# Enhancing Domain-Specific Web Crawling with CLIP

Nuradil Zhambyl

20220862

Enhancir
Aleksandra Pshenova

Julien Dupont

Matthew Heffernan

Team #15

20244673

20256086

20256098

# Abstract

We propose a novel method of targeted Web crawling and classifying pages using both textual content and images from each webpage. Our proposed method is based on 'BERT Based Topic-Specific Crawler' [1], and we replace SBERT textual embeddings with CLIP for both text and image.

- **Problematic:** Web page classification only considered the scraped textual content, ignoring the images scraped from each website.
- **Proposed Solution(s):** add CLIP features to the pipeline to process image output and adapted the pipeline for live Web search.

# Motivation

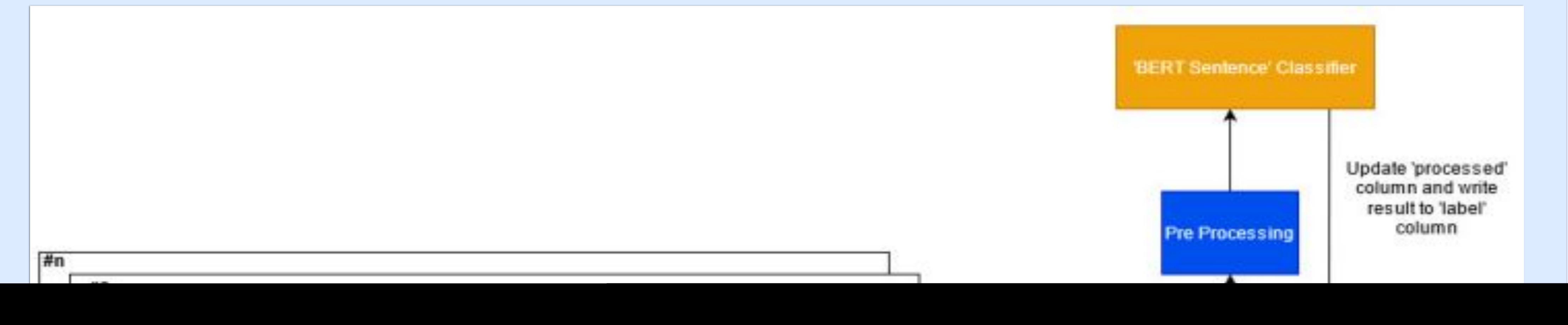
By enhancing the base SBERT-based model with keyword-driven Web search, we aim to create a useful tool with specific application in domains where precise categorization of the Web search results is required.

#### I. Contributions

- Enabled visual embedding via CLIP
- Switched from SBERT to CLIP textual embeddings
- Implemented **image classification** and **filtering** system for CLIP
- Advanced existing system to **keyword-guided live Web search**, instead of hardcoded URL list
  - → created a useful topic-specific search tool
- Evaluated the modified model using chunks embedding and LLM aided text summarization
- Optimized and refactored the original code

## III. Related Work & Method

This project is based on 'BERT Based Topic-Specific Crawler [Tawil, Alqaraleh, in Innovations in Intelligent Systems and Applications Conference (ASYU), 2021].



In order to enable image encoding in the current model, we added an image encoder from CLIP. We also changed the text encoder to CLIP for uniformity.

