

Database Design

<p>url-id , id-url</p> <p>Mapping Indexes:</p> <p>URL \leftrightarrow page ID</p> <p>*) URL \rightarrow page ID</p> <p>$\boxed{\text{string}} \rightarrow \text{integer}$</p> <p>int getEntry(string)</p> <p>void addEntry(string, int)</p> <p>*) page ID \rightarrow URL</p> <p>$\boxed{\text{int}} \rightarrow \text{string}$</p> <p>string getEntry(int)</p> <p>void addEntry(int, string)</p>	<p>word-page</p> <p>Inverted Indexes</p> <p>word \rightarrow page ID, freq</p> <p>$\boxed{\text{string}} \rightarrow \text{HashMap} < \text{ID}, \text{freq} >$</p> <p style="margin-left: 150px;">int int</p> <p>HMap getHMap(string)</p> <p>void addFreq(string, int)</p>	<p>page-word</p> <p>Forward Indexes</p> <p>page ID \rightarrow words</p> <p>$\boxed{\text{string}} \rightarrow \text{string}$</p> <p>void addEntry(int, string)</p> <p>string getEntry(int)</p>
<p>page-prop</p> <p>Page Properties</p> <p>pageID \rightarrow { Title, ^{last}modified, ^{content}length }</p> <p>$\boxed{\text{int}} \rightarrow \text{HashMap} < \text{tag}, \text{content} >$</p> <p style="margin-left: 150px;">string string</p> <p>void addProps(int, HMap)</p> <p>Vector<string> getResult(int)</p>	<p>parent-child , child-parent</p> <p>Link Based Indexes</p> <p>pageID (parent) \leftrightarrow pageID (child)</p> <p>$\boxed{\text{int}} \rightarrow \text{ArrayList} < \text{IDs} >$</p> <p style="margin-left: 150px;">int int</p> <p>void addEntry(int, int)</p> <p>ArrayList<int> getID(int)</p>	<p>url-url</p> <p>Mapping Indexes:</p> <p>URL \rightarrow URL</p> <p>$\boxed{\text{string}} \rightarrow \text{string}$</p> <p>string getURL(string)</p> <p>void addEntry(string, string)</p>

Above is the sketch of the database used in my search engine project. Every database works independently in a sense that every database does not contain any of the 8 databases as the value. Below is the explanation for each databases design:

1. URL \rightarrow PageID [Mapping Index] = this database contains every URL mapped with an unique ID to it

ex.

http://www.cse.uh.hk/	\rightarrow 0
http://hkust.edu.hk/news	\rightarrow 1
\vdots	\vdots

2. PageID \rightarrow URL [Mapping Index] = this database is the reversed of the first database design, which maps an ID into an URL

ex.

0	\rightarrow http://www.cse.uh.hk/
1	\rightarrow http://hkust.edu.hk/news
\vdots	\vdots

3. word \rightarrow (PageID, frequency) [Inverted Index] = this database maps every word into every pageID that has the words in it, together by how often it appears in the website as the frequency

ex.

add	\rightarrow { 0=10, 1=22, ... }
the	\rightarrow { 1=2, 10=107, ... }
\vdots	\vdots

4. PageID \rightarrow words [Forward Index] = this database maps pageID to every words that reside inside the page

ex.

0	\rightarrow "The Department of ..."
1	\rightarrow "News Hong Kong ..."
\vdots	\vdots

5. PageID \rightarrow page properties = this database maps pageIDs into the property of the page such as the page title, the content length, and the last modified date.

ex.

0	\rightarrow { Title = ..., Last-Modified = ..., ... }
1	\rightarrow { Title = ..., Last-Modified = ..., ... }
\vdots	\vdots

6. Parent → Child [Link Based Index] = this database contains every parent-child relation of every indexed urls

ex.

0	→ [1, 2, 3, ...]
1	→ [1, 3, 17, 18, ...]
⋮	⋮

7. Child → Parent [Link Based Index] = this database contains every child-parent relation of every indexed urls

ex.

0	→ [3, 4, 6, ...]
1	→ [0, 1, 2, ...]
⋮	⋮

8. URL → URL [Mapping Index] = this database maps every URL into its corresponding actual URL. Actual URL means that if there is any redirection from the website (status code= 3xx), then the actual URL is the link inside the Location properties of a HTML response

ex.

http://www.cse.uif.hk/	→ https://cse.hkust.edu.hk/
⋮	⋮