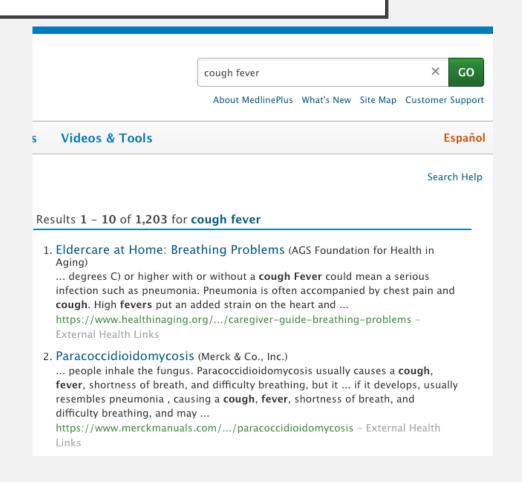# DISEASE-

A simple search engine, with a simple web crawler that focus on the disease

Rujun Yao

# PROBLEMS

- Current pandemic of COVID-19

- It is hard to search for symptoms from search engine

- Low efficiency of Hospital

- "Shelter In Place" Order

# SOLUTION

- Web crawler to crawl a medical website (depth = 3)
- Html parser to parse links and text (a tag href attr)
- Store texts locally (offline)
- Indexed locally stored documents (inverted index term:docids)
- Query keywords
- Tf-idf cosine score rank system

# SEED

- https://www.cdc.gov/
    - Many unrelated information
    - Pages in other language
    - Disease page, symptoms page are separated
- https://www.who.int/
    - Same problem as cdc
- https://medlineplus.gov/
    - Information are gathered together (related info)

# CRAWLER

- Initialization: visited_set(), front_queue(), seed

- For each link in the queue:

  - Check if the link is valid. (external links? Spanish?)

  - Extract links, add all valid link into front_queue

  - Extract metadata and texts, save texts as local docs

# INDEXER

- Remove pos of term in inverted index

- Form: {term1: {docID, docID, …}, term2: {docID, docID,…}, …}

- Other information: vector length, term frequency, document frequency and inverse document frequency, stop-word list

- Since there are up to 10,000 documents, The indexing process will take lot of time and memory, all of the information can be cached in disk storage.

- Use pickle.dump() and pickle.load() to save or load dictionary data structure in disk as binary file.

# QUERY

- Cosine similarity and tf-idf

- Could load index from file system

- Use subprocess to open the retrieved documents.