

UNIVERSITY OF ARKANSAS AT LITTLE ROCK

DOCTORAL THESIS

---

**A Framework for Collecting, Extracting  
and Managing Event Identity  
Information from Textual Content in  
Social Media**

---

*Author:*

Debanjan Mahata

*Supervisor:*

Dr. John R. Talburt

*A thesis submitted in fulfilment of the requirements  
for the degree of Doctor of Philosophy*

*in*

Integrated Computing  
Information Quality Track  
Department of Information Science

April 2015

# **Declaration of Authorship**

I, Debanjan Mahata, declare that this thesis titled, 'A Framework for Collecting, Extracting and Managing Event Identity Information from Short Social Media Text' and the work presented in it are my own. I confirm that:

- This work was done wholly or mainly while in candidature for a research degree at this University.
- Where any part of this thesis has previously been submitted for a degree or any other qualification at this University or any other institution, this has been clearly stated.
- Where I have consulted the published work of others, this is always clearly attributed.
- Where I have quoted from the work of others, the source is always given. With the exception of such quotations, this thesis is entirely my own work.
- I have acknowledged all main sources of help.
- Where the thesis is based on work done by myself jointly with others, I have made clear exactly what was done by others and what I have contributed myself.

Signed:

---

Date:

*“Torture the data, and it will confess to anything.”*

Ronald Coase, Economics, Nobel Prize Laureate

## *Abstract*

With the popularity of social media platforms like Facebook, Twitter, Google Plus, etc, there has been voluminous growth in the digital footprints of real-life events in the Internet. The user generated colloquial and concise textual content related to different types of real-life events, produced in these websites, acts as a hotbed for researchers and organizations for extracting valuable and meaningful information. There has been significant improvement in natural language processing techniques for mining formal and long textual content often found in blogs and newspaper articles. But, it is still a challenging task to mine textual information from the social media channels producing terse, informal and noisy text with an unusual structure. For an event of interest it is necessary to detect and store event-specific signals from the noisy social media channels that allows to distinctively identify that event among all others and characterizes it for drawing actionable insights. These event-specific cues also forms its identity in the unstructured domain of social media. This identity information when mined and analyzed in a timely manner has tremendous applications in the areas of real-life event analysis, opinion mining, reference tracking, recommendation engines, cyber security, event management, among others. Thus, there is a need of a generic framework that can collect short textual content related to real-life events, extract information from them and maintain the information persistently for performing data analytics tasks, and tracking newly produced content as an event evolves. The patent pending work presented in this thesis establishes the design and implementation of an extendable framework enabling collecting, extracting and persistently managing identity information of real-life events from short textual content produced in social media. Towards this objective a pipeline of data processing components going through repeated processing cycles - *Event Identity Information Management Life Cycle* (EIIM) is proposed. A novel persistent graph data structure - *EventIdentityInfoGraph* representing the identity information structure of an event is implemented that forms the core component of the EIIM cycle. Mutually reinforcing relationships between event-specific social media posts, hashtags, text units, URLs and users, forming the vertices of the graph and denoting *event identity information units*, are defined and quantified. An iterative and scalable algorithm - *EventIdentityInfoRank* is proposed that processes the vertices of the graph and ranks them in terms of event-specific informativeness by leveraging the mutually reinforcing relationships. The ranked *event identity information units* are further used in tracking new event related content and extracting valuable event-specific information. Different components of the framework are tested and validated for real-time event related content generated in social media. The work is concluded by discussing about its novel contributions, practical applications in various other domains and envisaging future directions.

## *Acknowledgements*

I would like to express the deepest appreciation to my committee chair Dr. John R. Talburt, who has shown the attitude and the substance of a genius. He continuously and persuasively conveyed a spirit of adventure in regard to research and scholarship, and an excitement in regard to directing innovation towards practical problems. Without his supervision and constant support this dissertation would not have been possible.

I would like to thank my committee members, Dr. Elizabeth Pierce, Dr. Ningning Wu, Dr. Russel Bruhn and Dr. Mathias Brochhausen, whose high quality contributions in the field of Information Science and Information Quality have inspired me to set high standards in my work, and kept me motivated. I would specially thank Dr. Mathias Brochhausen for devoting his valuable time for discussing about possible applications of ontologies in representing real-life events and the related information content in social media. I strongly consider it as one of the future directions of my research.

In addition, I thank Dr. Vivek Kumar Singh and his team from Banaras Hindu University, India, for collaborating with me and helping me to execute the necessary evaluation tasks in an unbiased way, including manual annotations and feedback. I also acknowledge the support of Mr. Jeff Stinson and Ms. Glediana Rexha for financially supporting the major part of my PhD by allowing me to work as a Graduate Assistant at TechLaunch, University of Arkansas at Little Rock. I would also like to thank Dr. Nitin Agarwal, who supported me in the initial days of my PhD.

