



Articles » General Programming » Internet / Network » Network

Making a Search Engine

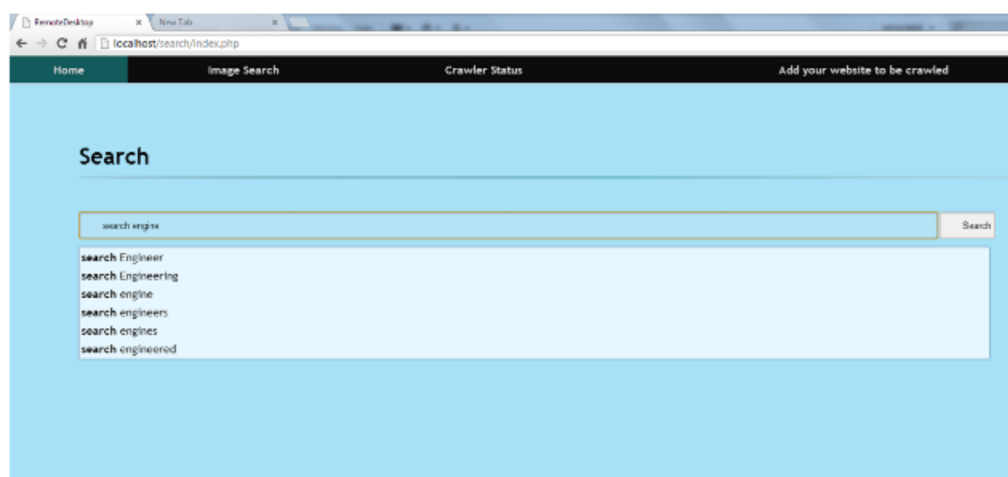
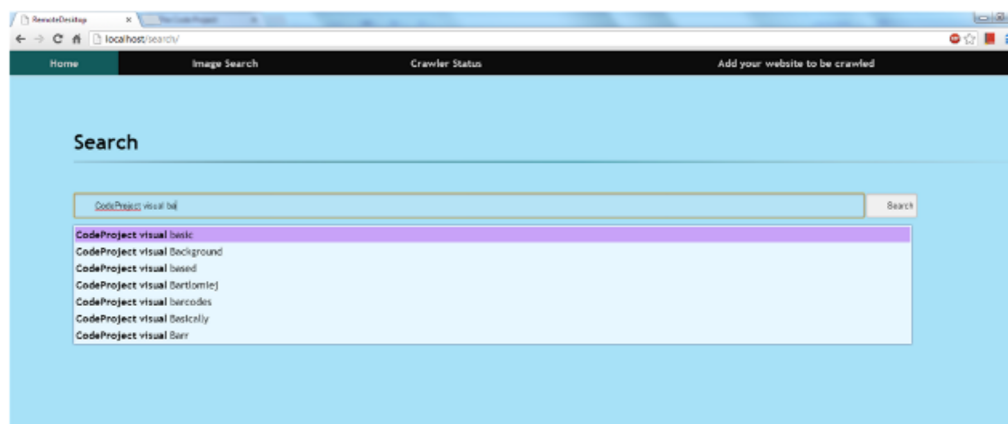
kburman6, 4 May 2013

