

# Stack Overflow Search

Shruthi Raghuraman

Github Source Code: <https://github.ccs.neu.edu/sraghuraman1/6200-StackExchangeSearch>

## Introduction:

StackOverflow is an extremely valuable tool in knowledge sharing between programmers. From their website, it is an “open community for anyone who codes” that focuses on answering tech questions and much more. It mainly acts as a question and answer site for programmers and other technology enthusiasts.

As we learned this semester, community-based question answering systems provide a lot of advantages. One benefit is that users can find answers to complex or obscure questions. Since the answers are from other people and not pieced together by algorithms, the search can take advantage of an archive of previous questions and answers. However, this comes with a few complications. One of the most evident problems with these systems is that questions are not guaranteed an answer and some of those answers may be wrong.

It is interesting to note that a lot of these complications are mitigated by measures taken by the StackOverflow ecosystem. With community monitoring, StackOverflow ensures content quality by making sure that all questions on there are tightly focused on a specific issue. Any questions that are vague/non-technical are rejected and closed. Duplicate questions are linked to other posts. Lastly, answers are also monitored with comments, votes, and other discussion. These functionalities make this site quite useful and relevant to the needs of its target audience.

Even though this ecosystem helps with curating what questions and answers get published on the site, the search functionality on StackOverflow is less than ideal and ultimately holds the site back from fully utilizing the formidable aspects of its community driven ecosystem. As a user, I have experienced quite a few issues with StackOverflow’s search functions. These problems make it incredibly difficult to find the right information in a timely manner. Although the dataset on the site is comprehensive, users have noticed StackOverflow’s search does not retrieve relevant results on queries asked. By going through community postings on this topic, I discovered that the current ranking system of the results places extreme importance on defining relevance based on text within the title and not enough on showing answered/active questions. The following proposal improves StackOverflow’s search by incorporating a new ranking and retrieval system that measures relevance based on a comprehensive and processed text snippet from a post’s metadata rather than just the title.

## Problem Identified:

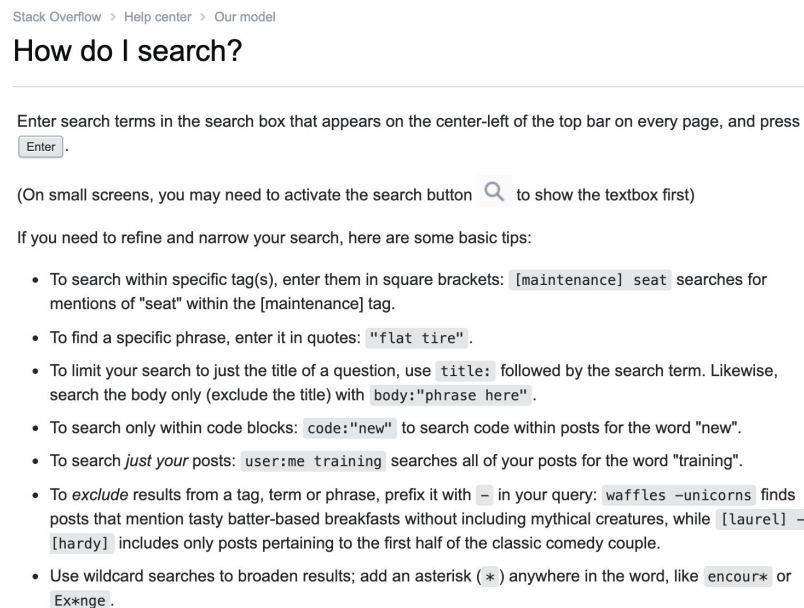
When using a search engine, users expect an interface that is easy to navigate as well as for the top ranked results of the search to give an appropriate solution to the question asked. StackOverflow has a plethora of viable solutions to many programmers’ common and obscure questions. Even though the information a

user needs is typically available on the site, the website tends to fall short by giving a different and less viable solution to the question posed by the user.

Problems with the current search can be classified into two components: the user interface on StackOverflow's web page and the retrieval from the system.