I am extremely thankful to Dr. Abhijit Bhattacharyya (Associated Dean, Donaghey College of Engineering and Information Technology), for providing me with advise and encouragement from time to time. This acknowledgement page would be incomplete without thanking the immense support of my friends and family. I thank my parents, wife and friends (specially Pathikrit Bhattacharya, Subhashish Duttachowdhury and Meenakshisundaram Balasubramaniam) for not only their support but for their constant interest in my work and the discussions that I had with them. The conversations with them helped me to understand the information seeking behavior of various people from social media, with different perspectives.

Lastly, I thank University of Arkansas for providing me with the facilities, funds and a congenial environment for working towards my goal of PhD. I also acknowledge the Board Of Trustees Of The University Of Arkansas for filing a provisional patent of my work and encouraging me to pursue a path of innovation.

# Contents

<b>Declaration of Authorship</b>	i
<b>Abstract</b>	iii
<b>Acknowledgements</b>	iv
<b>Contents</b>	v
<b>List of Figures</b>	viii
<b>List of Tables</b>	ix
<b>1 Dissertation Overview</b>	1
<b>2 Social Media and Real-life Events</b>	7
2.1 Social Media . . . . .	7
2.2 General Challenges in Social Media Mining . . . . .	7
2.3 Events from Different Perspectives . . . . .	7
2.3.1 Topic Detection and Tracking . . . . .	7
2.3.2 Automatic Content Extraction . . . . .	7
2.3.3 Multimedia Event Detection . . . . .	7
2.4 Events in Social Media . . . . .	7
2.5 Problem of EIIM in Social Media . . . . .	7
<b>3 Literature Review</b>	8
3.1 Identifying High Quality Informative Content in Social Media . . . . .	8
3.2 Entity Resolution . . . . .	10
3.3 Event Identification in News Text . . . . .	13
3.4 Event Identification in Social Media . . . . .	14
<b>4 Event Identity Information Management (EIIM) Life Cycle for Social Media</b>	16
4.1 Identity Integrity . . . . .	17
4.2 Event Reference Collection . . . . .	17
4.2.1 Blog Reference Collection . . . . .	18
4.2.2 Microblog Reference Collection . . . . .	19

4.3	Event Reference Preparation . . . . .	19
4.4	Event Information Quality . . . . .	20
4.5	Event Identity Information Capture . . . . .	23
4.6	Event Identity Information Structure . . . . .	24
4.7	Event Identity Information Processing . . . . .	24
4.7.1	TwitterEventInfoRank . . . . .	28
4.8	Event Reference Resolution . . . . .	37
4.9	Event Analytics . . . . .	37
	Top Five Event-specific Informative Hashtags for Sydney Siege Event . . . . .	37
	Top Five Event-specific Informative Text Units for Sydney Siege Event . . . . .	37
	Top Five Event-specific Informative URLs for Sydney Siege Event . . . . .	38
	Top Five Event-specific Informative Tweet Excerpts for Sydney Siege Event . . . . .	38
	Three Randomly Selected Tweets for Top Three Event-specific Informative Users posting about Sydney Siege Event . . . . .	39
	Top Five Event-specific Informative Hashtags for Millions March NYC Event . . . . .	39
	Top Five Event-specific Informative Text Units for Millions March NYC Event . . . . .	40
	Top Five Event-specific Informative URLs for Millions March NYC Event . . . . .	40
	Top Five Event-specific Informative Tweet Excerpts for Millions March NYC Event . . . . .	40
	Three Randomly Selected Tweets for Top Three Event-specific Informative Users posting about Millions March NYC Event for a particular hour . . . . .	41
5	<b>Potential Applications of the EIIM Framework</b>	42
5.1	Event Monitoring and Analysis . . . . .	42
5.2	Event Information Retrieval . . . . .	43
5.3	Opinion and Review Mining . . . . .	43
5.4	Recommender Systems . . . . .	44
5.5	Event Management and Marketing . . . . .	45
5.6	Social Media Data Integration . . . . .	46
6	<b>Conclusion and Future Work</b>	47
6.1	Conclusion . . . . .	47
6.2	Future Work . . . . .	47
6.2.1	Summarizing Event Related Content . . . . .	47
6.2.2	Identifying Insightful Opinionated Content Related to Events . . . . .	47
6.2.3	Event Topic Modeling . . . . .	47
6.2.4	Event-specific Recommendations . . . . .	47
6.2.5	Distributed Processing of EventIdentityInfoGraph . . . . .	47
6.2.6	Event Ontology for Social Media . . . . .	47

<b>A Appendix Title Here</b>	<b>48</b>
------------------------------	-----------

