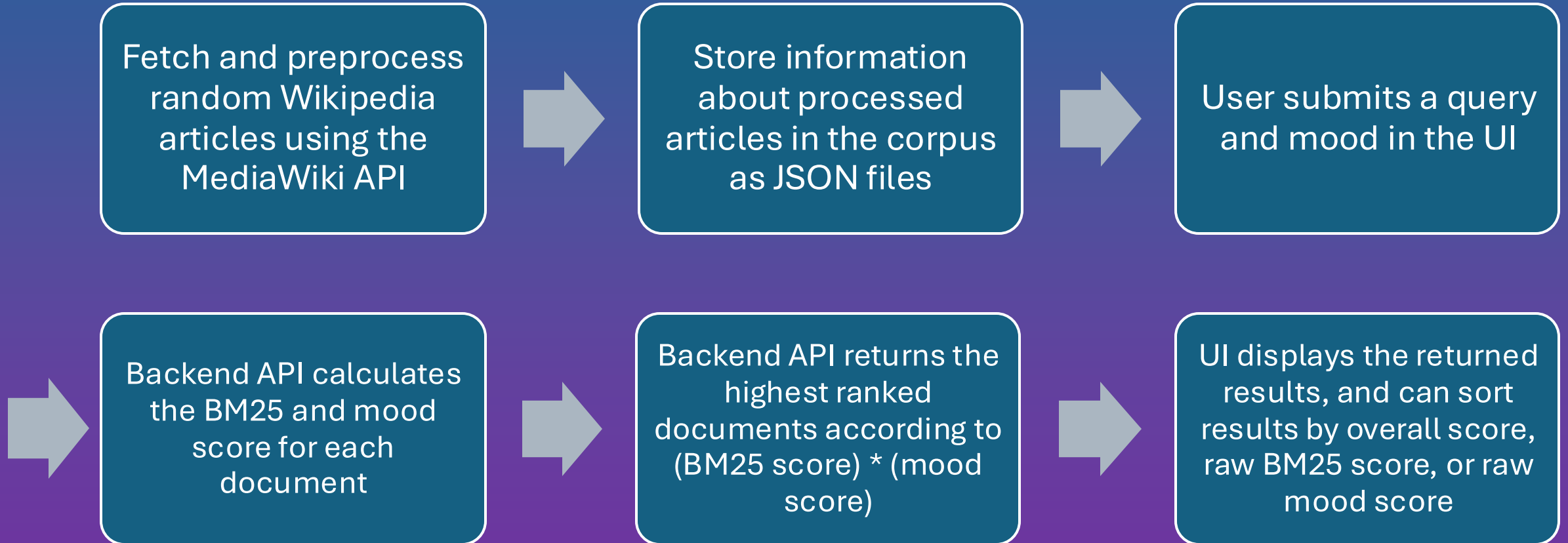# WikiMood

A tool to find Wikipedia articles

Caleb Frye

# Overview of the WikiMood Tool

- Consists of:
  - A backend API which fetches articles, performs document processing, and scores documents based on BM25 and mood score
  - A frontend UI where the user can input a search query and mood, and see articles which match the query and mood
- Motivation:
  - I like browsing random Wikipedia articles when I'm bored
  - I looked online but surprisingly didn't find any website like this
  - I also wanted to add something that made it more advanced than just a simple BM25 search, so I added the mood score

# System Flow

Fetch and preprocess random Wikipedia articles using the MediaWiki API

Store information about processed articles in the corpus as JSON files

User submits a query and mood in the UI

Backend API calculates the BM25 and mood score for each document

Backend API returns the highest ranked documents according to (BM25 score) * (mood score)

UI displays the returned results, and can sort results by overall score, raw BM25 score, or raw mood score

# Algorithm Choice: BM25

- Does a decent job of scoring documents while allowing for fast performance and efficient document storage (with the preprocessing that I've done)

- See the file 'bm25.py' for the functions I use to calculate the BM25 score

- But how to combine it with a mood score?

# Combining BM25 and Mood Score

- First, I defined lists of mood words for each mood in 'constants.py'
- When fetching articles, the system checks each term to see if it is in one of the mood word lists
- If so, it adds +1 to that article's '[mood]_frequency' value
- A document's mood multiplier is then:

  ([mood]_frequency) ÷ (number of mood words in [mood])

- <u>Overall score</u> = (BM25 score) × (mood multiplier)
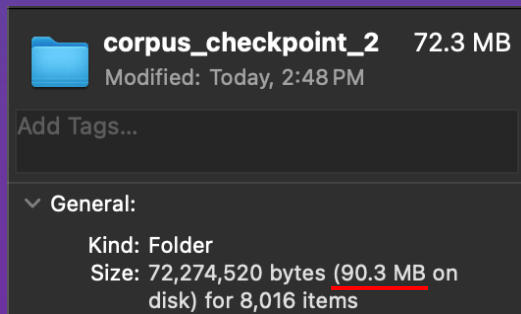
# Combining BM25 and Mood Score cont.

