

## Phase 1 RocksDB Database Scheme

RocksDB name	Mapping
wordToWordID	<p>word &lt;-&gt; wordID</p> <ul style="list-style-type: none"> <li>word: <code>String</code></li> <li>wordID: <code>int</code></li> <li>Uses the Java hashCode() function to hash a keyword to a unique wordID</li> <li>Contains all words found in all documents indexed by the crawler</li> </ul> <p>Design Explanation/notes:</p> <ul style="list-style-type: none"> <li>Before the word is inserted into the database, it is first stemmed using Porter's Algorithm. The word is not the original keyword that appears in the document, but instead the stemmed version of the word</li> <li>During implementation for queries, we plan to stem all words using the same algorithm to ensure consistency in retrieving the correct wordID from the same words.</li> </ul>
urlToPageID	<p>url &lt;-&gt; pageID</p> <ul style="list-style-type: none"> <li>Url: <code>String</code></li> <li>pageID: <code>int</code></li> <li>Uses the Java hashCode() function to hash a url to a unique pageID</li> <li>Contains all urls found in all documents indexed by the crawler</li> </ul>
pageInfo	<p>pageID &lt;-&gt; Page</p> <ul style="list-style-type: none"> <li>pageID: <code>int</code></li> <li>Page: a <code>Page</code> object <ul style="list-style-type: none"> <li><u>The <code>Page</code> class contains the following fields</u></li> <li>page title: <code>String</code></li> <li>url: <code>String</code></li> <li>Last modified date: <code>String</code></li> <li>Page size: <code>int</code></li> <li>Child links: <code>ArrayList&lt;String&gt;</code></li> </ul> </li> </ul> <p>Design Explanation/notes:</p> <ul style="list-style-type: none"> <li>The Page class is used to store meta information about each page, so we added this RocksDB to make it easy to access this information directly from the pageID.</li> </ul>

invertedIndex	<p>wordID -&gt; pageList(a list of pages that the word appears in)</p> <ul style="list-style-type: none"> <li>• wordID: <code>int</code></li> <li>• pageList: <code>HashSet&lt;Integer&gt;</code></li> </ul>
forwardIndex	<p>pageID -&gt; wordList(wordIDs of words that appear in the webpage, and the position of the words in that page)</p> <ul style="list-style-type: none"> <li>• pageID: <code>int</code></li> <li>• wordMap: <code>HashMap&lt;Integer, ArrayList&lt;Integer&gt;&gt;</code> <ul style="list-style-type: none"> <li>○ wordMap uses the wordID as the key. The value that correspond to the key is an ArrayList that contains the positions of that wordID for that particular pageID</li> </ul> </li> </ul> <p>Design Explanation/notes:</p> <ul style="list-style-type: none"> <li>• We stored the positions of the word along with the word itself. By retrieving the words in addition to their positions, it is easier to support searching by phrases</li> </ul>