<b>Bibliography</b>	<b>49</b>
---------------------	-----------

# List of Figures

1.1	Event Identity Information Management (EIIM) Life Cycle for user generated textual content in social media . . . . .	3
4.1	Identity Integrity component of the EIIM life cycle. . . . .	16
4.2	Event Reference Collection component of the EIIM life cycle. . . . .	17
4.3	Event Reference Preparation component of the EIIM life cycle. . . . .	20
4.4	Event Information Quality component of the EIIM life cycle. . . . .	21
4.5	Content characteristics of informative and non-informative tweets related to events. . . . .	22
4.6	Event Identity Information Capture component of the EIIM life cycle. . .	22
4.7	Event Identity Information Structure component of the EIIM life cycle. .	25
4.8	Event Identity Information Processing component of the EIIM life cycle. .	25
4.9	Mutual Reinforcement Chains in Twitter for an event. . . . .	27
4.10	Performance comparison of ranking techniques using NDCG scores. . . . .	33
4.11	Performance comparison of ranking techniques using NDCG scores. . . . .	34
4.12	Performance comparison of ranking techniques using NDCG scores. . . . .	34
4.13	Performance comparison of ranking techniques using NDCG scores. . . . .	35
4.14	Performance comparison of ranking techniques using precision scores. . . . .	35
4.15	Performance comparison of ranking techniques using precision scores. . . . .	36
4.16	Event Reference Resolution component of the EIIM life cycle. . . . .	36
4.17	Event Analytics component of the EIIM life cycle. . . . .	37

# List of Tables

4.1	Details of Data Collected . . . . .	18
4.2	Details of data collected for analyzing event related tweet content. . . . .	19
4.3	Tweet features for content informativeness. . . . .	23
4.4	Evaluation measures for logistic regression model. . . . .	23
4.5	Affinity scores of edges between vertices of TwitterEventInfoGraph . . . . .	26
4.6	Avg IIC scores and total avg scores of annotations for Millions March NYC event. . . . .	32
4.7	Avg IIC scores and total avg scores of annotations for Sydney Siege event.	33

*Dedicated to my parents, wife and my entire family for their  
endless love, support and encouragement.*

# Chapter 1

## Dissertation Overview

Social media has brought a paradigm shift in the way people communicate with each other. It has gone from being just a medium to a global medium of communication between people. Different types of social media platforms provide multiple venues to people for sharing first-hand experiences and exchange information about real-life events. It has become an indispensable source for disseminating news and real-time information about current events, using websites like Twitter, Facebook, Instagram, Flickr, Youtube, Vine, etc, that allow users to post short textual messages accompanied with images and videos. At the same time users also share their detailed citizen journalistic experiences in the form of diaries through different blogging platforms like Blogger, Wordpress, Medium, etc. Studies have shown the importance of different social media platforms as a news circulation service [1], and a source for gauging public interest and opinions [2–5]. Its efficacy as a real-time citizen-journalistic source of information has been recently harnessed in detection, extraction and analysis of real-life events [6–8]. The activities of users producing content in social media has also been studied for gaining deep insights about how they group together to form communities around topics related to real-life events [9–11], and lead to collective action [12, 13].

With the popularity of social media there has been proliferation of unstructured textual content about different real-life events, in the Internet. The information gained by identifying and tracking social media content expressing live reporting of an event, recent updates related to the event, insightful opinion about the different named entities (people, place, organization, etc) directly or indirectly involved with the event, summarization of content, among others, could prove to be extremely valuable for monitoring and gaining deeper actionable insights. There are tremendous applications in the areas of real-life event analysis, event management, opinion mining, reference tracking, online

targeted marketing, recommendation engines, cyber security, enterprise data integration, among others. Thus, there is a need of a generic framework that has the following capabilities:

- can collect different types of textual content produced in social media related to an event
- extract information that acts as an identity of the event used for characterizing it
- maintain the extracted event identity information persistently for resolving constantly produced new content and discovering important event-specific information.

The problem of collecting and extracting event identity information from social media is very similar to the task of event detection and tracking from newswires [14, 15]. However, in this thesis, we add new components of creating identity structures of an event and managing the tracked information persistently over time. In order to make our task well defined we avoid the task of detecting unidentified events, and instead track a pre-specified set of events. Also, the domain of social media poses additional challenges. News articles most often adhere to grammatical, syntactical and formal structures of writing, that are not common in the realm of social media. The user generated content in social media is most often colloquial, short, noisy and lack proper grammatical structures. This makes it a challenging task for the state-of-the-art natural language processing techniques to extract useful information and perform tasks like entity extraction and parts-of-speech tagging that lies at the core of the previous research on event detection and tracking.

The work presented in this thesis establishes the conceptual design and implementation of a framework capable of collecting, extracting and persistently managing event identity information from user generated textual content shared in social media (shown in Figure 1.1). The approach of the presented work is from the perspective of Entity Identity Information Management (EIIM) [16], with basic tenets of information quality at its core. Towards this objective, different challenges of mining high quality information from social media text is discussed and a patent pending novel approach to tackle the challenges for identifying event-specific informative content is explained, which lies at the heart of the framework. It further explores the applications of the research and concludes by pointing to different future directions of the work.