## 1. User Experience with Search Bar:

StackOverflow's basic search engine has quite a few features that are both useful and cumbersome for new users. The instructions from the website's Help Center for the basic search (not including the more verbose Advanced Search Options) are shown below.



Although this is a well documented constraint-based search, the amount of very specific instructions acts as a very high barrier of entry when using the search engine. This has resulted in a general consensus among users to not utilize the website's search bar and, instead, use Google's search bar knowing that Stack Overflow will be a top ranked choice when asking code specific questions using Google. A question on StackExchange asking "What are some good tips for searching Stack Overflow?" from a beginner user garnered the following responses:

- "Use Google to search Stackoverflow."
- "You can always use [Google](#), which is really Stack Overflow's home page."
- "In all honesty, your best bet is to use google instead when searching."
- "Search now *HEAVILY* weights title in the results, since people seemed to really like that approach (used on the /ask page, which searched exclusively on title alone)."

Some users suggested a good (but roundabout) way to have the information need met would be to post a question on the site and see if it is marked as duplicate based on the title.

Although the suggestion to use Google to query StackOverflow is a solution, I believe that improving the StackOverflow search engine is beneficial to the community. From the user interface point of view, the engine could benefit from this proposal of implementing a less constraint-based search.

## 2. Retrieval from a User Query:

Another major problem with the current search engine is how it retrieves relevant information. For example, if I want to start learning about Android Development, and I would like to do a basic search such as “passing in android activity intent” into the StackOverflow search engine, the results are only loosely related to the topic. Although some of the words searched appear in the titles of the search results, the content of these posts is complicated and greatly deviates from the user’s information need. It is frustrating that the majority of the words searched appear in the title, while the corresponding posts and answers do not satisfy the question asked.

The example for the above query is shown below:

74  
votes

7  
answers

**Q: Passing android Bitmap Data within activity using Intent in Android**  
I have a Bitmap variable named bmp in Activity1 , and I want to send the bitmap to Activity2 Following is the code I use to pass it with the **intent**. **Intent** in1 = new **Intent**(this, Activity2.class ... ); in1.putExtra("image",bmp); startActivity(in1); And in Activity2 I try to access the bitmap using the following code Bundle ex = getIntent().getExtras(); Bitmap bmp2 = ex.getParcelable("image ...  
android android-intent bitmap  
asked Jun 13 '12 by adi.zean

0  
votes

0  
answers

**Q: Send data from Android Activity to Xamarin Form**  
(this, typeof(?)); myIntent.PutExtra("greeting", "Hello from the Second **Activity!**"); SetResult(Result.Ok, myIntent); Finish(); } I am creating the **Intent** before Finish button in **android activity** but what type of class should I use in **Intent** and in which event I will get the result of **Intent**. ... I want to send data from **android activity** to Xamarin form. I have one content page, from where I am opening the **android activity** and **passing** the data from Xamarin app to **Android activity**. Now I need ...  
xamarin xamarin.forms xamarin.android  
asked May 23 '19 by Imrankhan

1  
vote

1  
answer

**Q: passing string to another activity without using intent extras in Android**  
Could I pass a string to another **activity in Android** without using intents? I'm having troubles with **intent** extras... They don't always seem to work! So, is there any other way? What I've tried with ... intents: String id = intent.getStringExtra("id"); String name = intent.getStringExtra("name"); But it gets two kinds of string every time I start the **activity**. The first time is different from the other times. Can I pass that second **activity** any other things without using intents? ...  
android android-intent android-activity  
asked Jun 24 '13 by Amin Keshavarzian

The same query done through Google gives much more relevant results:

stackoverflow.com › questions › how-do-i-pass-data-be... ▼

## How do I pass data between Activities in Android application ...

Mar 15, 2017 — The easiest way to do this would be to **pass** the session id to the signout **activity** in the **Intent** you're using to start the **activity**: `Intent intent = new Intent(getApplicationContext(), SignoutActivity.class); intent.putExtra("EXTRA_SESSION_ID", sessionId); startActivity(intent);`

52 answers

[Android: how to pass an activity as extra within an ...](#) 3 answers Dec 12, 2019

[How to pass an object from one activity to another ...](#) 32 answers Mar 12, 2016

[How to pass a value from one Activity to another in ...](#) 7 answers Sep 21, 2012

[How to send an object from one Android Activity to ...](#) 36 answers Nov 5, 2011

[More results from stackoverflow.com](#)

You visited this page on 12/9/20.

