

Assignment 1

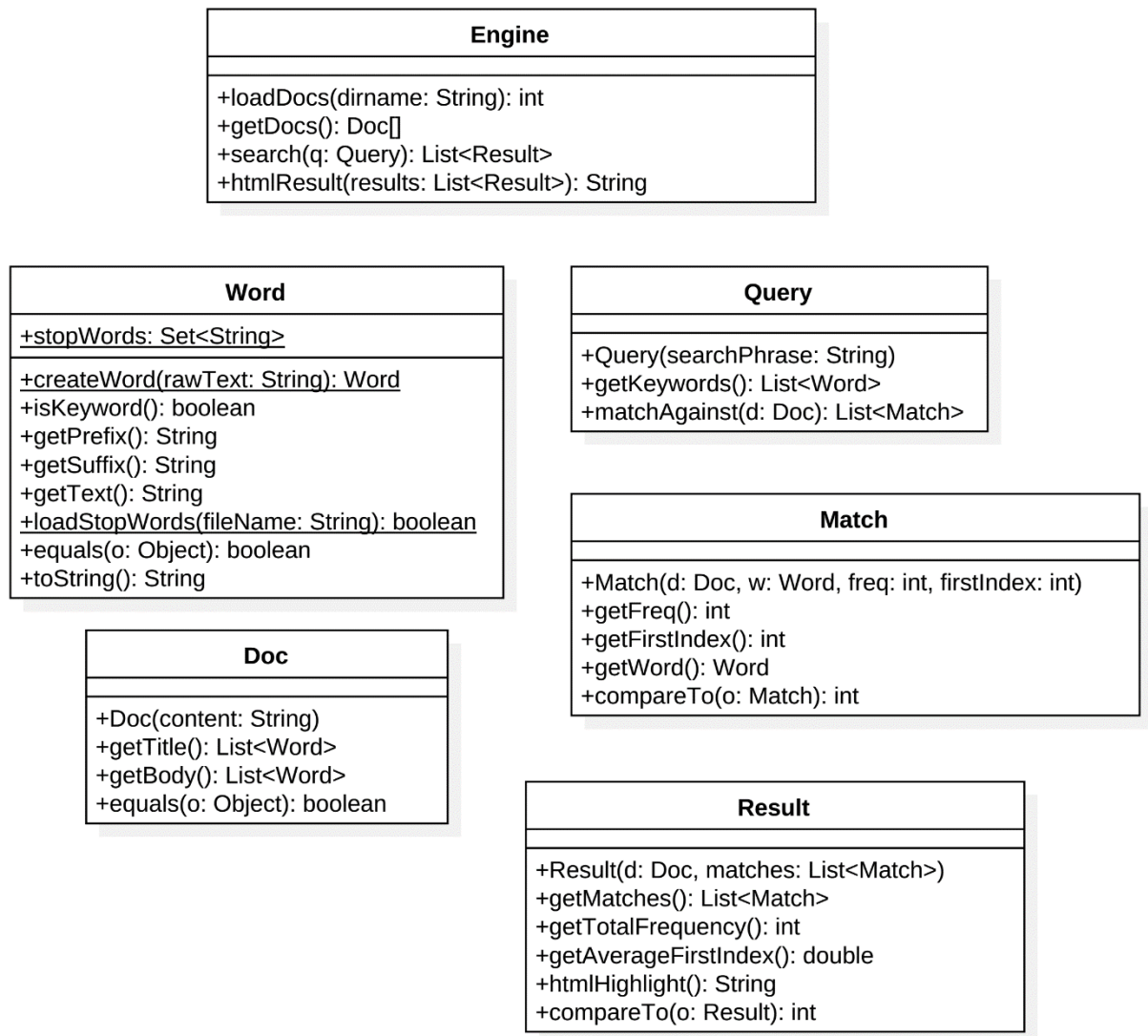
1. Requirements

In this assignment, you're going to build a simple document search engine. This search engine can run queries against a collection of documents. It only searches for the words in queries (called keywords) that aren't in stopwords list. The result is ranked by the frequency of keywords.

This app has these main features:

- Load documents from text files.
- Search for documents by a key phrase.
- Rank search results based on several criteria.
- Generate HTML search results with highlighted keywords.