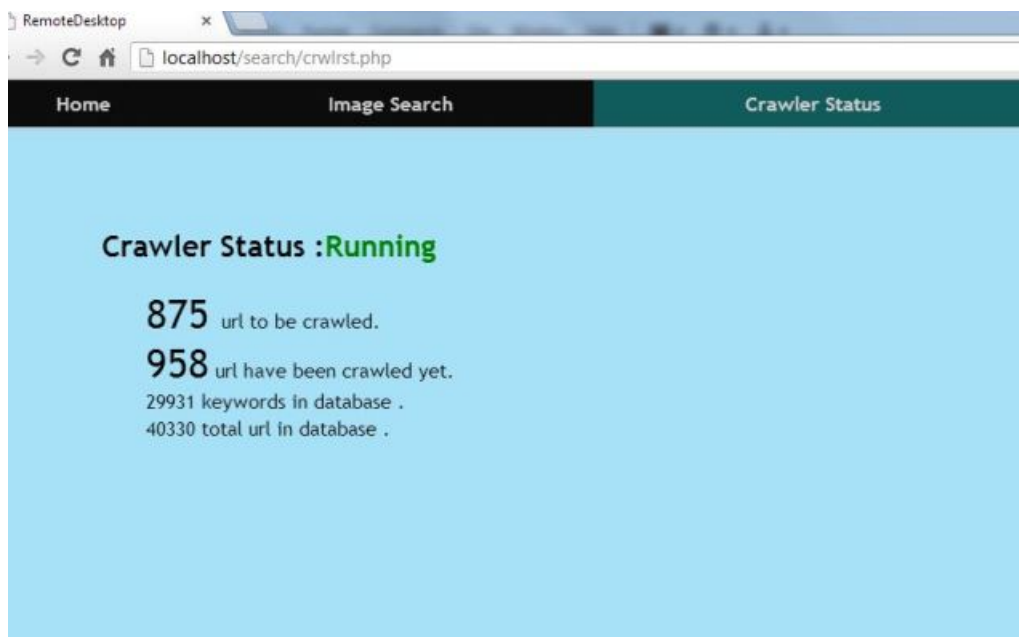
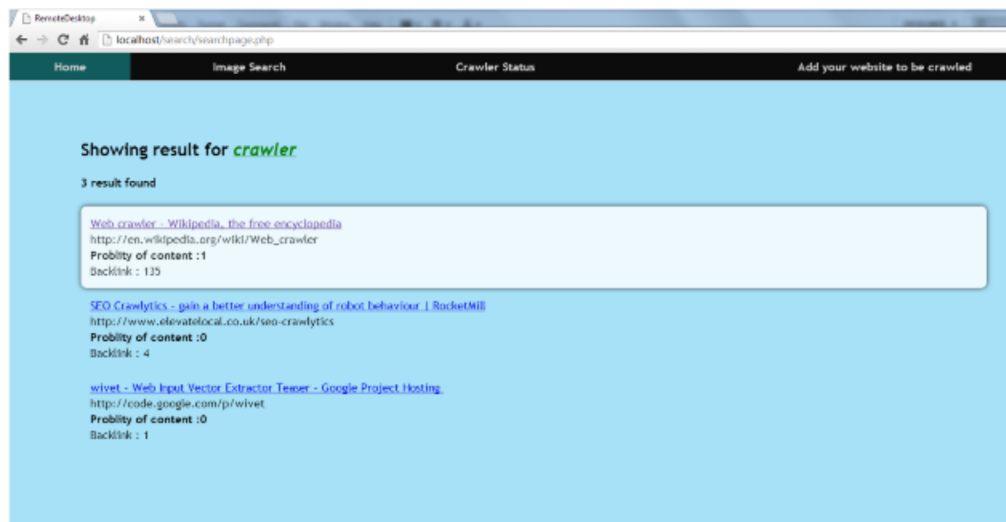
This article discusses the making of a search engine.

 **Download Executable files - 178.3 KB**

 **Download source -21.5 KB**

 **Download PHP file - 6.83 KB**





Introduction

This project is still not complete. You must use it to crawl thousands of URLs because you may find that it crawls the same URL for the last 100 times. (This is because of some unidentified problem in conversion of relative to absolute URL.)

Like most search engines, this one also has a crawler whose basic aim is to retrieve the source code of a given URL and then break the content into words with which we can create an array of tag words which will represent the content of the site. It is not a fool proof method, but can work for sites with lots of words in it, like a blog or article or a discussion forum, etc.