The retrieved information directly matches what I am looking for not only based on title but also based on phrase content within the posts. A complication brought on by a majority of the users not being able to correctly utilize the StackOverflow search is the idea of duplicates. If people are not able to find the information they need on the website, then they post a question. However, the google results show that the same questions often get posted repeatedly. Though in the above example, comments in several of the posts like them together and identify the duplication, it seems as if the community monitoring has to make up for the issues from the search.

### Solution Proposed:

Based on the issues identified, I believe that the relevance of the search engine output can be significantly improved if the ranking relevance decreases importance on the title and incorporates text from other metadata. In the improved ranking, the main difference is putting an emphasis on the community activity on the post rather than the post's title. The query will be processed to retrieve results that match based on phrases in the "Body", "Title", "Answer", and "Comment".

A first iteration of the proposal suggested that additional attributes be taken into account such as "CreationDate" and the "LastActivityDate". This way the ranking prioritizes the retrieval of posts created recently with the most up to date answers. However, the corresponding retrieval system prioritized recent and sparsely active posts. This occurs for the use case of StackOverflow because a recent post is often not answered and marked duplicate. Due to the technical and comprehensive question and answering system on this website, it is better to base relevance on community involvement in a post rather than it's recency.

Based on a detailed look at the dataset, this proposal will implement a search where relevance is based on the similarity between query text and title, post content, comments on posts, answers on posts, and comments on answers.

In the image below, parts of a StackOverflow post are circled with the content that was indexed for ranking:

## Python tokenizing strings

Asked 6 years, 10 months ago   Active 6 years, 10 months ago   Viewed 4k times

**Build the products that get people jobs.**  
Work remotely from anywhere in the United States.

**indeed**

Report this ad

I'm new to python and would like to know how I can tokenize strings based on a specified delimiter. For example, if I have the string "brother's" and I would like to turn it to ["brother", "s"] or a string "red/blue" to ["red", "blue"], what would be the most appropriate way to do this? Thanks.

python   regex   tokenize

share   improve this question   follow

asked Feb 5 '14 at 6:15

Mozbi  
1,179 ● 1 ● 13 ● 21

I would start with `pydoc str` and work from there. – larsks Feb 5 '14 at 6:17

add a comment

3 Answers

Active   Oldest   Votes

You would use the `split` method:

```
>>> 'red/blue'.split('/')
['red', 'blue']
>>> "brother's".split("'")
['brother', 's']
```

share   improve this answer   follow

answered Feb 5 '14 at 6:18

jterrace  
55.6k ● 20 ● 136 ● 181

Thanks. How about if I had something like "brothers" with the single quote after "brothers" and I want it to be ["brother", 's']? – Mozbi Feb 5 '14 at 6:31

add a comment

## Dataset Collected:

The data ecosystem of StackOverflow is quite large and cumbersome with over 40GB of data. The same archives (maintained by StackExchange, which oversees StackOverflow as well as a network of other webpages) however have smaller collections of data on other topics that are structured in the same way. For the scaling purposes of this project, the data utilized for the search engine came from codeReview StackExchange. The size of the data queried and indexed is approximately 1GB altogether.

The second complication in collecting the data is that it is separated into several different xml files, each of which represent a component of a post on StackExchange. I have outlined the overview of the relevant data with their corresponding files below:

- Comments.xml
  - Id

- Text, e.g.: "@Stu Thompson: Seems possible to me - why not try it?"
- Posts.xml
  - Id
  - PostTypeId
    - 1: Question
    - 2: Answer
  - ParentID (only present if PostTypeId is 2)
  - AcceptedAnswerId (only present if PostTypeId is 1)
  - Title=
  - AnswerCount
  - CommentCount
  - Body

Other aspects of the data can be seen on the archive documentation; however, since they span over two pages in structure, for the purposes of brevity I have only included the portions of the dataset the corpus for the search index was built on.