Program must be written in **Java 8**. The program must contain these classes and methods:



Class	Description
-------	-------------

Word	<p>When you split a text by the space (" "), you get an array of strings which do not contain space. Some of these strings may be empty. <code>Word</code> is the basic class to store each of these strings. Word objects can later be merged together to obtain the original text.</p> <ul style="list-style-type: none"> • The raw text of a word may contain the actual word (the text part) and its surrounding parts (prefix and suffix). • Words are divided into valid and invalid words. • Valid words are further divided into keywords and stop words (function words). <p>The text part of a valid word must have at least one letter and is composed of letters, hyphens and apostrophes. Examples of valid text parts: <code>Software</code>, <code>up-to-date</code>, <code>don't</code>, <code>doesn't</code>. There is an exception to the above rule. In case the text part ends with <code>'s</code> or <code>'d</code>, it should be considered as part of the word's suffix. To clarify this, the string <code>(other's)</code> should be converted into a <code>Word</code> object which belongs to the valid words category, with the prefix <code>(</code>, the text part <code>other</code> and the suffix <code>'s)</code>. Empty string and non-empty strings which do not have the text part which follow the above rules are invalid words.</p> <p>Prefix and suffix of a word are composed of characters which are not alphanumeric. Therefore, the raw texts such as <code>,se2021.</code> and <code>se,se,</code> should be rendered as invalid words. Invalid words will have empty prefix, empty suffix and text part is the raw text.</p> <p>For example, the phrase <code>in book's "Chapter 5", I've</code> can be split into 4 sequences <code>in</code>, <code>book's</code>, <code>"Chapter</code> and <code>5"</code>. The 1st sequence is the word <code>in</code>, a stop word with empty prefix and suffix. The 2nd sequence is the word <code>book</code> with empty prefix and the suffix <code>'s</code>. The 3rd sequence is the word <code>Chapter</code>, a keyword with the prefix <code>"</code> and empty suffix. The 4th sequence is an invalid word. It should be stored as a <code>Word</code> object with empty prefix, empty suffix and the text part of <code>5"</code>.</p> <p>Public attributes</p> <pre>public static Set<String> stopWords;</pre> <p>A set of stop words loaded by the <code>loadStopWords()</code> method.</p> <p>Public methods</p> <pre>boolean isKeyword()</pre> <p>Returns true if the word is a keyword.</p> <pre>public String getPrefix()</pre> <p>Returns the prefix part of the word.</p> <pre>public String getSuffix()</pre> <p>Returns the suffix part of the word.</p> <pre>public String getText()</pre> <p>Returns the text part of the word.</p> <pre>public boolean equals(Object o)</pre> <p>Two words are considered equal if their text parts are equal, case-insensitively.</p> <pre>public String toString()</pre> <p>Returns the raw text of the word.</p>
------	--

	<pre>public static Word createWord(String rawText)</pre> <p>Construct and return a complete <code>Word</code> object from raw text.</p> <pre>public static boolean loadStopWords(String fileName)</pre> <p>Load stop words into the set <code>Word.stopWords</code> from the text file whose name is specified by <code>fileName</code> (which resides under the project's root directory). This method returns <code>true</code> if the task is completed successfully and <code>false</code> otherwise.</p>
Doc	<p><code>Doc</code> is the class to represent a document which has a title and a body. The title and body of a document are lists of <code>Word</code> objects.</p> <p>Public methods</p> <pre>public Doc(String content)</pre> <p>A constructor which receives the raw text of a document and extracts the title and body parts from that. Documents are provided as text files (<code>.txt</code>) in the <code>docs</code> directory under the project's root directory. To reduce the difficulty of this assignment, each text file contain two lines. The first line is the title and the second line is the body.</p> <pre>public List<Word> getTitle()</pre> <p>Returns the document's title as a list of <code>Word</code> objects.</p> <pre>public List<Word> getBody()</pre> <p>Returns the document's body as a list of <code>Word</code> objects.</p> <pre>public boolean equals(Object o)</pre> <p>Two <code>Doc</code> objects are equal if their titles and bodies contain the same words in the same order. To determine if two words are equal, use the <code>equals()</code> method from the <code>Word</code> class.</p>
Query	<p><code>Query</code> is the class to represent a user's search query. A <code>Query</code> object should store a list of keywords internally.</p> <p>Public methods</p> <pre>public Query(String searchPhrase)</pre> <p>A constructor which receives the raw search phrase from user, then extract only keywords from it.</p> <pre>public List<Word> getKeywords()</pre> <p>Returns a list of the query's keywords in the same order as they appear in the raw search phrase.</p> <pre>public List<Match> matchAgainst(Doc d)</pre> <p>Returns a list of matches against the input document. Sort matches by position where the keyword first appears in the document. See the <code>Match</code> class for more information about search matches.</p>
Match	<p>A <code>Match</code> represents a situation in which a <code>Doc</code> contains a <code>Word</code>. The search engine's job is to find all documents that are related to a <code>Query</code>. Matches are the building blocks of the relationship between documents and a search query. Each <code>Match</code> object only stores information about which <code>Doc</code> contains which <code>Word</code> but also keeps the number of times that the <code>Word</code> appears in the <code>Doc</code> (frequency) as</p>

