# Corpus Interface Plan

## First Interface: Data Analysis and Keyword Search

**Purpose:** Enable users to analyze data and search for keywords efficiently.

**Features:**

- **Dataframe Display:** Includes text, likes, and other features from the original dataset.
- **Category Labeling:** Drop-down menu for labeling categories.
- **Text Output Limit:**
  - Display up to **100 tokens** per text entry for concise results.
  - Users can click **"See More"** to view full content.
- **Search Functionality:**
  - Enables users to explore the corpus and retrieve relevant results easily.
- **Visualization Options:**
  - Helps users analyze data patterns with **frequency distributions** and **counts**.
- **Export Function:**
  - Users can download or export search results in **CSV** or **JSON** format.

---

## Second Interface: Interactive Text Chunking with BERT

**Purpose:** Provide an interactive environment for text chunking using BERT or other ML models.

**Features:**

- **User Input:** Interactive text box for user-provided text.
- **Model Output:** Predicts and displays labels for the input text chunk.
- **Real-time Interaction:** Ensures quick model responses for an efficient user experience.

---

## Technical Stack

**Framework:**

- **Streamlit** – For building interactive web interfaces.

**Packages:**

- **Whoosh** – Full-text indexing and search.
- **Pandas** – Data manipulation and analysis.

# Interfaces for Reference:

SHE Corpus: https://genealogies.mvm.ed.ac.uk/webcli/



displaCy Named Entity VIsualizer:
https://demos.explosion.ai/displacy-ent