

# Social Networks Course Project KML (Knowledge Markup Language)

Amit Kumar Verma (mentor)  
Paras Kumar  
Nitin Gandhi

# Contents

- 1 Introduction to KML
- 2 Components
- 3 KML-Compressor-Decompressor
- 4 KML-for-Github
- 5 Github-Spider
- 6 Github-Spider-Proxy-Rotated
- 7 User-Agent-Spider
- 8 Wiki-Satck-KML-Downloader
- 9 Analysis

# Introduction to KML

KML stands for Knowledge Markup Language, a standard format for storing the data of all the Knowledge Building Portals. We are trying to propose a new standard format for all these types of portals such that the analysis is easy.



# Components

- KML-Compressor-Decompressor
- KML-for-Github
- Github-Spider
- User-Agent-Spider
- Wiki-Satck-KML-Downloader

Source Code: <https://github.com/csl-622/KML>

# KML-Compressor-Decompressor

These programs compresses and decompresses a KML file using diff algorithm inspired from git. It can compress the KML file by 30 percent or more depending upon the edits made on an article.

```
1 diff --git a/app.js b/app.js
2 new file mode 100644
3 index 0000000..144ec7f
4 --- /dev/null
5 +++ b/app.js
6 @@ -0,0 +1 @@
7 +module.exports = require('./lib/index');
8 diff -git a/index.js b/index.js
9 deleted file mode 100644
10 index 144ec7f..0000000
11 --- a/index.js
12 +++ /dev/null
13 @@ -1 +0,0 @@
14 -module.exports = require('./lib/index');
```

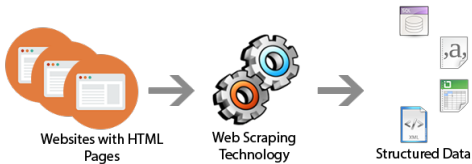
# KML-for-Github

We had to make KML for github and we had KML for wikipedia earlier. The attempt of this program is to retain the structure of KML for wikipedia while representing all the data in a git workflow in this newly generated KML file.



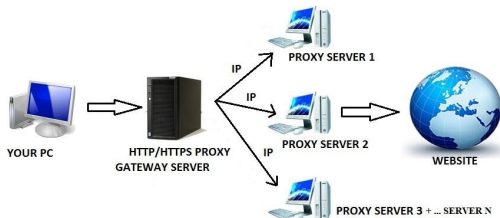
# Github-Spider

This multithreaded web-spider follows the rules of `/robots.txt` file and is useful for fetching small sized GitHub repositories. It outputs the results in JSON format.



# Github-Spider-Proxy-Rotated

Github-Spider-Proxy-Rotated is an advanced version of the previous spider; it contains pipelines, middlewares, and throttling parameters, which is in the middleware and settings file. It outputs the results in JSON format.





# User-Agent-Spider

This spider gets the latest user-agent headers from an online forum, output of this file is directly put into use in other two spiders namely Github-Spider and Github-Spider-Proxy-Rotated. The result of this spider is a text file containing a list of HTTP headers which can be useful for spoofing browser's HTTP request activity.

method	path	protocol
GET	/tutorials/other/top-20-mysql-best-practices/	HTTP/1.1
Host: net.tutsplus.com		
User-Agent: Mozilla/5.0 (Windows; U; Windows NT 6.1; en-US; rv:1.9.1		
Accept: text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=		
Accept-Language: en-us,en;q=0.5		
Accept-Encoding: gzip,deflate		
Accept-Charset: ISO-8859-1,utf-8;q=0.7,*;q=0.7		
Keep-Alive: 300		
Connection: keep-alive		
Cookie: PHPSESSID=r2t5uvjq435r4q7ib3vtdjq120		
Pragma: no-cache		
Cache-Control: no-cache		

**HTTP headers as Name: Value**

# Wiki-Satck-KML-Downloader

Dowloads xml from Wikipedia and Stackexchange for any given topic from user then coverts it to KML. The topics for wikipedia can be name of any page on wikipedia like UNO while for stackexchange the topics can any of the sub-site name of stackexchange like beer.



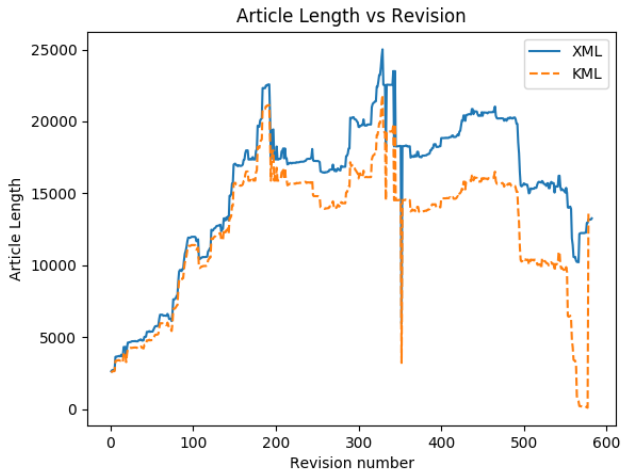
**WIKIPEDIA**  
The Free Encyclopedia

# Analysis

Graphs are generated to compare compressed KML and XML on various performance metrics like change in ArticleLength, Number of revisions, Number of Images, Number of Nouns, Number of sections etc vs. Revision number.