	<p>well as the first position which the <code>Word</code> appears (first index). When matching a document against a word, the title and body of the document should be combined into a single list of words, with the title placed before the body.</p> <p>This class must implement the <code>Comparable<Match></code> interface.</p> <p>Public methods</p> <pre>public Match(Doc d, Word w, int freq, int firstIndex)</pre> <p>A constructor to initialize a <code>Match</code> object with the document, the word, the frequency of the word in the document and the first position of the word in the document.</p> <pre>public int getFreq()</pre> <p>Returns the frequency of the match (as explained above).</p> <pre>public int getFirstIndex()</pre> <p>Returns the first index of the match (as explained above).</p> <pre>public int compareTo(Match o)</pre> <p>Compare <code>this</code> with another <code>Match</code> object by the first index. This method obeys the standard behavior specified by Java. <code>Match</code> object <code>A</code> is greater than <code>Match</code> object <code>B</code> if the first index of <code>A</code> is greater than the first index of <code>B</code>.</p>
Result	<p>For a <code>Query</code>, the search engine may find a number of related documents. Each document found is represented by a <code>Result</code> object. A <code>Result</code> object stores information about a related a document, a list of matches found in that document and also three derived properties:</p> <ul style="list-style-type: none"> • match count: the number of matches, indicated by the size of the list of matches. • total frequency: the sum of all frequencies of the matches. • average first index: the average of the first indexes of the matches. <p>This class must implement the <code>Comparable<Result></code> interface.</p> <p>Public methods</p> <pre>public Result(Doc d, List<Match> matches)</pre> <p>A constructor to initialize a <code>Result</code> object with the related document and the list of matches.</p> <pre>public List<Match> getMatches()</pre> <p>The method's name explains itself.</p> <pre>public int getTotalFrequency()</pre> <p>The method's name explains itself.</p> <pre>public double getAverageFirstIndex()</pre> <p>The method's name explains itself.</p> <pre>public String htmlHighlight()</pre> <p>Highlight the matched words in the document using HTML markups. For a matched word in the document's title, put the tag <code><u></code> and <code></u></code> around the word's text part (the <code><u></code> tag should not affect the word's prefix and suffix). For a</p>

	<p>matched word in the document's body, surround the word's text part with the tag <code></code> and <code></code>.</p> <pre>public int compareTo(Result o)</pre> <p>These are criteria to determine if <code>Result A</code> is greater than <code>Result B</code> (in descending order of priority):</p> <ul style="list-style-type: none"> • <code>A</code> has greater <i>match count</i> than <code>B</code> • <code>A</code> has greater <i>total frequency</i> than <code>B</code> • <code>A</code> has lower <i>average first index</i> than <code>B</code>
Engine	<p>This class represents the search engine.</p> <p>Public methods</p> <pre>public int loadDocs(String dirname)</pre> <p>Loads the documents from the folder specified by <code>dirname</code> (which resides under the project's root folder) and returns the number of documents loaded. Refer to the <code>Doc</code> class for more information about a <code>Doc</code> object.</p> <pre>public Doc[] getDocs()</pre> <p>Returns an array of documents in the original order.</p> <pre>public List<Result> search(Query q)</pre> <p>Performs the search function of the engine. Returns a list of sorted search results. Refer to the classes above to know the expected search results.</p> <pre>public String htmlResult(List<Result> results)</pre> <p>Converts a list of search results into HTML format. The output of this method is the output of <code>Result.htmlHighlight()</code> combined together (without any delimiter). Refer to the 3rd line of the file <code>testCases.html</code> for a specific example.</p>

You are provided with a project folder containing:

- A folder named `docs` which contains 10 documents saved as text files.
- A `src` folder which contains a package named `engine`. Inside the package, `App.java` has been provided to check your solution locally. This java program checks your solution against an incomplete set of test cases.
- A file named `stopwords.txt` containing the required stop words that you have to load into the search engine.
- A file named `testCases.html` which is used by `App.java`

(*) You shouldn't modify any of the provided files.

Apart from the required public attributes and methods, you are free to add more attributes and methods as you see fit. Refer to the test cases provided in the `App.java` program for specific examples of method outputs.

2. Submission

The provided java source files are put in the `engine` package folder. Rename the `engine` package folder into `a1_sid` where `sid` is your student ID. Put the package folder into a zip file so that when the zip file is opened, you'll see the `a1_sid` folder. **Failure in naming the file or use incorrect file structure as shown will result in no marks being given.**

NO PLAGIARISM: If plagiarism is detected, 0 mark will be given!