## Example of mood words:

```
mood_words = {
    "Curious": [
        "scientific", "scientifically", "learn", "learning", "explore", "exploring",
        "investigation", "investigative", "investigate", "discovery", "discovering",
        "groundbreaking", "breakthrough", "research", "researching", "curiosity",
        "curious", "hypothesis", "hypotheses", "theory", "theories", "uncover",
        "uncovering", "analyze", "analysis", "question", "questioning", "knowledge",
        "knowledge-seeking", "phenomenon", "phenomena", "experiment", "experimental",
        "observation", "observed"
    ],
    "Uplifting": [
        "inspiring", "inspiration", "hopeful", "hope", "positive", "positivity",
        "motivating", "motivation", "persevere", "perseverance", "encourage",
        "encouragement", "uplifting", "resilience", "resilient", "overcome",
        "overcoming", "strength", "strong", "achieve", "achievement", "fulfilling",
        "fulfilled", "joyful", "joy", "gratitude", "grateful", "success", "successful",
        "optimistic", "optimism", "passion", "passionate", "thriving", "flourish",
        "flourishing", "progress", "progressing"
    ],
    "Entertaining": [
        "movie", "movies", "fascinating", "fascination", "fascinate", "story",
        "stories", "entertaining", "entertainment", "amusing", "amusement",
        "funny", "humor", "dramatic", "drama", "comedy", "comedies", "captivating",
        "captivate", "captivatingly", "engaging", "thrill", "thrilling", "adventure",
        "adventurous", "epic", "characters", "character", "plot", "action", "scenes",
        "spectacle", "charismatic", "charm", "charming", "delightful", "excitement",
        "excite", "exciting", "mystery", "mysteries"
    ],
    "Sad": [
        "death", "deaths", "famine", "genocides", "genocide", "murder", "murders",
        "lost", "loss", "losses", "grief", "grieving", "sorrow", "sorrowful", "mourn",
        "mourning", "tragedy", "tragic", "heartbreak", "heartbroken", "disaster",
```

## Mood term frequency information in the JSON:

```
{} 777.json > ...
   1    {
   2        "title": "Annual plant",
   3        "content-preview": "An annual plant is a plant that co
   4        "url": "https://en.wikipedia.org/wiki/Annual_plant",
   5        "pageid": 777,
   6        "title-terms": {
   7            "annual": 1,
   8            "plant": 1
   9        },
  10        "content-terms": {
  11            "annual": 19,
  12            "plant": 13,
  13            "completes": 1,
  14            "(Other terms etc ... )": ""
  15        },
  16        "title-length": 2,
  17        "content-length": 433,
  18        "curious_frequency": 3,
  19        "uplifting_frequency": 0,
  20        "entertaining_frequency": 0,
  21        "sad_frequency": 0,
  22        "general_frequency": 0
  23    }
```

# The Corpus

- The system stores document information as JSON after stripping out stopwords

- This allows for very efficient document storage and fast search performance

  - At the time of making this, the corpus is 8,000 documents and takes up only 90MB:

Example document JSON file:

# Adding Documents to the Corpus

- I wrote a python script which automatically fetches, preprocesses, and adds new articles to the corpus

- Parameters are configurable in 'constants.py' and explained in the README on Github

- Simply use the script as seen below:

```
[(.venv) → wikipedia-mood-backend-flask git:(main) × python3 fetch_articles.py ]
Fetching 500 Random Wikipedia articles...
Processing Articles:    8%|█              | 38/500 [00:07<01:08,  6.77article/s]
```

# Using the UI

# Using the UI cont.

## WikiMood Mood-Based Search

| information retrieval | Curious ⌄ | Search |

Sort by: **Mood-Weighted Score** | Raw BM25 Score | Mood Multiplier

BM25 score is then weighted using the mood multiplier

When a mood is selected, the articles have mood multipliers

**Search Results**

Show Score Information: ☑
Showing 100 results

### Edward Y. Chang

Mood-Weighted BM25 Score: **11.612**        Raw BM25 Score: **10.451**        Mood Multiplier: **1.111**

*Edward Y. Chang is a computer scientist, academic, and author. He is an adjunct professor of Computer Science at Stanford University, and Visiting Chair Professor of Bioinformatics and Medical Engineering at Asia University, since 2019. Chang is the author of seven books, including Unlocking the Wis...*

### Communication theory

Mood-Weighted BM25 Score: **10.645**        Raw BM25 Score: **3.794**        Mood Multiplier: **2.806**

*Communication theory is a proposed description of communication phenomena, the relationships among them, a storyline describing these relationships, and an argument for these three elements. Communication theory provides a way of talking about and analyzing key events, processes, and commitments tha...*

### Research question

Mood-Weighted BM25 Score: **10.451**        Raw BM25 Score: **1.881**        Mood Multiplier: **5.556**

*A research question is "a question that a research project sets out to answer". Choosing a research question is an essential element of both quantitative and qualitative research. Investigation will require data collection and analysis, and the methodology for this will vary widely. Good research qu...*

# Evaluation

- The BM25 algorithm returns relevant articles decently well, but…
- Since I don't take term positions into account the results can be slightly unrelated
  - For example, searching "tennis player" with any mood other than 'general' returns a lot of articles about video games, since the word "player" appears in them frequently, and the articles that are actually about tennis players don't have high mood scores
- However, I think this is ok! The goal of the system is to help you find articles to read, so if it exposes you to more topics, you're more likely to find something interesting to read
- That said, if I made the system again I would definitely take term positions into account

# Evaluation cont.

- The mood multiplier does affect which documents are ranked higher, especially noticeable for the 'sad' and 'entertaining' moods

- For example, see the top result for searching 'America' with no mood vs. 'sad' mood:

- It's not perfect, but the mood does make a difference

# Conclusion

- Even though it's not perfect, the system does return interesting Wikipedia articles based on your query and mood selection

- The considerations I made in document preprocessing allow for fast searches and efficient document storage

- I know the tool works, because there were many times while I was working on it that I got distracted and started reading the articles returned by the tool

# Thank you!