For example:

TAG CLOUD of <http://www.google.co.in>

account advanced advertising bengali blogger books business calendar co com drive
feeling fields gmail **google** gujarati hidden hindi history
images india kannada language lucky malayalam maps marathi news offered options
orkut photos play privacy programs punjabi reader **search**
settings sign solutions tamil telugu terms tools translate web youtube

TAG CLOUD of <http://www.wikipedia.org>

activities archive areas **articles** bahasa city
community current **dorgon** email emperor encyclopedia
english featured foundation francais **free**
free-content game history hosts languages list main manchu
ming navigation nederlands news norsk **page** pictured
political pope portugues projects qing range recent rule son
successful system third turkce view **wikipedia**
wikipedia world years

TAG CLOUD of <http://www.w3schools.com>

ajax asp **browser** building certification certified character code
color com complete **CSS** css3 dom editor
examples experiment **html** html5
javascript jquery **learn** net offer pages
php picker popular **quiz** razor read references result
services sets sql statistics svg tests topics **tutorial**
validate w3schools **web** website xhtml
xml xpath xsl-fo xslt

It is clear how tag cloud can highlight the key words which could describe a given URL.

Now these keywords are stored in a database and then used to find the relevant URL for given keywords.

Background

It all started when my friend showed me his search engine with 4 URLs in an XML file. It seemed like an auto complete rather than a search engine, but later that night it was 2:00 am, and I couldn't sleep at all because of that auto complete feature with which I was too impressed. I wanted my own... there was a thunder storm of ideas in my mind. After 3 sleepless nights, on the 3rd day at 6:00 a.m., I was ready with my search engine working with 100 URLs in database... that was the time when I finally slept comfortably and full of satisfaction. It took me 3 days because every day I started from the beginning because I was not satisfied with the performance of the crawler or there was some problem.

Using the Code

I have basically divided every task into small parts so that the work could become easy. So you would find lots of classes in the project.

Some important classes are given below:

- **Panda**-> It retrieves the source code of a given URL using a `get_sourcecode(url)` from Module '`func`' and then passes it to **juicer**(class) for extracting the useful information and after the work is completed, reports to the **panda_manger** (boss of all the **pandas**) and again assigns a new job to **panda** if any.
- **panda_manger**-> It manages all the **pandas** and from this class we assign any web URL for crawling and it will automatically assign this work to any free **panda** and if no **panda** is free, then it creates a new **panda** and assigns the work to it. When a **panda** finishes, it works and reports back to **panda_manger** then **panda_manger** checks whether there is any URL left to be crawled, if any then it gives the command to crawl that URL to the **panda**.
- **juicer**-> This could be said to be the main class for the whole crawler which extracts keywords from the source code of any website and then saves it into the database.
- **juicer**->**extract_juice**-> This class gets the source code and then converts it into `HtmlAgilityPack.HtmlDocument` using the method `LoadHtml(source)`. It is used because we need an HTML parser to scrap the text out of HTML (building a custom one will require a lot of time) and since it supports xpath for retrieving any element(s) from source, it simplifies our work a lot.
Now we pass the `HtmlDocument.DocumentNode(.SelectNodes)` to the various methods of this class for extraction of some type of information.

```
Dim doc As New HtmlAgilityPack.HtmlDocument()  
doc.LoadHtml(source)  
process_texttag(doc.DocumentNode.SelectNodes("//meta"))  
process_anchor(doc.DocumentNode.SelectNodes("//a"))  
process_image(doc.DocumentNode.SelectNodes("//img"))
```

- **juicer**->**process_metatag**-> Currently, it just processes the meta tag containing keywords and then assigns every word a 45% to total word count in the document and 30% of total word count to the words which are mentioned in the keywords but can't be found in the document.
- **juicer**->**process_anchor**-> This is the most important part if you want the search engine to automatically move to the next link without manually entering every link to the crawler and the hardest part was to convert the relative URL to absolute URL. But Microsoft saved me with the `URI` class which can easily be used to convert the `URI`.

Database Structure

At the starting point, we have a 4 tables stored in a database named as "**Crawler**" by default.

