

2019

Towards linked data for Wikidata revisions and Twitter trending hashtags

Paula Dooley

Bojan Bozic

Follow this and additional works at: <https://arrow.tudublin.ie/scschcomcon>



Part of the [Computer Sciences Commons](#)

This Conference Paper is brought to you for free and open access by the School of Computer Sciences at ARROW@TU Dublin. It has been accepted for inclusion in Conference papers by an authorized administrator of ARROW@TU Dublin. For more information, please contact arrow.admin@tudublin.ie, aisling.coyne@tudublin.ie, gerard.connolly@tudublin.ie.



This work is licensed under a [Creative Commons Attribution-NonCommercial-Share Alike 4.0 License](#)

Towards Linked Data for Wikidata Revisions and Twitter Trending Hashtags

Paula Dooley
Technological University Dublin
Dublin, Ireland
paula.dooley@mydit.ie

Bojan Božić
Technological University Dublin
Dublin, Ireland
bojan.bozic@dit.ie

ABSTRACT

This paper uses Twitter as a microblogging platform to link hashtags, which relate the message to a topic that is shared among users, to Wikidata, a central knowledge base of information relying on its members and machine bots to keeping its content up to date. The data is stored in a highly structured format, with the added SPARQL Protocol And RDF Query Language (SPARQL) endpoint to allow users to query its knowledge base.

Our research, designs and implements a process to stream live Twitter tweets and to parse existing Wikidata revision XML files provided by Wikidata. Furthermore, we identify if a correlation exists between the top Twitter hashtags and Wikidata revisions over a seventy-seven-day period. We have used statistical evaluation tools, such as 'Jaccard Ratio' and 'Kolmogorov-Smirnov' to investigate a significant statistical correlation between Twitter hashtags and Wikidata revisions over the studied period.

CCS CONCEPTS

• **Information systems** → **Data extraction and integration; Wikis; Presentation of retrieval results; Social networks; Information retrieval query processing; Similarity measures.**

KEYWORDS

Wikidata, Twitter, Hashtags, SPARQL, Trending, Microblogging, Kolmogorov-Smirnov, Jaccard Ratio

ACM Reference Format:

Paula Dooley and Bojan Božić. 2019. Towards Linked Data for Wikidata Revisions and Twitter Trending Hashtags. In *The 21st International Conference on Information Integration and Web-based Applications Services (iiWAS2019), December 2–4, 2019, Munich, Germany*. ACM, New York, NY, USA, 10 pages. <https://doi.org/10.1145/3366030.3366048>

1 INTRODUCTION

Information on the World Wide Web is available through home computers and mobile phones, and with continuous advancements in technology, people have become increasingly more electronically connected. Along with this information, there has come many powerful innovation services facilitating both how people access information and how they connect with one another. Social networking sites, such as Twitter and Facebook have evolved alongside

wiki-sites containing huge amounts of information, such as Wikidata and Wikipedia. However, the broad variation of platforms makes it hard to determine whether, and how, current trends and topics are cross-related and whether what information a user consumes depends on the platform. Therefore, the aim of our research is not only to implement a system for streaming tweets and parsing Wikidata revisions, but also to investigate correlations of trends.

There are two main parts in our paper. The first part extracts the data from both Twitter and Wikidata. Twitter data are tweets posted by individuals consisting of hashtags, URLs, plain text and user names. The focus of this study will look at Twitter hashtags for comparison. Wikidata like Wikipedia is an encyclopedia of information [9] [16] which has evolved over time through authors continually revising the data to keep the information current. A revision is considered any one of insert, delete or substitution of data to an article [13]. This data is cleaned and prepared for comparison with Wikidata revision article titles. The top Wikidata revision articles and Twitter hashtags are identified over a seventy-seven-day period.

The second part of the paper compares the Wikidata revisions and Twitter hashtags to identify if a correlation exists between the hashtags posted and Wikidata revisions made. Statistical formulae, Kolmogorov-Smirnov & Jaccard's Ratio, will compare the text-ranked results from each group to determine if a statistically significant correlation exists. Visualisation analytics will be used to provide insight into the results of the Twitter trends and Wikidata revisions over the studied period.

The main research objective is to determine if trending topics in the English language Wikidata, identified by the title of the most frequently edited pages, show a statistically significant correlation to the real-time streaming data top-trending hashtags on Twitter, over the studied period, using the statistical analysis tools 'Jaccard Ratio' and 'Kolmogorov-Smirnov'. The research question and research hypothesis aim to support the objective defined as:

- Research Question: Is there a correlation between Wikidata revisions and trending topics hashtags on Twitter?
- Null hypothesis (H0): a correlation does not exist between Wikidata revisions and trending hashtags on Twitter.
- Alternative hypothesis (H1): a correlation exists between Wikidata revisions and trending hashtags on Twitter.

This research incorporates both primary and secondary methods. Initially, secondary research was conducted on existing literature which examined studies focused on Wikidata and Twitter data processing and analysis. It provided insight on both the current techniques for processing and analysing data and on the statistical analysis methods for text comparisons. Primary research was

ACM acknowledges that this contribution was authored or co-authored by an employee, contractor or affiliate of a national government. As such, the Government retains a nonexclusive, royalty-free right to publish or reproduce this article, or to allow others to do so, for Government purposes only.

iiWAS2019, December 2–4, 2019, Munich, Germany

© 2019 Association for Computing Machinery.

ACM ISBN 978-1-4503-7179-7/19/12...\$15.00

<https://doi.org/10.1145/3366030.3366048>

conducted through streaming live twitter data over the studied period, where the hashtag lists within each tweet were extracted for analysis. Secondary research also incorporated extracting revisions from Wikidata downloads that were used for further analysis. An experimental research method has been used on both sets of data to quantify whether a statistically significant correlation exists.

This project has four main objectives that will test the hypothesis:

- To retrieve streamed Twitter data, extracting its hashtag items per tweet. The data will be cleaned. Up to four n-grams will be applied and the data will then be ranked based on the volume of tweets over the study period.
- To extract Wikidata page details and revision data from Mediawiki data dumps and, using the SPARQL API endpoint, to retrieve the individual revision page titles. The data will then be cleaned by removing all spaces before counting and ranking the number of page titles based on the number of revisions occurring per page title over the studied period.
- To identify if a statistically significant correlation exists between both the top revised Wikidata pages and the top trending hashtags on Twitter. The statistical techniques to be used in identifying the presence of correlation are Jaccard's Ratio and Kolmogorov-Smirnov.
- To provide additional insights in to the data results, using bar graphs for visualisation.