After obtaining the data from the XML files, a basic tokenization was implemented on the “Body” of the posts. This involved removing stop-words that did not contribute heavily to the information retrieval. Punctuation nor html tags were removed due to the nature of the posts. The Code Review StackExchange supports several different languages as well as several different writing conventions to express questions. Both punctuation and non-english words play a significant role in such a technical forum.

Comments on posts and answers were also processed by removing the @ annotation on the text field in the comment. The @ with a username links to someone else in the conversation is a major writing tag on the website. However, there is already a constraint currently in the StackOverflow website that can search user content solely based on the user @ handle. Since this engine is based on the content rather than user constraints, notations with @ handles in the comments were removed.

Lastly, the data was filtered to exclude all posts that do not have a title. Not all posts contain a title. StackOverflow currently has a system that allows for users to look at duplicate posts with similar titles at time of post creation. Therefore, it is less likely that a user duplicates a post if they make sure to include a title since they will have alternate suggestions to ensure that their questions have not already been answered. The alternate suggestions also allow the user to look at other similar questions and phrase their post appropriately for the topic. As a result, for the purposes of increasing the quality of the data set used for the indexing, posts without titles were filtered out.

## Implementation:

This project is a python application. The requirements.txt specifies all the libraries utilized in the application:

- Flask==1.1.2
- elasticsearch
- bs4

- bleach

All services that make up the application in order to start the Elasticsearch cluster are specified in the docker-compose.yml. The dataset is stored as xml files and then parsed into dictionaries. Elasticsearch API python client is utilized to implement the search query. Lastly a Flask server is utilized for the backend.

## **Backend:**

### **Data restructuring and Index Creation:**

After data collection, restructuring is done to map all the contents of the post to its corresponding post id. In the first round of indexing, each comment is processed and mapped in a dictionary to its "PostId". In the second round of indexing, the Posts.xml data is split into posts that are questions and posts that are answers (the original dataset combined both). The boolean distinction between the two types of posts is indicated by the "PostTypeId" field. If the post is a question, it is mapped to its "Id". If the post is an answer, it is mapped to its "ParentId". This way both questions and answers can be mapped to post content. The third indexing is done to "Title" If the post is a question and, if a title tag is included, then the "Title" is mapped to the "Id".

The result of the above restructuring allowed for four different indexes: comments, title, post and answers. Utilizing the id attribute of each post, a combined dictionary is created of the processed comments, title, post, and answers to allow for a nested mapping of the Id of a post to its content as defined for this proposal. A simple combined index of {'id': 'processed text of corresponding post'} really helped to simplify the data for the search query. It matched all post content and nested the titles with questions and answers and comments. This restructuring was a key component in including not just the title but also the valuable community engagement within a post relating to the title.

### **Search Query:**

Utilizing the combined dictionary, the Elasticsearch API is utilized to store the index. The indexing here is simple due to the significant restructure and index creation of the data. The resulting index matched the id as the unique "Id" of each post and the "text" as the combined text of each post. For the query, the result is determined by directly matching the query (not necessarily as a phrase) to the processed index of each post. From the API library, the match query is able to search for analyzed tokens rather than exact terms so parts of phrases can also be matched.

The result (result['hits']['hits']) of the match query returns the number of documents that match the given ask. The scoring for the document is based on the field match between the query and the text indexed. This is derived from Lucene's Practical Scoring function, which Elasticsearch utilizes under the hood. This scoring is comparable to the TF-IDF scoring in vector space. A future work in this component of the backend would be to redefine the scoring function to be less computationally expensive if the dataset (in the case of the original StackOverflow dataset) is too large or expensive.

