

Implement a Java application that builds an Inverted Index

What's an inverted index

An inverted index is an index data structure storing a mapping from content, such as words or numbers, to its locations in a document or a set of documents. The purpose of an inverted index is to allow fast full text searches.

You have millions of documents or webpages or images anything that we may need to retrieve later. Here's an example of an inverted index, matching a word with pages in a book that contain that particular word.

Index

Page numbers in **bold face** refer to key term definitions

Page numbers in *italics* refer to images or diagrams

Page numbers followed by a "t" indicate a table

A

absolute temperature scale, **350–351**

absolute zero, **351**

acceleration of gravity, **A.231**

accuracy, **A.3**

acetic acid (CH_3COOH)

 buffers, **575–576**, **581–582**

 conjugate acid-base pairs, **540**

 ionization constant, **523**, **554t**

 manufacture of, **431**

 titrations, **590–592**

 as weak acid, **144t**, **145**, **551–552**

acid-base pairs, conjugate, **540–544**

acid-base reactions, **538**

 autocatalysis of water, **545–547**

 gas-forming exchange, **153–151**

 net ionic equations for, **148–150**

 neutralization, **146–150**, **561–566t**

 of salts, **146–151**, **561–566**

acid-base titrations, **587–594**

acidic solutions, **546**, **650–652**, **A.32t–A.33t**

acid ionization constant expressions, **550–551**

acid ionization constants (K_a) *See* ionization constants,

 acids (K_a)

acidosis, **576**

acids, **143**. *See also* acid-base reactions; ionization

 constants; acids (K_a); specific entries, e.g. carboxylic acids,

 lactic acids

 Brønsted-Lowry concept, **538–544**

 buffer solutions, **575–585**

 conjugate acid-base pairs, **540–544**

 equilibrium constants, **473t**

 ionization constants, **550–561**, **A.28t**

 Lewis, **566–568**

 organic, **544**

 pH scale, **547–550**

 properties, **143–145**

 solubility of salts, **557–598**

 solutions, **546**, **650–652**, **A.32t–A.33t**

 strengths, **145–146**, **555–556**

 titrations, **587–594**

 water's role, **540**

actinides, **55**, **250–251**

activated complex, **433**

activation energies (E_a), **434**, **438–440**, **445**, **447**, **449–450**

active sites, **449**

activities, **406**, **547**, **703–710**

activity series, metal, **159–160**

actual yields, **121**

addition, **45**, **A.6**, **A.9**

addition, significant figures in, **A.6**

air, **342–343**, **366–370**, **380–381**, **706**

alkyl groups, **70–71**

alcohols, **64**, **505–507**

aldehydes, **273–279**

alkali metals, **55**, **106**

alkaline batteries, **670**

alkaline earth metals, **55**

alkaline fuel cells, **674**

alkalosis, **576**

alkanes, **68–71**, **277–278**, **A.25–A.26**

alkenes, **280–283**, **A.26–A.27**

alkyl group, **70–71**, **A.25–A.26**

alkynes, **281**, **A.27**

allotropes, **22–24**, **206**, **403–405**

alpha particles, **38–39**, **693–694**, **697**, **699–700**

alpha radiation, **693**

alpha rays, **36–37**

aluminum (Al), **7**, **8t**, **103**, **634–635**, **682**

amines, **544–545**

ammonia (NH_3)

 amino, **545**

 Brønsted-Lowry base, **538–539**

 complex ions, **567**

 ionization constant, **554t**, **561**

 standard molar enthalpy of formation, **210t**

 structure, **64t**

 synthesis, **107–108**, **462**, **494–495**

 VSEPR model, **312–313**

 as weak base, **145**, **539**

ammonium ions, **74**, **77**

amorphous solids, **390**

ampheres (A), **658**

antiprotic species, **540**

amphoteric metal hydroxides, **567**, **602**

amplitude, **223**

amu (atomic mass unit), **46**

analytical chemistry, **114**

Anderson, Carl, **696**

angular geometries, **314**

anions, **72**, **76**, **78**, **257–258**

anodes, **37**, **654**. *See also* electrochemistry

anodic inhibition, **683–684**

antibonding molecular orbitals, **298**

antimatter/antileptons, **696**

aqueous equilibria. *See also* acid-base titrations; buffers

 factors affecting solubility, **577–602**

 precipitation, **603–604**

 solubility product constant, **594–597**

aqueous solutions (aq), **100**. *See also* buffer solutions;