Table Name : "Keyword_index"

This table is used to provide the suggested result in the search box. It contains all the words which the crawler has encountered till now and no word is repeated.

```
mysql> select * from keyword_index where word like 'a%';
```

word	freq
A	20
About	1
ac	5
According	2
accumsan	3
Achievements	1
across	1
address	4
adipis	2
adipiscing	4
adiscing	1
ado	2
Airtel	3
alais	1
aliquam	2
aliquet	1
All	6
amet	9
an	3
and	10
annihilation	2
answer	2
ante	1
antiparticle	3
api	2
arcu	3
are	2
as	2
asp	13
at	7
atl	13
Atom	1
augue	2
Available	2
Axiata	3
AXIS	3

```
36 rows in set (0.00 sec)
```

Search

A
About
ac
According
accumsan
Achievements
across

For multiple keywords, we first break the keyword from " ", then find the urlhash for the word and then find what the other words are that urlhashes contain and show the words. (Not yet implemented because I was getting an unidentified error and if I used SQL `join`, then it took more than 1 minute to search the database for 1 keyword and time increased exponentially.)

Table Name : "Keyword_list"

This table contains all the words in a tag cloud of the given URL. And to identify this word belongs to which URL, we store the urlhash of that website with the word.

```
mysql> desc keyword_list;
```

Field	Type	Null	Key	Default	Extra
word	varchar(100)	NO		NULL	
urlhash	char(32)	NO		NULL	
freq	int(3)	NO		NULL	
per_in_site	int(3)	NO		NULL	

```
4 rows in set (0.01 sec)
```

```
mysql> select * from keyword_list where word like 'a%' limit 10;
```

word	urlhash	freq	per_in_site
A	ead7b7cab3aa7349913987b7b176bb15	1	0
All	ead7b7cab3aa7349913987b7b176bb15	1	0
at	ead7b7cab3aa7349913987b7b176bb15	1	0
AT	ead7b7cab3aa7349913987b7b176bb15	1	0
address	34a84aa0676b69fc8826328bdaa9967d	1	0
Address	34a84aa0676b69fc8826328bdaa9967d	1	0
all	34a84aa0676b69fc8826328bdaa9967d	1	0
abelardo	4d04edf6d9edc6d7f1b36b2b65028523	1	0
About	4d04edf6d9edc6d7f1b36b2b65028523	1	0
adara	4d04edf6d9edc6d7f1b36b2b65028523	1	0

```
10 rows in set (0.00 sec)
```

Table Name : "url_webpage"

This table is used to store all the links which it finds in the pages that it crawled till now. In this table, we also use MD5 hash of URL to refer to that URL instead of the original URL because no one knows how long the URL which we may find can be. So we change it to MD5 because it does not matter how long the URL is, the MD5 will always be 32 chars long.

```
mysql> desc url_webpage;
```

Field	Type	Null	Key	Default	Extra
urlhash	char(32)	NO	PRI	NULL	
url	varchar(250)	NO	UNI	NULL	
state	int(2)	NO		NULL	
crawl_date	datetime	YES		NULL	
backlink	int(3)	NO		NULL	
priority	int(1)	NO		NULL	
title	varchar(500)	NO		NULL	

```
7 rows in set (0.01 sec)
```