The rest of this paper is structured as follows: Section 2 contains details of related work and examined existing research in the areas of Wikidata and Twitter data processing. Section 3 summarises the three phases of the Design and Implementation process of this work. Section 4 discusses the Results and Evaluation of the experiment, testing the research hypothesis and examines the strengths and weaknesses of the results and evaluation. Finally, Section 5 contains the Conclusion, summarising the results found and examining exciting areas of future work that could be completed.

2 RELATED WORK

Trending topics are the most popular talked about items at any point in time over a social media network [18]. As events are more frequently talked about, they become more popular for a period of time where it then peaks and falls. There are a number of areas to be considered when deciding on the approach to use for trend analysis. The data studied may be streamed or static data and may even be a combination of both. The data to be used in the study impacts which Natural Language Processing (NLP) techniques are selected, varying depending on whether the data is structured or unstructured. In addition, the data selected for analysis determines which statistical measures are best suited in identifying text similarity. The following section will examine previous research completed in these areas.

Microblogging sites are a platform used by individuals to share information and voice opinions on any topic, such as current events, products or services. Real-time analysis of social media data is increasingly studied due to the use of social media in sharing information and connecting people, assisting companies to make decisions [11] and gain insight in to their customers' views on their products to help improve such products [20]. There is a large amount of unstructured data available today on microblogging sites like Twitter,

review sites and information articles. There are two hundred million members which produce approximately four hundred million tweets daily, [19] sharing their thoughts, views and opinions [10]. In recent years there have been many studies completed on Twitter data for analysis in areas such as, predicting stock behaviour [15]; book recommendations from twitter feeds [4]; sentiment analysis [1] [11]; burstiness [2]; longevity of trending topic with predictions [18]; and trend identification [6].

Wikidata launched in 2012 as a knowledge base of the Wikimedia foundation, storing its knowledge in the structured format of subject-predicate-object statements [12] organized and structured into pages [8]. Wikidata content is language independent supporting four-hundred-and-ten languages [14]. "The data model of Wikidata is based on a directed, labelled graph where entities are connected by edges that are labelled properties." [5]. There are two types of entities including items and properties. Each item entity has a page relating to a subject area, for example, a city, person or a university where it's data can be entered, edited or viewed[8].

Full streaming of twitter data is used in studies, such as trend identification [15], [6], [23] and sentiment analysis [20], and will be used within this study. The approach to retrieving data from Twitter has varied across studies including examining historic data by topic [18], [1], as well as streaming the data by topic [24], [4]. In one study, streaming twitter data by the topic over a ten-month period monitoring lifetime of trending topics found, if a topic had six hundred or more tweets each day in the first week it would last a month, and how positive and negative sentiments impacted whether they would trend for more than one month [18]. Twitter provides a Streaming API that allows for the collection of publicly available tweets and this approach will be used to retrieve Twitter data. Wikidata dump files are made available through their website and come in a number of forms. The full Wikidata revision information can be downloaded and the SPARQL endpoint API can be used to extract additional information. SPARQL is a powerful API to access linked data collections that allow for retrieval of precise and insightful information in to the data [5]

There are a number of statistical analysis techniques to be considered when comparing text lists. When considering the statistical measures, the list characteristics are an important consideration. In the case of trend lists, in this study they are non-conjoined lists, where the lists may have different items within their lists. The lists are top-weighted, therefore, the top items of the list are more important than the lower ranked items and indefinite ranking will not be considered where a percentage of items will be examined. The following studies look at list similarity using statistical techniques:

- A study completed examining the correlations of search engine results URL's included Jaccard Ratio similarity distribution measure with different sizes for set similarity that included both with and without confidence levels, find a low overlap of two major search engines where 80% of queries had less than three search engine overlaps [7].
- In a study examining the likeness of Wikipedia pages for near duplicate detection Jaccard's similarity measure was used with a finding of a large amount of duplication within the Wikipedia page content [21].

- Use of Jaccard Coefficient to determine the association between words was implemented in the language Python where it was found to be performing well when measuring the similarity of words [17].
- "Weighted Kendall's Tau is the number of swaps we would perform during the bubble sort in such a way to reduce one permutation to the other" [7] however this does not apply to this research as we not have the same items in each list.

Visualisation is a frequently used technique to display and explain results in a visual format and includes representation of data in formats such as a word cloud for visual representation of most frequent words, [11]; Time Series to show trends over time [4], [3]; moving average to show the tweet rate [4]; and analysis bar graphs [6].

3 IMPLEMENTATION

This chapter details the design, implementation and statistical analysis performed to identify if a correlation exists between Twitter hashtags and Wikidata revisions. The overall process has been split in to three phases as outlined in Figure 1, where the details of each phase's implementation and processing details are outlined.

In phase one, data is streamed from Twitter and its hashtags are extracted and cleaned, applying n-grams before determining the top hashtags tweeted over a seventy-seven-day period. Secondly, for the same period, the Wikidata revisions are extracted from its available data dumps. The Wikidata titles are retrieved using SPARQL, identifying the top revision pages. Finally, statistical comparisons are completed on the top hashtags and Wikidata revisions to identify if a correlation exists. The edit-distance statistics will calculate the similarity between the text items in each list and a statistically significant correlation will be determined on the overall similarity of the text lists. The results are displayed through visualisation techniques.

3.1 Twitter Data Mining

During phase one, Twitter data is streamed to identify the top trending tweets by hashtag. The Twitter real-time data is accessed through its Streaming API using tokens OAuth to ensure secure Authorization data requests. The Streaming API returns the data and notifications in real-time from its public stream result in a JSON format [15].

3.1.1 Storing the Data. The data is stored in JSON format files. The full tweets are retrieved where they contain at least one hashtag (#) and are of locale English where they are stored in batches, with file name labels based on date and time of file creation. When large numbers of tweets were stored in files it was found that the process slowed down, therefore batches were created of five-hundred per file.