The fields retrieved from the Elasticsearch query were result["hits"]["hits"][["\_source"]]["text"] (with the list of documents that were matched) and the result["hits"]["hits"][["\_score"]] (with the scoring of each document). A third field is retrieved with result["hits"]["hits"][["highlight"]]["text"]. This is returned as a

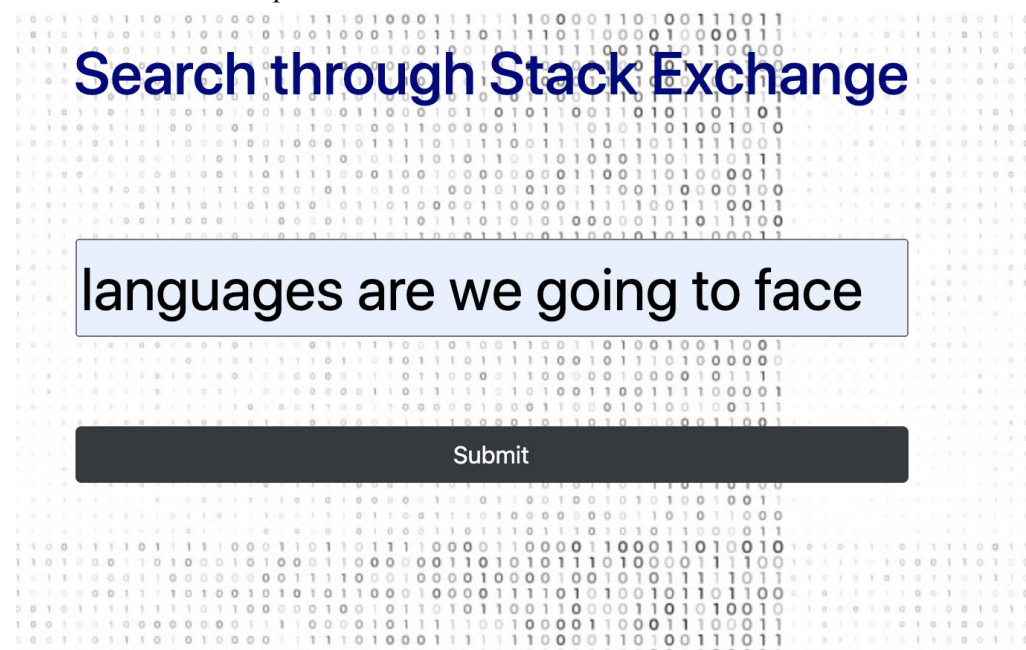


customization of the query to include html tags around the phrases that the elastic search query matched. This is done after receiving feedback from the professor on the UI during the in class demo.

From the results of the search, a heap structure is utilized to rearrange the results in descending order based on the ElasticSearch score of each ranking.

### User Interface:

The primary UI page is very simple and mimics the base engine format. Feedback from the demo to include a large search bar has been incorporated. This encourages users to use larger queries and not limit themselves to smaller phrases.



Once all documents are returned, additional processing is done to check if data returned contains html. At point of rendering, the UI originally had been cluttered with post content. It is very difficult to understand why something is matched if the UI displayed the whole post as text. The html tags help identify which part of the result contains the matches. However, rendering with dependency on the highlight html tags meant that other parts of the post can break the UI.

As a result, further checking is done to see if ElasticSearch matched with english phrases within the post or if it matched within the html metadata inside a post. The complicated nature of the wide variety of posts in StackExchange makes it difficult to dismiss data if matched within the html. This is especially important in terms of “tags” within posts. Since this solution allows for matching all post content, tags within posts could also be a reason for relevance. Checking if the highlighted match result is a HTML or a phrase allows for better rendering on the UI. The following images are results from searches performed during the in-class demo. The ranking function has not changed but the UI has been improved to allow for easier user understanding of why a result is matched.

