



crawler

1. fetch with BFS (use HTML parser)
2. extract links to recursively process more pages.
3. build the file structure containing the parent/child link relation.

Pages

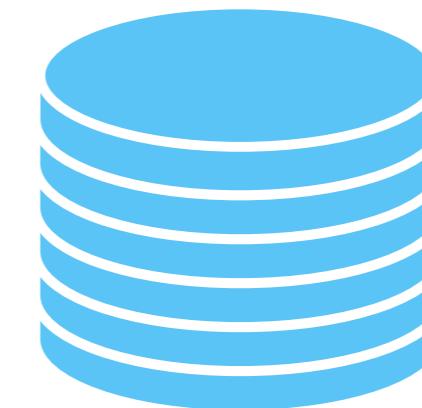
crawler service

indexer service

indexer

1. remove all stop words from the file.
2. transform words into stems using Porter's algorithm (Perform stemming to the words)
3. Store the stems into inverted titles and inverted bodies
4. support phrase search e.g. "hang kong" same type

Database Manager



Map db

Database Schema

- one-to-one map
URL \rightarrow pageID
Word-ID \rightarrow {Page-ID, Freq, pos}
- many-to-many
parent \rightarrow child
representing the link structure
- inverted-titles inverted-bodies
Word-ID \rightarrow {Page-ID, Freq, pos}
- lightweight forward index
page ID \rightarrow {keywords}
- Page properties
PageID \rightarrow {title, URL, last-modified, size, content}

\hookrightarrow top 10 stemmed keywords only.
we can check this before

phase 1 tester program

txt

Bonus

- F 1. relevance feedback
F 2. Allow user to:



- F 3. UI-friendly (Theme: neworphism, colour: bg: white, primary colour: #9966CC, yellow, Anethryse, secondary colour: #ff66ff, font: Poppins)
for highlighting

5. Other search engine func:

1. word pos B
2. result sort by time F

B 6. Ranking (Page rank + Term based)

B 7. do caching, threading (speed)

FIFO, LFU,
LRU?