3.1.2 Tweet Structure. The entity item hashtag list 'text' values, stored in JSON format, are extracted from the tweet and stored in a .CSV of five-thousand tweets per file for further cleaning and processing. For example, the hashtag 'Florida' is extracted from the hashtag list:

```
"entities":{
  "hashtags":[{"
```

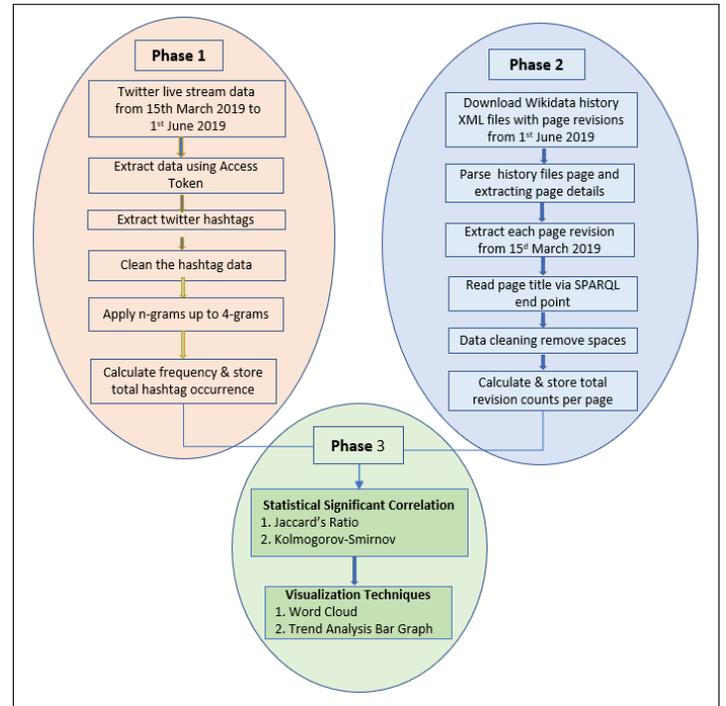


Figure 1: Three project phases of Wikidata and Twitter processing.

```
"text": "Florida",
"indices": [80, 88]},
"urls": [{"url": "https://t.co/Z98KvO6nhB",
"expanded_url": "https://twitter.com/i/web/status/1112821872926777345",
"display_url": "twitter.com/i/web/status/1\u2026",
"indices": [117, 140]}],
"user_mentions": [],
"symbols": []
}
```

3.1.3 Cleaning the Tweet. For each hashtag text extracted, all non-ASCII characters are removed, where only a-z characters remain. This includes removing foreign language characters, numerical data, punctuation etc. For example, hashtags like "text": "trump2020" is updated to 'trump' removing the digits '2020'. The tweet hashtags were split in to words for further processing.

3.1.4 Removing Stop Words from the Tweet. The remaining tweet text is updated to lower case. Stop words are removed using 'nltk.corpus' of the English language. All tweets that are less than two characters are omitted from further processing.

3.1.5 Applying n-grams to the Tweet. Firstly, an n-grams pre-processing step was added to split large hashtags containing five or more words in to smaller groupings of words. For example, if a hashtag contained five words it is split in to three words and two words where, as outlined in the next steps, n-grams are applied.

This process applied n-grams up to 4-grams to each of the extracted tweets as follows:

('social', 45315)
('bbm', 41759)
('stop', 40230)
('bts', 23621)
('love', 14153)
('tweet', 12591)
('exo', 12263)
('mtv', 11799)
('game', 10927)
('thrones', 9818)
('army', 9770)
('day', 9582)
('music', 9505)
('got', 9228)
('chen', 8733)
('play', 8611)
('zubair', 8564)
('fandom', 8400)
('maga', 8090)
('cool', 7881)
('follow', 7871)
('black', 7639)
('fashion', 7542)

Figure 2: Cleaned counted and ordered hashtags 1-grams

- The full hashtag has been split in to words where in the first sample 1-grams is applied to the full Twitter hashtag corpus. This involves taking any split hashtag with more than one word and splitting it in to individual words for processing.
- The process applies 2-grams to each of the applicable extracted tweets as follows. One-word hashtags are included, and two-word hashtags are included. For all hashtags greater than two, the hashtag is split and added for additional processing. This process required, in the case of a three-word hashtag, a twofold process. Firstly, the first two words and the third word are extracted and added and secondly, that the first word and the last two words are extracted and added. In the case of a four-word hashtag, the first two words and second two words were added.
- The process applies 3-grams to each of the applicable extracted tweets as follows. One-word up to three-word hashtags are included without change. For all hashtags greater than three, the hashtag is split and added for additional processing. This process required, in the case of a five-word hashtag, a twofold process. Firstly, that the first three words and the last two words are extracted and added and secondly, that the first two words and the last three words are extracted and added. In the case of a six-word hashtag, the first three words and the last three words were added.

3.1.6 *Counting the Tweets.* For all tweets collected, a count of each tweet occurring in the data set is stored in a .CSV file for further processing (see Figure 2).

3.2 Wikidata Mining and Understanding

In phase two, the English language Wikidata files containing full revision history are downloaded, parsed and prepared for analysis as detailed below.

3.2.1 *Wikidata History Revision Files.* The English language Wikidata compressed files containing full revision history are downloaded and parsed for analysis with a name format 'Wikidata-date-stub-meta-history[num].xml'. These Wikidata dumps are released at regular intervals and available on the Wikidata site. The selected revision files for this study contained the required revision

```
<page>
  <title><Text></title>
  <id><Page Identifier</id>
  <revision>
    [First revision]
  </revision>
  <revision>
    [Second revision]
  </revision>
  [Additional revision information]
</page>
```

Figure 3: Wikidata history file revision structure

information with minimal page data, for example wikidatawiki-20190601-stub-meta-history1.xml.gz. The twenty-seven metadata history files from 1st of June 2019 were downloaded for revision analysis. These stub files contain the page and revision data without text content. These files contained the required revisions and were on average 1.8 GB each when compressed. When uncompressed these files were approximately 12 GB in size, except for the final file wikidatawiki-20190601-stub-meta-history27.xml.gz, with a total size of 15.7 GB when compressed and approximately 78 GB when uncompressed. This final file contains all the revisions since the previous release of the wiki-media-history files containing a larger volume of data to the other twenty-six files. This is the intended design of revision output by Wikidata with this final file continuing to grow where other files should not [22]. Once the files were decompressed the revision data per page were ready to be extracted from each XML file as detailed in the next section.

3.2.2 *Wikidata Download Process.* The basic structure of a page revision is shown in Figure 3 containing the page details and its related revisions outline.

The revision history metadata file consists of many page elements and revision elements of relevance in this study.

