for application startup.
INFO: Application startup complete.
INFO: Uvicorn running on https://0.0.0.0:8001 (Press CTRL+C to quit)
INFO: 127.0.0.1:58456 - "GET / HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/index-6e28cf60.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/index-45a964fd.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /info HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /theme.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /favicon.ico HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/Blocks-9a8b13e0.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/Button-6f70b409.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/Button-c1e6319b.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/Blocks-0733f3b3.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/index-95c9f556.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/index-d3c2a888.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58460 - "GET /assets/index-312ceac2.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/index-320f25d0.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58459 - "GET /assets/index-3d71b090.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/StaticColumn-97e4efdc.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/StaticColumn-203a3031.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58460 - "GET /assets/Radio-38bb7cb6.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/StaticMarkdown-733a2e11.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/UploadButton-03d58ab8.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58459 - "GET /assets/index-07f64bdc.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/Table-4bac0ba6.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/Textbox-dde0f8cc.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/index-57819924.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/StaticForm-3812b7f1.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/index-53bad8d6.css HTTP/1.1" 200 OK
INFO: 127.0.0.1:58459 - "GET /assets/Radio-8cc59070.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/BlockTitle-9e3db0fc.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58460 - "GET /assets/index-2fce551f.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/UploadButton-07793b08.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/info-bb061690.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/index-bd0e0f8c.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58459 - "GET /assets/Table-c0811abf.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/utills-c3e3db68.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58460 - "GET /assets/Upload-1ced8f8a.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/StaticMarkdown-4afac2a3.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/dav-076afacd.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/index-3bb0d053.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58459 - "GET /assets/ShareButton-497a067a.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/CommitButton-1360929c.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58460 - "GET /assets/Copy-d179d5ea.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/BlockLabel-7794b0ff.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58459 - "GET /assets/index-9a3272d9.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/index-04fedabc.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/Textbox-27167496.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/InteractiveTextbox-5c1e1ef9.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58459 - "GET /assets/index-83260fdd.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58456 - "GET /assets/StaticForm-f0f2a219.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "GET /assets/index-fa962f4c.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/index-b79eddbb.js HTTP/1.1" 200 OK
INFO: 127.0.0.1:58457 - "GET /assets/Img-ea87d6cf.svg HTTP/1.1" 200 OK
INFO: 127.0.0.1:58458 - "GET /assets/api-logo-63a6f193.svg HTTP/1.1" 200 OK
INFO: 127.0.0.1:58461 - "POST /upload HTTP/1.1" 200 OK
INFO: (127.0.0.1, 58446) - "WebSocket /queue/join" [accepted]
INFO: connection open
Pasting documents [enter mode]: 100%
Generating embeddings: 20% | 1/1 [00:00:00.00, 3.75it/s] | 1030/5200 [00:48:02:07, 32.70it/s]
```

Img3: Query:-What are the top Ranked Hospital



Img4:Query Can you provide details on the methodology used to rank these hospitals



Query You can explore more here is a List:-

1. What are the top-ranked hospitals on the list, and where are they located?
2. Can you provide details on the methodology used to rank these hospitals?
3. Are there any notable trends or changes in the rankings compared to previous years?
4. What criteria were used to assess the quality of these hospitals?
5. How many hospitals were considered in this ranking?
6. Can you provide information about the best hospitals in a specific country or region?
7. Are there any specialized or top-performing departments or services within these hospitals?
8. Do the rankings take into account the hospital's patient outcomes or other performance indicators?
9. Are there any specific medical breakthroughs or innovations associated with the top-ranked hospitals?
10. Can you provide information on the services and specialties offered by these hospitals?
11. Is there a way to access more detailed profiles or information about each hospital on the list?
12. Are there any specific areas of medical expertise or research where these hospitals excel?
13. What is the geographical distribution of the top hospitals? Are there notable clusters in certain regions?

14. Are there any partnerships or collaborations these hospitals have with other healthcare institutions?
15. Is there information about the accessibility and affordability of healthcare services in these top-ranked hospitals?

8. Conclusion

In conclusion, the "World's Top 50 Hospitals Data Collection and Private GPT Model Training" project successfully collected data from leading hospitals, cleaned and processed the data, and trained a Private GPT model for healthcare-related tasks. Challenges were documented, and model performance was evaluated. The provided replication guide allows others to reproduce the project's results.