urlhash	url	state	crawl_date	backlink	priority	title
02a504b16154bb1b136aeb0b081d1	http://mitunscasseroles.tumblr.com	2	(null)	0	0	Zagazo
c66a50425514b11b6a754b7cb020dc5	http://www.youtube.com/?hl=en&gl=IN&as=H	2	(null)	0	0	YouTube
fe341ad7303b85f02cc12f3d032ec0cb	http://www.youtube.com/?gl=IN&as=w1	2	(null)	0	0	YouTube
0026c50ef8ba23b999225a110431019	http://fucksocial.tumblr.com	2	(null)	0	0	your stupid
02132d1671b7ba445dbbd04a8c97675d	http://somuchfortakingchances.tumblr.com/p	2	(null)	4	0	You have to love something before you
04c072d76623a043b670014eb01a29b	http://somuchfortakingchances.tumblr.com	2	(null)	7	0	You have to love something before you
056eac9b50cd1d39693d0d5dc26092	http://somuchfortakingchances.tumblr.com/p	2	(null)	4	0	You have to love something before you
009e190917268377465a2f8a8a40a44c	http://smile-youre-amazing.tumblr.com	2	(null)	0	0	you are loved
046d17a3396ad4325ed2b775d9fa48	http://forevertillinfinity.tumblr.com	2	(null)	1	0	YOU ARE BEAUTIFUL J
006b0c0a098c612c38e94862d7607c	http://megstalejon.tumblr.com	2	(null)	1	0	Ya did clean
021c30e2f0012890e0e838e10e8e	http://neogrest.tumblr.com/post/4487167022	2	(null)	2	0	WOW/GREAT - Joe Winter (via Primer
013163eb3f6e697635c3295039ce7	http://televholvstheuniverse.tumblr.com	2	(null)	0	0	Worms - oh my god WORMS
01262b0bce1e0b7cfd3473850a967	http://selecionar.tumblr.com/post/453072911	2	(null)	3	0	Will you stay?
00ce42502fc5694b33d21e09bc801d	http://wildthemes.tumblr.com	2	(null)	0	0	wildthemes
005c2b136c001f7c5eb022a5e1b61db	http://live-in-these-moments.tumblr.com	2	(null)	1	0	Wherever you are, be all there.
006b0e4b6950b32139d7b1b967d0f	http://wishes-assbutts.tumblr.com/post/4531	2	(null)	2	0	Where you invest your love, you invest
019264a023d77c0506c493475780f03e	http://wishes-assbutts.tumblr.com	2	(null)	19	0	Where you invest your love, you invest
01d2f65da3d04b65e9e0c000efc10d3	http://wishes-assbutts.tumblr.com/post/4531	2	(null)	3	0	Where you invest your love, you invest
00612e3596d0413ee50e41d5f94b44	http://mylightinthedarknesskanis.tumblr.com	2	(null)	2	0	Where there is love there is life
046675e377e38e02d37672e6a540b	http://this-fire.tumblr.com/post/4504322235	2	(null)	0	0	Where Do I Put This Fire? (Isurbe:
039709585c62f1aeb0da918bc70dc	http://this-fire.tumblr.com	2	(null)	1	0	Where Do I Put This Fire?
0001e9e204a8d6a0c292f99d64578b	http://feelsadthensmile.tumblr.com	2	(null)	3	0	When you feel sad... just smile :D
001d1335c2d9d2975b416156c572ae	http://coolhips.tumblr.com/post/4530659884	2	(null)	4	0	whatever
0371eb56b0a13ec772907e198b8f	http://coolhips.tumblr.com	2	(null)	3	0	whatever
0188111a2e36d5d23ad6c187d3d177	http://whatstlov3.tumblr.com	2	(null)	1	0	What
024884e0d0e8f3c20a0018a17633156a	http://cheezburger.com/7132712448	2	(null)	16	0	What Wise I Thinking? - Cheezburger
00b95a2eada08d4b38e8deff594397b	http://eoy-un-dango.tumblr.com	2	(null)	8	0	What is love?
010ef309e1876147d4b2274b5b82be	http://eoy-un-dango.tumblr.com/post/453083	2	(null)	0	0	What is love?