The page element <page> contains information about the Wikidata page with its sub elements revisions. This element is used to determine the start of the next page for its revisions to be considered. The sub elements of the page are as follows:

- The page title element <title> is the string representation of its identifier containing a number value. This is added to the output file as 'pagetitle'.
- The element <id> represents the page identifier and is stored as 'pageid' in the output file.
- The <revision> list element contains each revision made to a page and many of its attributes are of relevance in this study to determine the total number of edits applied to a page.
 - The revision represents one revision item <revision> applied to a page.
 - This identifier relates to the revisions identifier and is stored as 'revisionid' in the output file
 - The parent identifier is the <parentid> element links the previous revision and is stored in the output as 'parentid'.
 - The timestamp element is the date the revision occurred and is stored in the output file as 'timestamp'.
 - The comment element contains the summary comment from the user when the revision was introduced and is stored as 'comment' in the output file.

Figure 4 shows a sample of revision data extracted from Wikidata history files where page elements 'pageid' and 'pagetitle' are

pageid	pagetitle	label	revisionid	timestamp	comment	parentid
20604	Q17758	Buttigliera d'Asti	38303	2019-03-17T00:44:01Z	b/* wbsreference-add:2 */ [[Property:P2046]]: 15.76 square kilometre, #quickstatements; [[tool:abs:quickstatements/#/batch/9360] batch #9360]] by [[User:Underlying lk]]	885198340
20604	Q17758	Buttigliera d'Asti	38303	2019-03-16T12:32:09Z	b/* wbscreateclaim-create:1 */ [[Property:P1082]]: 2,564, #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	804491948
20604	Q17758	Buttigliera d'Asti	38303	2019-03-16T12:32:11Z	b/* wbssetqualifier-add:1 */ [[Property:P585]]: 1 January 2018, #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	884690338
20604	Q17758	Buttigliera d'Asti	38303	2019-03-16T12:32:13Z	b/* wbsreference-add:2 */ [[Property:P1082]]: 2,564, #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	884690367
20604	Q17758	Buttigliera d'Asti	38303	2019-03-16T12:32:15Z	b/* wbssetqualifier-add:1 */ [[Property:P459]]: [[Q:15911027]] #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	
20604	Q17758	Buttigliera d'Asti	38303	2019-03-17T00:44:03Z	b/* wbssetqualifier-add:1 */ [[Property:P585]]: 9 October 2011, #quickstatements; [[tool:abs:quickstatements/#/batch/9360] batch #9360]] by [[User:Underlying lk]]	884690423
20605	Q17759	Calamandran	38303	2019-03-17T00:44:05Z	b/* wbsreference-add:2 */ [[Property:P2046]]: 19.16 square kilometre, #quickstatements; [[tool:abs:quickstatements/#/batch/9360] batch #9360]] by [[User:Underlying lk]]	885198392
20605	Q17759	Calamandran	38303	2019-03-16T12:32:18Z	b/* wbscreateclaim-create:1 */ [[Property:P1082]]: 1,745, #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	804491936
20605	Q17759	Calamandran	38303	2019-03-16T12:32:20Z	b/* wbssetqualifier-add:1 */ [[Property:P585]]: 1 January 2018, #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	884690459
20605	Q17759	Calamandran	38303	2019-03-16T12:32:22Z	b/* wbsreference-add:2 */ [[Property:P1082]]: 1,745, #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	884690485
20605	Q17759	Calamandran	38303	2019-03-16T12:32:24Z	b/* wbssetqualifier-add:1 */ [[Property:P459]]: [[Q:15911027]] #quickstatements; [[tool:abs:quickstatements/#/batch/9352] batch #9352]] by [[User:Underlying lk]]	884690512

Figure 4: Wikidata revision with additional title information retrieved using SPARQL endpoint

extracted together with the revision element data. The revision element data includes its 'datetime' stamp if validated to be on or after 15th March 2019 together with its 'comment', 'parentid', and 'revisionid' all stored within .CSV files for additional processing.

The page title required for each revision is not available within the metadata revision history files and is required for processing in this work. However, each revision contains a 'pageid' in the format of Q-ID, that is a unique identifier relating to its page article title. Using SPARQL, its value is read and added to the field 'label' in the output file for later processing. The edit titles are cleaned and the total number of edits per title is recorded during processing.

3.2.3 Wikidata Processing and Assumptions. Python has been used to parse the XML files to extract the Wikidata revision data in to individual records within a .CSV file for additional processing. The attributes extracted per revision were 'pageid', 'pagetitle', 'label', 'revisionid', 'timestamp', 'comment' and 'parentid' for each revision after the data 15th March 2019, from when twitter data was streamed. The following assumptions have been made when processing this data:

Assumption 1: Items without a page identifier are omitted. There are a number of references in the Wikidata history files that do have a Q-ID defined and when retrieved via the SPARQL service from Wikidata, the page does not exist and returns an exception. For these values they are not included in the final result. It was confirmed that these titles did not exist by running the SPARQL query from the provided service.

Assumption 2: User items and contacts omitted. Entries such as 'user' or 'contact the developer' pages have also been omitted from this study. These entries do not have a page ID that can be retrieved

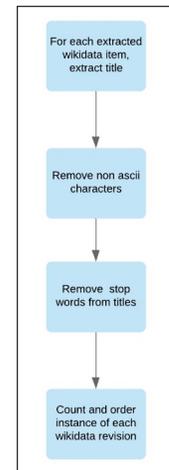


Figure 5: Wikidata additional processing flow diagram

by SPARQL and therefore have not relevance to the study and were omitted from the final analysis result.

3.2.4 Retrieving the Revision Article Title using SPARQL Endpoint. SPARQL is a powerful API with which to access linked data collections that allow for retrieval of precise and insightful information in to the knowledge graph of Wikidata linked data [5]. The revision page title is retrieved and stored per revision item by querying the SPARQL endpoint as shown below.

```

SELECT DISTINCT * WHERE {
  wd: ' + wiki_id + ' rdfs:label ?label .
  FILTER (langMatches(lang(?label), "EN"))
}
LIMIT 1
  
```

The following example returned from the Wikidata revision XML files contained the Q-id value of Q5561905 (the identifier for the Technological University Dublin).

```

SELECT DISTINCT * WHERE {
  wd:Q5561905 rdfs:label ?label .
  FILTER (langMatches(lang(?label), "EN"))
}
LIMIT 1
  
```

3.2.5 Additional Wikidata Processing. Once the Wikidata XML files were parsed, a number of cleaning steps were then required as shown in Figure 5 below.

The non-ASCII characters were extracted from the Wikidata page titles and stop words were removed. This used the same process, English language 'nltk' stop word corpus, that was applied to Twitter. To ensure the comparison with Twitter hashtag data was comparable, all spaces were also removed. Finally, the Wikidata revisions per page were counted to make them available for statistical analysis.

