



# MedExtractAI

—

Web Scraping and AI Model Training

## Project Overview

The goal of this project was to scrape data from the top 50 hospitals worldwide, clean and preprocess the data, and then train a **Private GPT** model on this dataset. The Private GPT model, sourced from GitHub, was fine-tuned to generate hospital-related content, including information about hospital services, departments, and general healthcare topics.

## Steps

1. Data collection: Scraped relevant data from the websites of the top 50 hospitals.
2. Data cleaning: Cleaned and structured the scraped data for model training.
3. Model training: Fine-tuned a **Private GPT** model using the cleaned data.

### 1. Data Collection

- We began by identifying the top 50 hospitals in India using [www.newsweek.com/rankings/worlds-best-hospitals-2024/india](http://www.newsweek.com/rankings/worlds-best-hospitals-2024/india)
- These rankings provided us with a list of hospitals across various specialties, which were then used to gather URLs for scraping.
- The script `extract_top_hospitals.py` was used to automatically fetch the URLs of the top 50 hospitals by scraping the list.
- Once we had the URLs of the hospitals, the next task was to extract relevant data from their websites.
- The script `hospital_data_scraper.py` was used to scrape.
- We used BeautifulSoup for scraping.
- The scraped data was saved in `scraped_data.json`.

### 2. Data Cleaning

- Load JSON Data: Read the data from `scraped_data.json`.
- Extract Text: Access the value of `'content' ['p']` for each entry.
- Remove HTML Tags: Use `BeautifulSoup` to strip any embedded HTML tags.
- Normalize Whitespace: Replace all types of whitespace (`\s+`) with a single space.
- Remove Escape Sequences: Clean up characters like `\t`, `\n`, and zero-width spaces (`\u200b`).

- Remove Special Characters: Delete any non-alphanumeric characters except punctuation marks.
- Convert to Lowercase: Standardize text to lowercase.
- Trim Spaces: Remove leading and trailing spaces.
- Write Cleaned Data to File: Save the cleaned text to `cleaned_dataset.txt`.

### 3. Model Training using Private GPT

- For this project, we used **PrivateGPT**, a private version of GPT that allows you to fine-tune a GPT model on your own dataset. PrivateGPT was sourced from GitHub and is designed to be flexible and customizable for various types of data.
- The fine-tuning was done on the cleaned hospital data to generate content relevant to healthcare and hospital-related inquiries.

### Findings and Performance of the Model

- **Quality of Data:** The quality of the scraped data was directly linked to the performance of the trained model. Hospital websites with structured and detailed information resulted in more accurate and coherent outputs from the model.
- **Model Performance:** After fine-tuning, the Private GPT model was able to generate coherent responses to healthcare-related queries, including hospital information, services, and general healthcare knowledge.

### Screenshots

```

Enter a query: What are the specialties at Fortis Hospital, Vasant Kunj?
The hospital offers services across various medical fields such as Cardiology, Neurology, Neurosurgery, Orthopedics, Pediatrics, Pulmonology, Urology and Vascular Surgery.

> Question:
What are the specialties at Fortis Hospital, Vasant Kunj?

> Answer (took 31.75 s.):
The hospital offers services across various medical fields such as Cardiology, Neurology, Neurosurgery, Orthopedics, Pediatrics, Pulmonology, Urology and Vascular Surgery.

> source_documents/scraped_data [MConverter.eu].txt:
"Fortis Hospital, Vasant Kunj's broad range of specialties and its extensive experience in treating millions of patients reflect its dedication to providing comprehensive, expert healthcare. The hospital's multidisciplinary approach ensures that patients receive tailored and effective care across a wide spectrum of medical fields.",
"Infrastructure and Capacity",

> source_documents/scraped_data [MConverter.eu].txt:
"Fortis Hospital, Vasant Kunj, is having facility with over 162 beds, including a substantial 65 critical care beds including Medical Surgery, Neonatal Paediatric ICU. This extensive capacity ensure s that the hospital can cater to a wide range of patient needs, from routine check-ups and elective surgeries to complex critical care. The hospital's infrastructure is designed to support various m edical procedures and patient care requirements efficiently. 6 well-equipped operating

> source_documents/scraped_data [MConverter.eu].txt:
"Fortis Hospital, Vasant Kunj, stands as a premier healthcare institution with a storied legacy of 18 years in delivering exemplary medical care. Located in the upscale locality of Vasant Kunj in De lhi, this hospital is recognized for its commitment to high standards of medical excellence, advanced technology, and comprehensive patient care. Since its establishment in 2006, Fortis Hospital, Va sant Kunj, has become a leading healthcare provider, having successfully treated over 2

> source_documents/scraped_data [MConverter.eu].txt:
"Fortis Hospital, Vasant Kunj, stands as a premier healthcare institution with a storied legacy of 18 years in delivering exemplary medical care. Located in the upscale locality of Vasant Kunj in De lhi, this hospital is recognized for its commitment to high standards of medical excellence, advanced technology, and comprehensive patient care. Since its establishment in 2006, Fortis Hospital, Va sant Kunj, has become a leading healthcare provider, having successfully treated over 2

```