### 1. Knuth's Web



[Go to post!](#)

## 2. Title: These mysterious \$A framework questions

[Go to post!](#)

### 3. Title: SiteSmashers - Reviewing Website designs?

[Go to post!](#)

4. Title: Can I remake other people's websites using HTML/CSS (only for education purposes) and post in Code Review?

[Go to post!](#)

## 2. elasticsearch java api

## Top 10 StackExchange Results

## 1. Title: Do we have some rebel tags that shouldn't exist on Code Review?

[Go to post!](#)

Post Context: `elasticsearch java api` not available in preview. However, post is relevant through links and tags so proceed to site.

## 2. Title: What's the [api] tag for?

[Go to post!](#)

Post Context: a REST *API*.

### 3. Title: Are implemented API reviews without code on-topic?

[Go to post!](#)

Post Context: To be clear, I'm not looking for a comparative "is this *API* better than this other *API*" or help

#### 4. Title: Should we burn [Google]?

[Go to post!](#)

Post Context: " class="post-tag" title="show questions tagged 'google-*api*'" rel="tag">google-*api* tag.

\*Note that elasticsearch is present as a tag in this dataset. The first result leads to this post:

<https://codereview.meta.stackexchange.com/questions/3673/do-we-have-some-rebel-tags-that-shouldnt-exist-on-code-review>

# Do we have some rebel tags that shouldn't exist on Code Review?

Ask Question

Asked 6 years, 1 month ago   Active 6 years, 1 month ago   Viewed 103 times

I just found a question with `finance` and there are 45 questions tagged with it. what does that have to do with code review?

3

This led me to the Tags page and I found a few that look out of place

1

- `arguments` -> could be `parameters` ? or not (it doesn't exist)
- `adventure-game`
- `asynctask` -> **note** I just retagged the only 3 questions with `asynchronous`
- `autoloader`
- `contextmenu` -> I added a tag wiki for that.
- `connection-pool`
- `covariance` probably a good tag that just isn't used often, no wiki.
- `data-importer`
- `device-driver` re-tag as `drivers` and add a wiki?
- `elasticsearch` I think we have a tag for Search algorithms, don't we? maybe this should be synonymized

From here, clicking on the tag `elasticsearch` leads to all the (controversially) tagged `elasticsearch` posts in the site.

## 3. Off-topic

### 1. Title: Should we have just one Meta tag for topicality?

[Go to post!](#)

Post Context: a href="/questions/tagged/*off-topic*" class="post-tag" title="show questions tagged '*off-topic*&#39;39

### 2. Title: Should answers to clearly off-topic questions be encouraged or discouraged?

[Go to post!](#)

Post Context: This question is clearly *off-topic*

### 3. Title: Move to SO or Programers

[Go to post!](#)

Post Context: it with a quick search; But: But there must already be a question here on meta about moving *off-topic*

### 4. Title: Should the FAQ on off-topic questions link to other SE sites?

[Go to post!](#)

Post Context: there is a *off-topic* not available in preview. However, post is relevant through links and tags so proceed to site.

### 5. Title: Gold tag-badge users closing questions as pre-written off-topic duplicates

[Go to post!](#)

Post Context: Gold tag-badge users closing questions as pre-written *off-topic* duplicates

It is important to note that phrases within the metadata including post body, tags, comments, and answers are being highlighted in the match phrase. This demonstrates the original goal of making community engagement a key identifier in ranking rather than just the title.

The last feature that is key in the UI is the “Go to post!” link. This link allows for the user to easily navigate to the corresponding post. The link is created through url normalization based on the same id utilized as the key in the combined index. For example the post with an id of 10 is created by concatenating "https://codereview.meta.stackexchange.com/questions/" + 10. This is very important because it is difficult to find the exact post on the current StackOverflow Engine.

## **Conclusion:**

The main purpose of this proposal is to improve on the current StackOverflow search engine by decreasing query matching with emphasis on the title and, instead, making community engagement the primary determinant of relevance. Community engagement here has been determined by a combined index of post body, title, answers, comments on the answers, and comments on the posts. The resulting ranking system now promotes retrieving less sparse, higher quality matches that can really benefit the user. Future work still needs to be done by exploring a larger data set and improving the relevance to better process html components and code within a post.