3.2.6 Wikidata Processing Issues. The Wikidata parsing process could not be started until the cut-off date of Twitter collected data and required the data dumps to be made available on the same date. The date selected was 1st of June 2019. During the parsing process, two of the twenty-seven Wikidata dump XML files were fully parsed and eleven were partially parsed. This resulted in the

Wikidata	Total	100%	50%	10%	0.1%
	1867281	270135	135068	27014	270
Twitter	Total	100%	50%	10%	0.1%
1-grams	N/A	52633	26317	5263	53
2-grams	N/A	145133	72567	14513	145
3-grams	N/A	132300	66150	13230	132
4-grams	N/A	128791	64396	12879	129

Figure 6: Page numbers analysed for Wikidata revisions and Twitter hashtags

collection of 1.8 GB of data revisions that occurred within the study period.

3.3 Data Preparation for Statistical Analysis

The statistical analysis process included applying Jaccard’s Ratio and Kolmogorov-Smirnov to a number of datasets formed on a percentage total of the full Twitter hashtags datasets in each n-grams and Wikidata page revisions. The language Python was used to implement Jaccard’s Ratio and Kolmogorov-Smirnov calculation functions executed against these datasets. The percentage of data examined included 0.1%, 10%, 50% and 100% of these datasets.

The volume of revision data collected from Wikidata was 1.8 GB and resulted in out-of-memory exceptions when attempting to run the Kolmogorov-Smirnov against the full dataset. As a result, the lowest frequently items of less than four occurrences were removed from the Wikidata dataset so that the process could be successfully run. As outlined in Figure 6, the total number of unique revisions, once ordered by the most frequent and counted in the full Wikidata dataset, is 1,867,281 unique pages. This was reduced to 270,135 unique pages, equating to 14.5% of the Wikidata unique revision pages, to allow for the Kolmogorov-Smirnov statistical formula to be run successfully. For all further references to 100% of Wikidata this relates to the revised dataset containing 270,135 unique Wikidata pages.

Initially, the data was analysed using the statistical tool Jaccard’s Ratio and Kolmogorov-Smirnov with 100% of the data but, when significant correlation was not found between Wikidata page revisions and Twitter hashtag frequencies, the lower percentage multiples of each data set were also examined. Figure 6 shows the breakdown of the number of both Wikidata items and Twitter hashtag for 100%, 50%, 10% and 0.1% of each dataset. Each counted item in the percentage groupings were counted based on frequency of occurrence. Therefore, each relate to unique references of both the Twitter hashtags and Wikidata pages.

3.4 Jaccard’s Ratio and Kolmogorov-Smirnov Statistical Measures Processing

3.4.1 Kolmogorov-Smirnov. Kolmogorov-Smirnov is a measure of distribution similarity with a range of $[0 - 2]$ where 2 indicates input distribution is equal [7]. This test is a statistical hypothesis test, determining if the two samples of Wikidata pages and Twitter hashtags follow the same distribution. A statistic value is used to determine the probability that the samples are from different distributions where exceeding a confidence level the original null hypothesis H_0 is rejected and so the two samples are from different distributions and thus accepting the alternative hypothesis H_1 . The Kolmogorov-Smirnov p-value is the probability of the null

hypothesis. Where the p-value is less than the significance level of 5% (0.05), the null hypothesis is rejected that both sets of data are from the same distribution, and the alternative hypothesis is accepted.

3.4.2 Jaccard’s Ratio. The statistical measure Jaccard’s Similarity is a statistical hypothesis test used to evaluate the similarity between unordered sets containing a list of items. The Jaccard’s Ratio (similarity) statistical measure was introduced in 1901 and is used to determine set similarity between the two trend lists with a range of $[0 - 1]$, where 0 represents no similarity and 1 indicates the same items exist in each list [7]. Jaccard’s similarity is the total number of items shared (intersection) across both datasets, divided by the total number of items in both datasets (union), to determine the similarity between the sample sets. The items in both lists are unique to the individual list. As a frequency count of both the Twitter hashtags and Wikidata revisions step has been completed as part of the data processing, all words in each dataset used to calculate Jaccard’s similarity are unique. An additional statistical measure Jaccard’s distance is also used within the study to measure dissimilarity between sets. This value is calculated as 1 minus Jaccard’s coefficient.

3.5 Visualisation Statistics

The data evaluation process takes an in-depth look at the results by examining visualisations of key areas in the data. Visualisations were implemented using the language R and Python ‘matplotlib’. The IDE RStudio with the R language was used to create word-cloud charts for the most frequently used Twitter hashtags and Wikidata pages, based on revision frequencies for the studied period. The Python ‘matplotlib’ package was used to create bar charts, giving insight in to the frequency of top trending Twitter hashtags and Wikidata page revisions, as well as to create clusters showing statistical analysis output.

3.6 Data availability, Project Links and Datasets

Using streamed Twitter data¹ meant being confined to the API limit restrictions made available through the Twitter Streaming API. While Twitter provides an enterprise Power Track API for paying customers, this resource cost could not be waived for this research project. The wikidata meta-data-history XML files² containing page revision details could only be parsed³ after the live streaming of twitter data had completed and the revision XML files⁴ were made available by wikimedia. The project ‘TwitterWikidata’ implementation code can be accessed on GitHub⁵.

4 EVALUATION

This section examines and discusses the results found from the statistical tools Jaccard’s Ratio and Kolmogorov-Smirnov, which use quantitative techniques to identify if a significant correlation exists between the top Wikidata revisions and Twitter hashtag

¹<https://drive.google.com/drive/folders/1UYsfniurVl8-uL5emlWXjzqD3JMQVmv4>

²<https://dumps.wikimedia.org/wikidatawiki/>

³<https://drive.google.com/drive/folders/13FnnsSSskVi11KNJptw9pWVB9WBrpuj>

⁴<https://drive.google.com/drive/folders/13FnnsSSskVi11KNJptw9pWVB9WBrpuj>

⁵<https://github.com/D01110788/TwitterWikidata>

The top twenty, 3-grams results are again reflective of the previous n-grams results with 37,302 tweets relating to the termination of the Blackberry messenger app and the TV show ‘Game of Thrones’ related tweets ranked as the third and fifth most popular hashtags over.

When examining 4-grams top ranked list, there is no difference in the top twenty output results where again termination of the Blackberry messenger app and the TV show ‘Game of Thrones’ related tweets ranked as the third and fifth most popular hashtags over the time period. This shows that the top trending hashtags were never greater than three words.