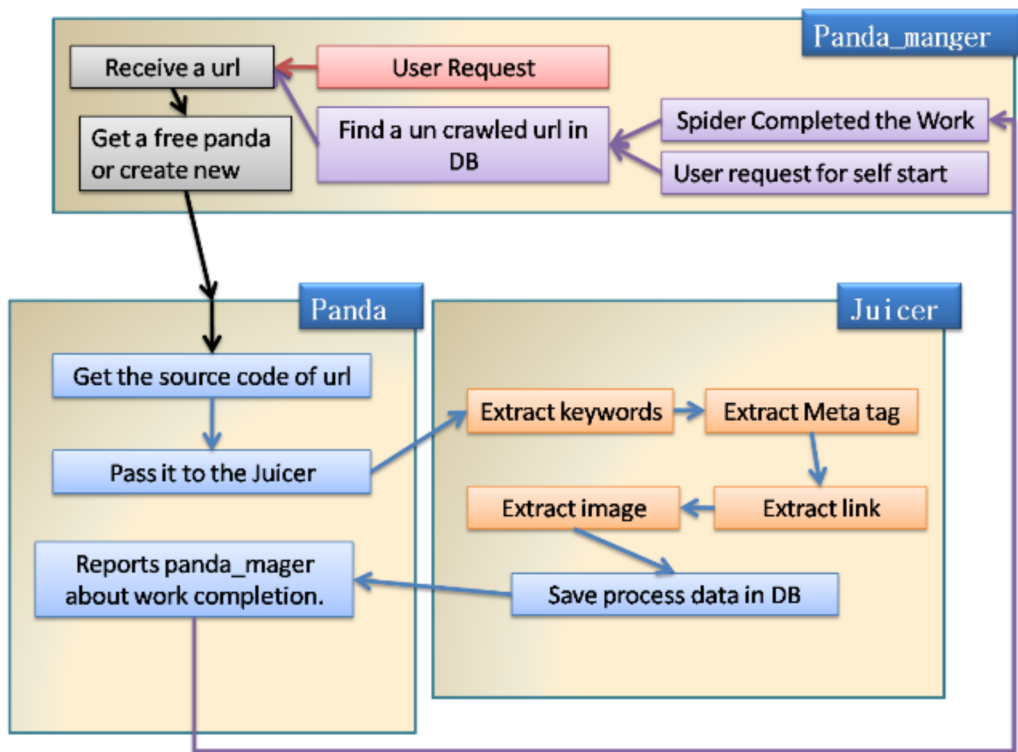
Table Name : "url_image"

This table is used to store all the links of the image which are found in the webpage it crawls. It will be used for image search, but I have not yet completed its processing and front end(PHP code).

So I will not discuss it right now. But I would like to tell that it will take a lot of processing power.