 solutions

 electrolysis, **675–678**

 electrolytes, **82–83**, **136**, **527–528**, **634–635**

 ionic compounds, **136–139**

 molarity, **166–168**

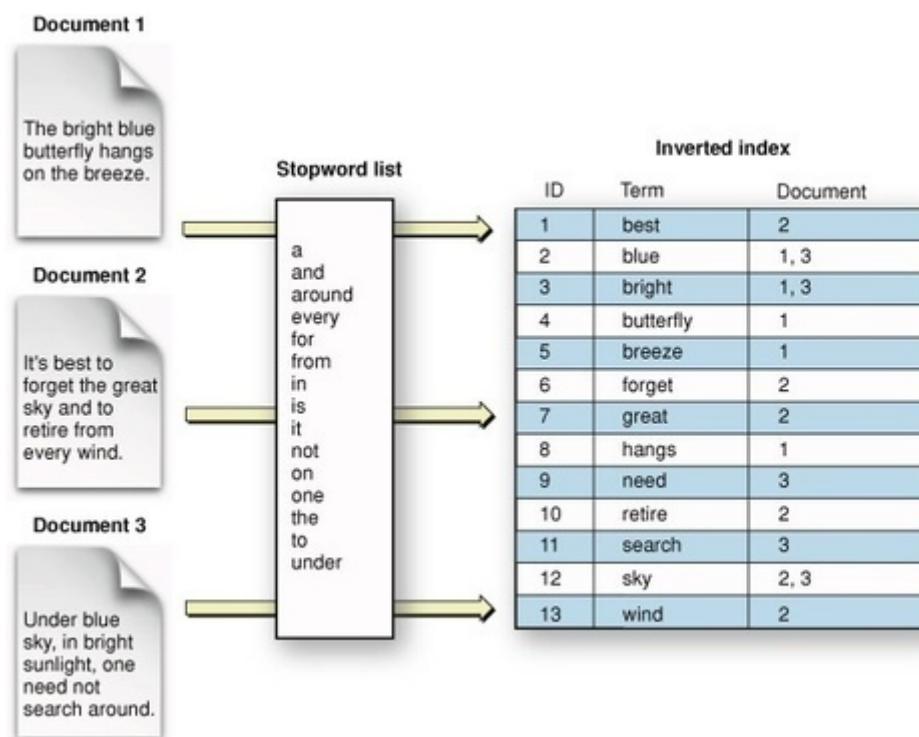
 standard reduction potentials in, **A.32t–A.34t**

 stoichiometric relationships in, **157**

This is an example from some random textbook. If you want to find something specific, say **activation energies**, you will open up this index and the inverted index will tell you the page numbers where that word is explained in a big bulk of a thousand pages. As a human, if you were to perform a regular linear search, you'll take hours to reach that page. But now it was hardly a matter of seconds.

How to build an inverted index?

Building an inverted index for maintaining any kind of searching system requires you to perform a series of steps while parsing the pages or documents. Let's have a walkthrough while constructing our own search engine.



Resources: [wikipedia.org](https://en.wikipedia.org), [quora.com](https://www.quora.com)

Homework requirements

I want to create a search engine Java application for the documents on my computer.

1. Fetching the documents

Write a Java application which takes a file path as input. This file path is a directory containing all the documents that need to be indexed.

The application must go through all the documents, parse them and build the inverted index.

Enter file path where documents are located:

[user input]: /home/user/MyDocuments

> Indexing process started...

> Indexing Algebra.txt

> Indexing Algorithms.txt

> Indexing Biomechanics.pdf

1. Fetching the Document

Focus on reading text file (.txt) documents.

Bonus: Parse doc and pdf documents. You'll need to use some libraries to retrieve their text.

2. Removing the Stop Words

Consider the above paragraph. What were the important words we may be looking for? "text", "libraries", "doc", "pdf", "retrieve", "parse". But most of the other words are just a waste. We denote the most occurring words as "**stop words**". This includes "I", "the", "we", "is", "an", "to".

See the attached **stopwords.txt** files which contain the stop words.

The Java application needs remove the stop words so that they are not to be included in the inverted index. This is to reduce the index size.

It should support English words.

Bonus: Support Romanian words

3. Stem to the Root Word

Then comes Stemming. Now whenever I want to search for "retrieval", I want to see a document that has information about it. But the word present in the document is called "retrieve" instead of "retrieval". To relate the both words, I'll chop some part of each and every word I read so that I could get the "root word". Retrieve may become "retriev". So will "retrieval". We have to be sure about the rules we use to chop the words. There are standard tools for performing this like "Porter's Stemmer". You can play around with a porter stemmer here : [Porter Stemmer Online](#)

Bonus: Support Romanian words

4. Record Document IDs

Now get ready for the main task - Indexing.

Every document I have has got an unique document name. We'll take each non-stop word, we will stem and we will save it in the inverted index.

Example:

retriev ==> document1

If the same word is in some other document, we'll store:

retriev ==> document1

retriev ==> document2

But very soon I've to combine them in a single list

retriev ==> document1 AND document2

This needs some specific data structures that simplify your job.

We need to further improve our algorithm by storing how many times did the word occur in each document.

retriev ==> [document1|2 times] AND [document2|5 times]

Using the number of occurrences, we can rank the search results, displaying document2 before document1.

5. Search capability

At this point, the indexing process has finished. We should have an in-memory data structure which holds the inverted index for our documents.

The inverted index will contain mappings between a single word to the documents where they appear.

We want to be able to search **multiple** words. The results should contain the document names that contains **all** the desired words (intersection).

See example below:

Enter file path where documents are located:

[user input]: /home/user/MyDocuments

> Indexing...

> Indexing Algebra.txt

> Indexing Algorithms.txt

> Indexing Biomechanics.pdf

> Indexing Precalculus.pdf

> Indexing Statistics.txt

> Indexing Mathematical Methods.doc

> Indexing process finished.

Enter search words:

[user input]: **abstract algebra**

> Your query matches the following documents, presented in a ranked order:

> /home/user/MyDocuments/Precalculus.pdf

> /home/user/MyDocuments/Statistics.txt

> /home/user/MyDocuments/Mathematical Methods.doc