The n-grams visualizations show a consistency across all 4-grams where the termination of the Blackberry messenger application was the most tweeted hashtag across all n-grams. Also, consistently the television show ‘Game of Thrones’ is always high on the frequency list and is spread across a number hashtag entries. This supports the possibility of introducing a bespoke bag of words to allow combining of related tweets like ‘gameofthrones’ occurring 9145 times and ‘got’ occurring 8016 times as shown in figure 9, in to one related hashtag item because they relate to the same topic. Similarly, a bespoke translator could convert ‘bbm’ to ‘Blackberry messenger’ for better comparison to Wikidata. A number of general words also included like ‘music’ and ‘fashion’ could be omitted from the study by the bespoke bag of words during cleaning for the twitter data.

When the Wikidata page items list was examined for ‘Game of Thrones’ related pages, three items were identified from the data extracted. These included revisions on the page ‘listofgameofthronescharacters’, seventy-five revisions on the page ‘gameofthrones’ and 9 revisions on ‘agameofthrones’. Similarly, the data retrieved from Wikidata pages was examined for references to blackberry with twenty-five revisions on the page ‘blackberry’.

4.3 Jaccard’s Ratio and Kolmogorov-Smirnov Statistical Measures Results and Evaluation

This analysis was completed by firstly separating the Twitter hashtags retrieved by its StreamingAPI and created n-grams up to 4-grams grouping of the split hashtag words. Initially the data was analysed using the statistical tool Kolmogorov-Smirnov with 100% of the data made up of 1,867,281 unique pages, but the number of pages included in the calculations was reduced to 14.5% of the overall data with 270,135 unique pages because of performance issues in running the calculation across the full Wikidata page revisions. Within each n-grams groupings the data was grouped by the percentage of data to be analysed. For each n-grams the following coverage split was completed 0.1%, 10%, 50% and 100% of the Twitter data per n-grams. The same split percentage was also applied to the Wikidata sets within each grouping. The data was evaluated using the statistical tools Kolmogorov-Smirnov and Jaccard’s Similarity, to identify if a correlation exists between Wikidata page revisions and Twitter hashtags. The number of unique Twitter hashtags and Wikidata pages are detailed in Section 3.3 Figure 6.

Based on the list characteristics of the Twitter hashtags and Wikidata pages the Jaccard’s Ratio and Kolmogorov-Smirnov statistical measures were used to evaluate the Wikidata revisions and trending Twitter hashtags to determine if a correlation strength existed between the two sets of variables. The finding has accepted the

null hypothesis and rejected the alternative hypothesis indicating a statistically significant correlation was not found between Wikidata page revisions and Twitter hashtags for the studied period when applied across a number of percentages of the datasets including Wikidata items and Twitter hashtag for 100%, 50%, 10% and 0.1% of each dataset. The following section discusses and evaluates the results.

4.3.1 Jaccard’s Ratio Statistical Measure. The Jaccard’s Ratio (similarity) statistical measure was used to determine set similarity between the two trend lists with a range of [0 - 1] where 0 represents no similarity and 1 indicates the same items exist in each list [7]. Jaccard’s Similarity is a statistical hypothesis test evaluating the similarity between unordered sets containing a list of items. In this study the two sets of items are examined each containing string-lists of Wikidata page titles and Twitter hashtags. The analysis for Jaccard’s Ratio was completed for the full corpus of both datasets and run against the four datasets with n-grams applied. Additionally, analysis was completed for Jaccard’s Ratio against 0.1%, 10% and 50% of both datasets. An additional statistical measure Jaccard’s distance is also computed against both list of text-strings used within the study to measure dissimilarity between sets. This value is calculated as 1 minus Jaccard’s coefficient. The results are shown below in Table 1.

Interpreting Jaccard Similarity results will have values in the range of 0-1 where 0 represents no similarity and 1 represents an exact match. Firstly, looking at the results for 1-grams across 0.1%, 10%, 50% and 100%, we can see there is no similarity of words when similarity was calculated on 0.1% of the datasets with a result of 0. This 0.1% of the dataset equated to top 53 unique hashtags from Twitter and the top 270 Wikidata pages ranked by most revisions. This value is also reflected in the Jaccard’s distance where the calculated value is 1 indicating the greatest distance. By increasing the size of the datasets to 10% for 1-grams this equates to 5263 Twitter hashtags and 27,014 Wikidata pages, we can see an increase in similarity to 0.0392 and a reduction in distance with a value of 0.9608. An increase in the similarity continues to occur up to 50% of the 1-grams data sample and reduces again as the dataset is analysed at 100% of the sample. This is an interesting pattern that is reflected across each of the n-grams where the similarity is low on 0.1% of the data in all n-grams datasets analysed and increases in similarity when 50% of the data is analysed, but after 50% the similarity decreases again when 100% of the data was analysed but that 100% distance value is always greater than the recorded 10% n-grams value. Similarly, the pattern established for Jaccard’s Distance as outlined for 1-grams above is consistent across all n-grams with a decrease in distance up to 50% of the sample and an increase again when 100% of the data is analysed for each of the n-grams. The lowest possible similarity was calculated for 1-grams and 4-grams with a value of 0 showing no similarity. The highest similarity was recorded for 1-grams when 50% of the data was examined. This equates to 26,317 unique top Twitter hashtags and 135,068 ordered unique Wikidata pages. A value of 0.0564 was recorded for similarity and a value of 0.9436 recorded for distance with this value being the only one that reached above the 0.05 threshold. The next closest similarity value measured for similarity was also identified within the 1-grams analysis a value of .0417

Test & % of data	1-grams (100%)	2-grams (100%)	3-grams (100%)	4-grams (100%)
Jaccard's Similarity (100%)	0.0417	0.0326	0.0331	0.0335
Jaccard's Distance (100%)	0.9583	0.9674	0.9669	0.9665
Test & % of data	1-grams (50%)	2-grams (50%)	3-grams (50%)	4-grams (50%)
Jaccard's Similarity (top 50%)	0.0564	0.0380	0.0395	0.0399
Jaccard's Distance (top 50%)	0.9436	0.9619	0.9605	0.9601
Test & % of data	1-grams (10%)	2-grams (10%)	3-grams (10%)	4-grams (10%)
Jaccard's Similarity (top 10%)	0.0392	0.0237	0.0256	0.0262
Jaccard's Distance (top 10%)	0.9608	0.9763	0.9744	0.9738
Test & % of data	1-grams (0.1%)	2-grams (0.1%)	3-grams (0.1%)	4-grams (0.1%)
Jaccard's Similarity (0.1%)	0.0	0.0024	0.0025	0.0
Jaccard's Distance (0.1%)	1.0	0.9976	0.9975	1.0