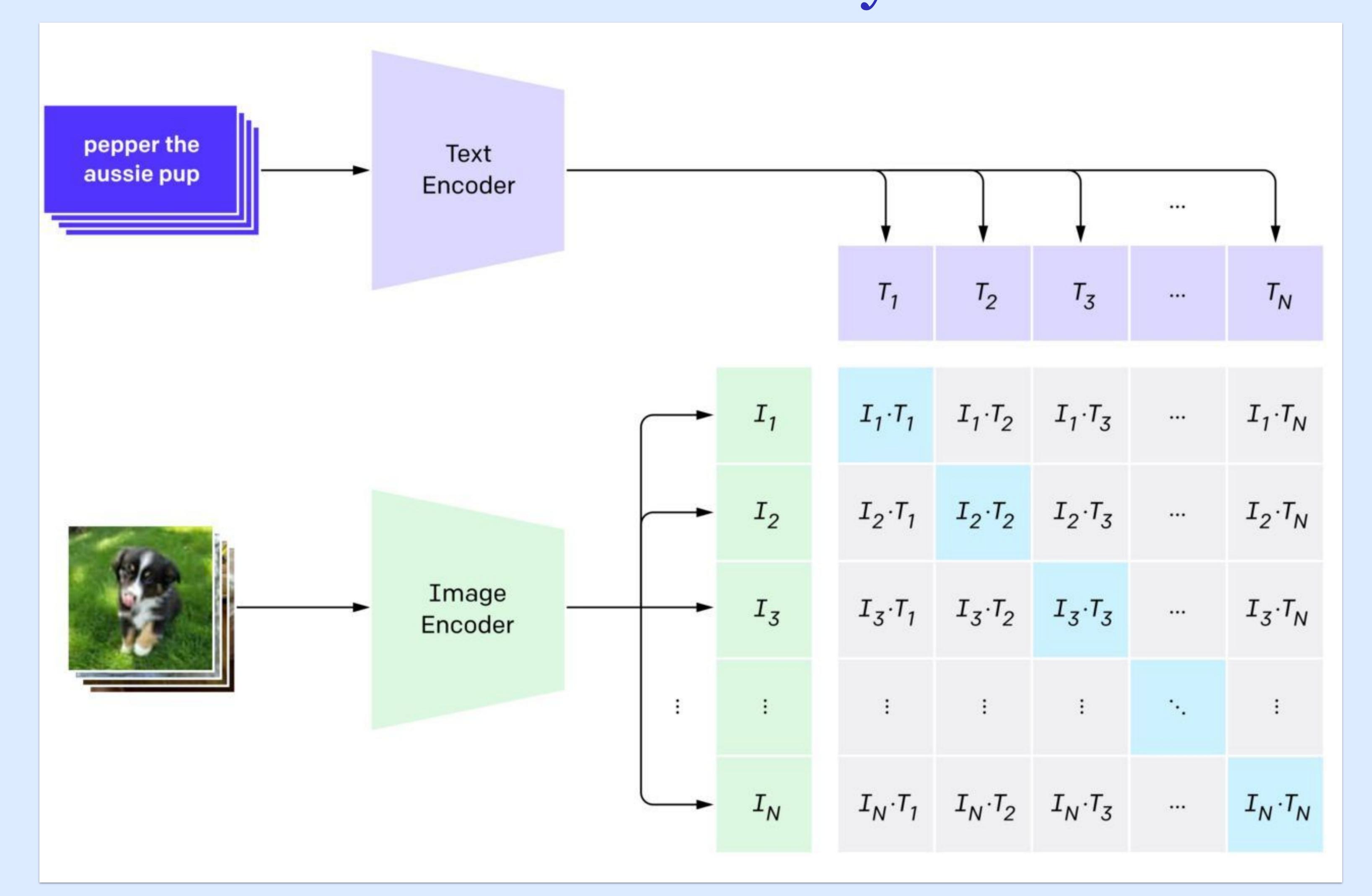


Fig. 3. CLIP Contrastive pre-training approach

To shape the system into a tool with real-world application, we swapped static URL database input with a keyword-based query, scraped the top relevant links, and classifies them. We also enabled an image filter for CLIP to dedupe and filter images irrelevant to the query.

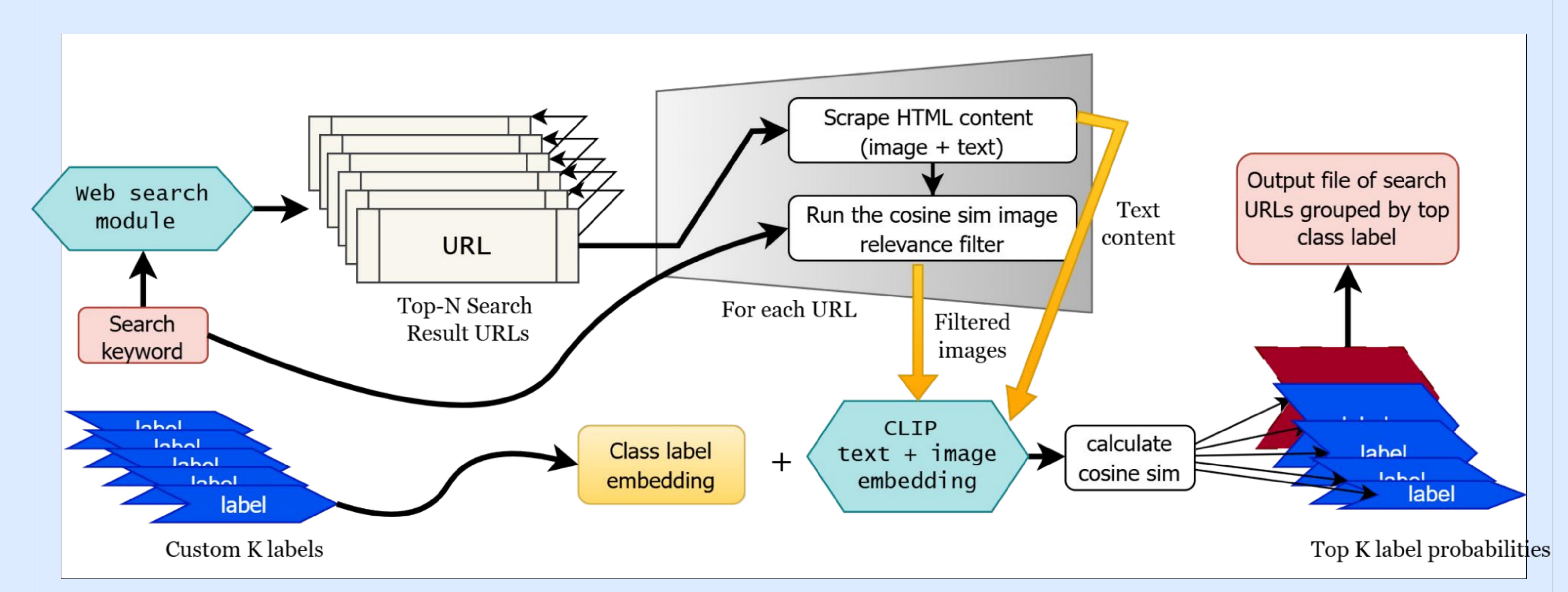


Fig. 4. Our proposed keyword-driven Web search pipeline

## IV. Evaluation & Results

We compare our proposed system results with the baseline model:

Total predictions: | Total predictions: | Total predictions: | Delition | Deli

| Total predictions: 4460 | Tech | History | Business | Politics |
|-------------------------|------|---------|----------|----------|
| Match                   | 1230 | 334     | 200      | 133      |
| Mismatch                | 654  | 516     | 21       | 38       |

Table 1. Number of matching and mismatching samples in the baseline system

| Total predictions: 451 | Tech | History | Business  | Politics |
|------------------------|------|---------|-----------|----------|
| Total madiations       | Took | Tictom  | Durainaga | Dalitiaa |