urlhash	url	state	crawl_date	backlink	priority	title
1c65170e11c70e47c1cc2b8ee8072	http://pixel.quantserve.com/pximg/19,1qE8egoZBM.gif	0	(null)	600	0	
0ce61001681d99704208a432470995	http://assets.tumblr.com/images/default_avatar_16.gif	0	(null)	227	0	
0c6b768f4d24451520e1d8f6e2240	http://static.tumblr.com/vl883dv/OUJhB58/epesox.gif	0	(null)	118	0	
2846406a849c4c60195691c4608d	http://assets.tumblr.com/images/sex.gif	0	(null)	129	0	
07d35b18e1810e614443879c0214	http://static.tumblr.com/uygh83v/OAGCwme/reblog.png	0	(null)	52	0	
6a0af5ae1375c0e32fec275d936393	http://i.ytimg.com/vt/ajmp/pxcel-vf3z5v9fn.gif	0	(null)	86	0	
4e1711ee5d1bbbf3bfe39cd1c691	http://out.cursor-4u.net/cursor.png	0	(null)	43	0	
7788859e1486c4d8c38585a3863e	http://b.cdn.americanexpress.com/p/617-w2k2-15742520&pxv=2_0&mpx=q=1	0	(null)	41	0	
602a0035e39ff79c9f9632026fc035	http://boring0a.akamaihd.net/profile_images/3007514026/0516e5cc3636eef82091468147	0	(null)	39	0	
0d3036a1776779c3bd001933e6ddcd	http://static.tumblr.com/cont585/O8lnu59/c4c7ryablogbutton.png	0	(null)	38	0	
315cc238c0e91471d2d945de19083	https://pixel.quantserve.com/pximg/19UteEDngoZBM.gif	0	(null)	35	0	
33c86736ac07f64848001d77da29a	http://assets.tumblr.com/images/frame_reblog_alpha.gif#74	0	(null)	34	0	
c50f05eb1494e59f0a60c76d4267537	http://assets.tumblr.com/images/default_avatar_40.gif	0	(null)	33	0	
84f78373479d3c02e515ba37851f8	http://i.ytimg.com/vimages/peasout.gif	0	(null)	29	0	
11b925ea4702b1452561c11b91c788	http://28.media.tumblr.com/tumblr_lpn6k5eZUG1qm57ime1_500.png	0	(null)	28	0	
8097a0eb6f68cd054f6c31050630	http://static.tumblr.com/v6ixh7w2/Gomkicv/rebloggy.png	0	(null)	28	0	
9eeb7592fc381a6c40476602085c	http://static.tumblr.com/v6ixh7w2/50mmibkzckleece.png	0	(null)	28	0	
ed5477141d969369c25a36690ab574	http://static.tumblr.com/lm6e6v6V/Yellowfir/heart_gill_3x7.png	0	(null)	28	0	
056a7432c0e9f7b4dd7c368983b1	http://uploads.netasm.com/users/themes/pel/media/loading.gif	0	(null)	26	0	
058ee77045e877cd9f5cc0e56e7a	http://i42.snyipic.com/23u9v.jpg	0	(null)	25	0	
98614742c34330d9804da5319808e	http://bchod4.discoveryeducation.com/images/acknowledgmentpaper.gif	0	(null)	21	0	
066c12bbe05038c330e780c00b42	https://s.chazbr.com/s/release_20130313/5img/chaz-logo-req.png	0	(null)	19	0	
24b1c178b2b119d808d817f7092f13	http://i40.twimg.com/profile_images/3007514026/0516e5cc3636eef82091468147404_nor	0	(null)	19	0	
474666307370174331c25264e6347c1	http://static.tumblr.com/jf1xwlp/25m8e6reblog.png	0	(null)	19	0	
40363037370174331c25264e6347c1	http://static.tumblr.com/jf1xwlp/25m8e6reblog.png	0	(null)	19	0	
9b773597474331c25264e634a8d	https://img0a.akamaihd.net/profile_images/2515226595/124ghoofdvruv8z2_normal.jpg	0	(null)	19	0	
3c3704801d2c9e47f78dd5b6b33d	http://i40.twimg.com/profile_images/2515226595/124ghoofdvruv8z2_normal.jpg	0	(null)	19	0	
2036414e5299267474bbf0a0c2172a	http://assets.tumblr.com/images/sex.gif	0	(null)	17	0	

How the Crawler Works



Points of Interest

It was very annoying to get rid of the relative URL. I tried several ways of resolving it, but every one ends up with some bug. Then I got a magic class, named as `Uri` which solved all my problems, but while writing this article I should crawl codeproject.com and show the result for CodeProject search, but then a problem hit my crawler after identifying that I have deleted my database of URL, otherwise I would have posted a screenshot of that. It was something like `/search.aspx` (some text)(same text repeated again)(and again, increasing with each crawl). It may be a problem with my code. I will try later to identify this problem and post the solution.

Currently Crawler is not following robots.txt so be careful while crawling.

It is you responsibility if crawling a site which is not to be crawled.

Sorry, but currently i am working on a new model of crawler.

License

This article, along with any associated source code and files, is licensed under [The Code Project Open License \(CPO\)](#)

About the Author



kburman6

Student

India

I just love coding. But due to my studies it became very tough for me to manage both.

Comments and Discussions

57 messages have been posted for this article Visit <https://www.codeproject.com/Articles/563869/Making-a-Search-Engine> to post and view comments on this article, or click [here](#) to get a print view with messages.

[Permalink](#) | [Advertise](#) | [Privacy](#) | [Cookies](#) | [Terms of Use](#) | [Mobile](#)
Web01 | 2.8.190205.1 | Last Updated 4 May 2013

Article Copyright 2013 by kburman6
Everything else Copyright © [CodeProject](#), 1999-2019