Table 1: Jaccard's Similarity and Jaccard's Distance statistical result

was calculated when 100% of the data was analysed. For remaining distance values calculated they were all less than 0.04

4.3.2 Kolmogorov-Smirnov Statistical Measure. Kolmogorov-Smirnov is a measure of distribution similarity with a range of [0 - 2] where 2 indicates input distribution are equal [7]. This test Kolmogorov-Smirnov is a statistical hypothesis test, determining if the two samples of Wikidata pages and Twitter hashtags come from the same distribution. To evaluate the samples with Kolmogorov-Smirnov, the null hypothesis H0 and H1 hypothesis is defined without knowledge of its result. The null hypothesis and alternative hypothesis were defined in this study as follows:

- Null hypothesis (H0): a correlation does not exist between Wikidata revisions and trending hashtags on Twitter determined by 'Jaccard Ratio' and 'Kolmogorov-Smirnov'.
- Alternative hypothesis (H1): a correlation exists between Wikidata revisions and trending hashtags on Twitter determined by 'Jaccard Ratio' and 'Kolmogorov-Smirnov'.

Next, the data, in terms of probability, is examined to determine if the hypothesis is rejected. A number closer to 0 indicates a likelihood the two samples are coming from the same distribution. If the probability that the samples are from different distributions exceeds a confidence level the original null hypothesis H0 is rejected and so the two samples are from different distributions and thus accepting the alternative hypothesis H1. To evaluate this, a statistic value is calculated using both datasets. The Kolmogorov-Smirnov p-value was also calculated as part of this study used to determine the probability of the null hypothesis. If the p-value is greater than the significance level of 5% (0.05) the null hypothesis is accepted. If the p-value is less than the significance level of 5% (0.05) the null hypothesis is rejected. A low p-value means that the two samples are significantly different. The results for the Kolmogorov-Smirnov statistic and p-value are shown below in Table 2.

When the statistic value and p-value from the Kolmogorov-Smirnov test are examined together, where a small statistic value together with a high p-value then the hypothesis that the distributions of the two samples are the same cannot be rejected. From the results we can see a high p-value across the majority of tested samples where its value is always greater than the 5% threshold of 0.05 as a result this supports the acceptance of the null hypothesis that there is not a statistically significant correlation between Wikidata page revision frequencies and Twitter hashtags for the period

and data evaluated. There is one exception to this when datasets of 2-grams when tested with 100% of the data resulted in a p-value of 0 that is slightly higher than the 0.0661 score calculated for the dataset. The Kolmogorov-Smirnov statistic p-values contained very high levels across all datasets examined. An additional test was completed against a sample of the data by reducing the dataset lists to be of the same length where the Kolmogorov-Smirnov was calculated but it was found reducing the lists to be the same size did not impact the p-value result significantly.

While the outcome of this study rejects the alternative hypothesis that a correlation exists between the data sets examined, improvements identified during this study may have a positive impact on the result. These main suggested improvements include:

- Increased processing power to allow statistical analysis calculations to be run over large datasets. In this study the Wikidata sample was reduced to 14% of the collected sample to run the calculation Kolmogorov-Smirnov without memory errors.
- Introduction of a bespoke bag of words may also improve the results by removing slang words, noisy data words and identifying similar meaning words so that they are combined.

4.4 Hypothesis Outcome

Having analysed Wikidata page titles of the most revised items against Twitter trending hashtags using the statistical tools Jaccard's Ratio and Kolmogorov-Smirnov, the null hypothesis (H0) is accepted, and the alternative hypothesis (H1) has been rejected. This result is based on having identified a high Jaccard's distance value, and a low Jaccard's similarity value between both lists across all data tests completed in the data. Additionally, when the data was examined with the Kolmogorov-Smirnov a high p-value was found together with a low statistic value across supporting acceptance of the null hypothesis.

5 DISCUSSION AND FUTURE WORK

This study has examined Wikidata revisions page titles and streamed Twitter trending hashtags over a seventy-seven-day period to identify if a correlation exists between both sets of data. The results from this study have accepted the null hypothesis that a correlation does not exist between Wikidata revisions and trending hashtags on Twitter validated by the results from the statistical measures

Test & % of data	1-grams (100%)	2-grams (100%)	3-grams (100%)	4-grams (100%)
Kolmogorov-Smirnov p-value (100%)	5.7264e-181	0.0	2.4486e-320	3.5797e-318
Kolmogorov-Smirnov statistic (100%)	0.0687	0.0661	0.0644	0.0648
Test & % of data	1-grams (50%)	2-grams (50%)	3-grams (50%)	4-grams (50%)
Kolmogorov-Smirnov p-value (top 50%)	1.1769e-102	8.2344e-172	4.6453e-154	3.0204e-103
Kolmogorov-Smirnov statistic (top 50%)	0.0557	0.0531	0.0514	0.0521
Test & % of data	1-grams (10%)	2-grams (10%)	3-grams (10%)	4-grams (10%)
Kolmogorov-Smirnov p-value (top 10%)	1.3052e-25	3.9473e-83	7.8451e-78	2.8083e-78
Kolmogorov-Smirnov statistic (top 10%)	0.0813	0.0647	0.0630	0.0634
Test & % of data	1-grams (0.1%)	2-grams (0.1%)	3-grams (0.1%)	4-grams (0.1%)
Kolmogorov-Smirnov p-value (top 0.1%)	0.4183	0.4269	0.1520	0.4183
Kolmogorov-Smirnov statistic (top 0.1%)	0.1294	0.0883	0.1185	0.1294

Table 2: Kolmogorov-Smirnov statistic and p-value results

‘Jaccard Ratio’ and ‘Kolmogorov-Smirnov’. This work has included the mining of live streamed data for a seventy-seven-day period and parsing of Wikidata history revision XML files.

There are many interesting areas where this work could either be extended or improved upon, that were not examined in this study because of limited access to data and time constraints. These are discussed below.

Improvements Through Data Availability. The volume of tweets studied relied on the available downloaded tweets through its publicly available Twitter Streaming API. However, if access was available to the enterprise Power Track API that is currently only available for paying customers this would allow access to a larger volume of steamed tweets to be used in the research.