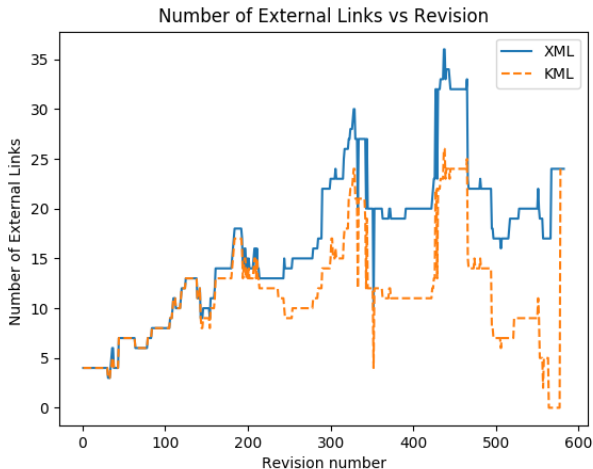
# Article Length vs Revision



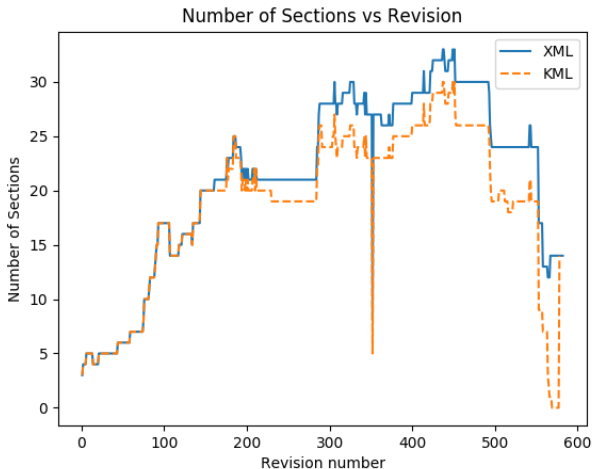
# Number of Proper Nouns vs Revision



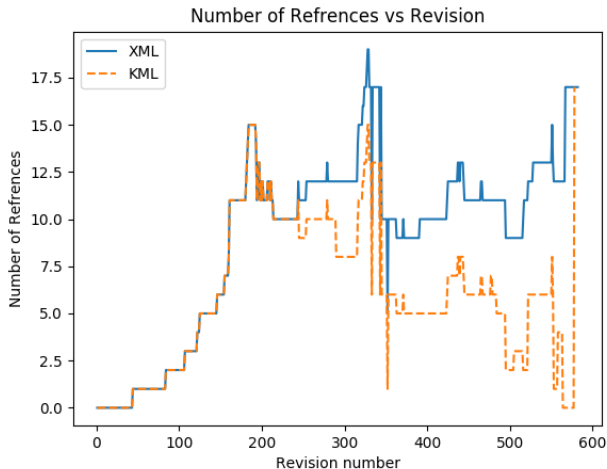
# Number of External Links vs Revision



# Number of Sections vs Revision

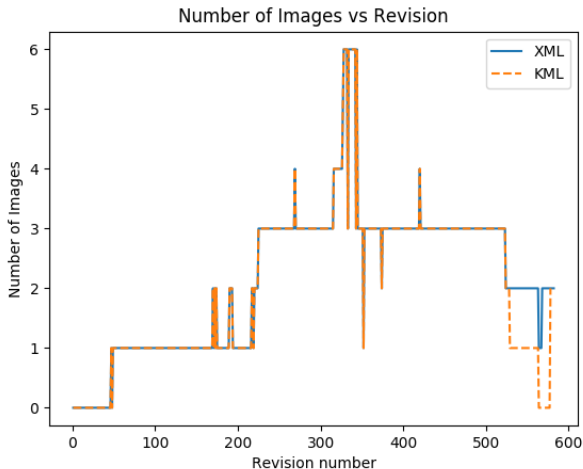


# Number of References vs Revision

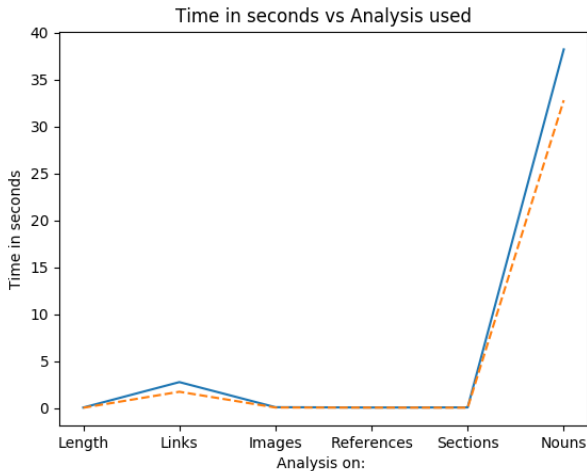




# Number of Images vs Revision



# Time in seconds vs Analysis used



## Link to our Repository

**URL:** `https://github.com/csl-622/KML`