Improvements Through Extending the Period Analysed. While the initial aim of this study was to download streamed data over a three-month period, the final study examined the tweet downloads over a seventy-seven-day period. Extending the corpus of tweets to the intended three-month period may increase the accuracy of this study; allow for improvement and alternative analysis with Wikidata; or analysis of other sources of available data, for example Wikipedia.

Extending the Techniques of Data Analysis. This work could be extended to include ‘like’ and ‘retweets’ per Twitter item. The impact of a trending hashtag can increase when a tweet is liked or retweeted by high profile individuals and could better identify correlations between trending hashtags and Wikidata revisions.

Creation of a bespoke bag of words to handle individual tweet parts containing slang words or abbreviations may also be added to the study to improve results accuracy.

REFERENCES

[1] Ahuja, S. and Dubey, G. (2017). Clustering and sentiment analysis on twitter data. In *2017 2nd International Conference on Telecommunication and Networks (iQ-Net)*, page 1âĂŞ5.

[2] Al Tamime, R., Giordano, R., and Hall, W. (2018). Observing burstiness in wikipedia articles during new disease outbreaks. In *Proceedings of the 10th ACM Conference on Web Science - WebSci âĂŞ18*, page 117âĂŞ126. ACM Press.

[3] Alsaadi, H. I., Almajmaie, L. K., and Mahmood, W. A. (2017). Forecasting of twitter hashtag temporal dynamics using locally weighted projection regression. In *2017 International Conference on Engineering and Technology (ICET)*, page 1âĂŞ4.

[4] Arulselvi, A. C., Sendhilkumar, S., and Mahalakshmi, S. (2017). Classification of tweets for sentiment and trend analysis. In *2017 International Conference on Intelligent Computing and Control Systems (ICICCS)*, page 566âĂŞ573.

[5] Bielefeldt, A., Gonsior, J., and Kröttsch, M. (2018). Practical linked data access via SPARQL: the case of wikidata. In *Proceedings of the WWW2018 Workshop on*

Linked Data on the Web (LDOW-18), volume 2073 of *CEUR Workshop Proceedings*. CEUR-WS.org.

[6] Doshi, Z., Nadkarni, S., Ajmera, K., and Shah, N. (2017). Tweealyzer: Twitter trend detection and visualization. In *2017 International Conference on Computing, Communication, Control and Automation (ICCUBEA)*, page 1âĂŞ6.

[7] DăĂAlberto, P. and Dasdan, A. (2011). On the weaknesses of correlation measures used for search enginesâĂŞ results (unsupervised comparison of search engine rankings). *arXiv:1107.2691 [cs, stat]*. arXiv: 1107.2691.

[8] Erxleben, F., GăAjnther, M., KrĂüttsch, M., Mendez, J., and VrandeĂiĂĈ, D. (2014). *Introducing Wikidata to the Linked Data Web*, volume 8796, page 50âĂŞ65. Springer International Publishing.

[9] Goldfarb, D. and Merkl, D. (2018). Visualizing art historical developments using the getty ulan, wikipedia and wikidata. In *2018 22nd International Conference Information Visualisation (IV)*, page 459âĂŞ466. IEEE.

[10] Hao, M., Rohrdantz, C., Janetzko, H., Dayal, U., Keim, D. A., Haug, L., and Hsu, M. (2011). Visual sentiment analysis on twitter data streams. In *2011 IEEE Conference on Visual Analytics Science and Technology (VAST)*, page 277âĂŞ278.

[11] HariPriya, A. and Kumari, S. (2017). Real time analysis of top trending event on twitter: Lexicon based approach. In *2017 8th International Conference on Computing, Communication and Networking Technologies (ICCCNT)*, page 1âĂŞ4.

[12] Heindorf, S., Potthast, M., Engels, G., and Stein, B. (2017). Overview of the wikidata vandalism detection task at wsdm cup 2017. *arXiv:1712.05956 [cs]*. arXiv: 1712.05956.

[13] Jhandir, M. Z., Tenvir, A., On, B.-W., Lee, I., and Choi, G. S. (2017). Controversy detection in wikipedia using semantic dissimilarity. *Information Sciences*, 418âĂŞ419:581âĂŞ600.

[14] Kaffee, L.-A. and Simperl, E. (2018). Analysis of editorsâĂŞ languages in wikidata. In *Proceedings of the 14th International Symposium on Open Collaboration - OpenSym âĂŞ18*, page 1âĂŞ5. ACM Press.

[15] Li, Q., Zhou, B., and Liu, Q. (2016). Can twitter posts predict stock behavior?: A study of stock market with twitter social emotion. In *2016 IEEE International Conference on Cloud Computing and Big Data Analysis (ICCCBDA)*, page 359âĂŞ364.

[16] Medelyan, O., Milne, D., Legg, C., and Witten, I. H. (2009). Mining meaning from wikipedia. *International Journal of Human-Computer Studies*, 67(9):716âĂŞ754.

[17] Niwattanakul, S., Singthongchai, J., Naenudorn, E., and Wanapu, S. (2013). Using of jaccard coefficient for keywords similarity. *Hong Kong*, page 6.

[18] Sundar, D. S. and Kankanala, M. (2015). Analyzing and predicting lifetime of trends using social networks. In *2015 International Conference on Computer Communication and Informatics (ICCCI)*, page 1âĂŞ7.

[19] Tajalizadeh, H. and Boostani, R. (2019). A novel stream clustering framework for spam detection in twitter. *IEEE Transactions on Computational Social Systems*, page 1âĂŞ10.

[20] Trupthi, M., Pabboju, S., and Narasimha, G. (2017). Sentiment analysis on twitter using streaming api. In *2017 IEEE 7th International Advance Computing Conference (IACC)*, page 915âĂŞ919.

[21] Weissman, S., Ayhan, S., Bradley, J., and Lin, J. (2015). Identifying duplicate and contradictory information in wikipedia. In *Proceedings of the 15th ACM/IEEE-CE on Joint Conference on Digital Libraries - JCDL âĂŞ15*, page 57âĂŞ60. ACM Press.

[22] Wikimedia (2018). Data dumps/faq - meta.

[23] Xie, W., Zhu, F., Jiang, J., Lim, E., and Wang, K. (2013). Topicsketch: Real-time bursty topic detection from twitter. In *2013 IEEE 13th International Conference on Data Mining*, page 837âĂŞ846.

[24] Zangerle, E., Schmidhammer, G., and Specht, G. (2015). wikipedia on twitter: Analyzing tweets about wikipedia. In *Proceedings of the 11th International Symposium on Open Collaboration, OpenSym âĂŞ15*, page 14:1âĂŞ14:8. ACM.