Some of the main contributions of the work are:

FIGURE 1.1: Event Identity Information Management (EIIM) Life Cycle for user generated textual content in social media



- Extending the Entity Identity Information Management model [16] from the closed world domain of Master Data Management (MDM) to the open and unstructured domain of social media.
- Design and implementation of an *Event Identity Information Management* framework that is capable of tracking and identifying event-specific information from long as well as short user generated textual content in social media. Towards this objective a data processing pipeline named *Event Identity Information Management Life Cycle* is developed (Figure 1.1), which is capable of :
  - collecting event related real-time content generated in social media
  - pre-processing them using natural language processing techniques
  - identifying high quality informative sources of information
  - extracting event-specific information in order to create *Event Identity Information Structures* (EIIS) for persistently storing and characterizing the salient and high quality event related information
  - identifying event-specific informative content produced in social media
- Implementation of a supervised classifier in the domain of short and informal social media textual content, for segregating high quality informative messages having higher chances of containing event related information from the low quality non-informative ones.
- Analysis of informative and non-informative event related content from 3.8 million short textual social media messages.

- A novel model that leverages mutually reinforcing relationships between blog posts and named entities mentioned in them, and simultaneously ranks blogs as well as the named entities, allowing identification of event-specific content and further analysis of event-specific information.
- A novel model based on principle of mutual reinforcement that takes into account the semantics of relationships between short textual *social media messages*, *hashtags*, *text units*, *URLs* and *users*, and represent them in a graph structure - *EventIdentityInfoGraph*. A scalable graph processing iterative algorithm -*EventIdentityInfoRank*, is implemented for ranking the nodes of the *EventIdentityInfoGraph*. The algorithm is capable of simultaneously ranking *social media messages*, *hashtags*, *text units*, *URLs* and *users* in terms of event-specific informativeness providing deeper insights into the identity of an event.
- Evaluation of the proposed techniques against popularly used baseline techniques using large scale datasets.

Already published work as well as upcoming publications that represents our contributions related to specific topics covered by the broad area of research as presented in this thesis are given below.

## Related Filed Patent

- A System for Collecting, Ranking and Managing Entity Identity Information from Social Media (US 62135258). Inventors: **Debanjan Mahata** and John R. Talburt, Assignee: The Board Of Trustees Of The University Of Arkansas.

## Related Award

- **Debanjan Mahata** and John R. Talburt. *Chatter that Matter : A Framework for Collecting, Extracting, and Managing Event Identity Information from Short Social Media Text*. Student Research and Creative Works Expo, Graduate Competition, University of Arkansas at Little Rock, April, 2015. (Awarded First Place in Engineering and Information Technology).

## Related Publications

- **Debanjan Mahata**, John R. Talburt and Vivek Kumar Singh; *Identifying and Ranking of Event-specific Entity-centric Informative Content from Twitter*. 20<sup>th</sup> International Conference On Applications Of Natural Language To Information Systems (NLDB 2015), Passau, Germany. 17<sup>th</sup> – 19<sup>th</sup> June, 2015.

- **Debanjan Mahata** and John R. Talburt; *A Framework for Collecting and Managing Entity Identity Information from Social Media.* 19<sup>th</sup> International Conference on Information Quality, Xi'An, China.
- **Debanjan Mahata** and Nitin Agarwal; *Identifying Event-specific Sources from Social Media.* Online Social Media Analysis and Visualization. Lecture Notes in Social Networks, Springer, Kawash, Jalal (Ed). January, 2015.
- Nitin Agarwal, **Debanjan Mahata**, and Huan Liu. *Time-and Event-Driven Modeling of Blogger Influence.* Encyclopedia of Social Network Analysis and Mining. Springer New York, 2014. 2154-2165.
- **Debanjan Mahata** and Nitin Agarwal. *Learning from the crowd: An Evolutionary Mutual Reinforcement Model for Analyzing Events.* Advances in Social Networks Analysis and Mining (ASONAM), 2013 IEEE/ACM International Conference on. IEEE, 2013.
- Nitin Agarwal, and **Debanjan Mahata**. *Grouping the Similar among the Disconnected Bloggers.* Social Media Mining and Social Network Analysis: Emerging Research (2013), 54.
- **Debanjan Mahata**, and Nitin Agarwal. *What does everybody know? identifying event-specific sources from social media.* IEEE Fourth International Conference on Computational Aspects of Social Networks (CASoN), 2012.
- **Debanjan Mahata** and Nitin Agarwal. *Analyzing Event-specific Socio-Technical Behaviors Through the Lens of Social Media.* The International Sunbelt Social Network Conference (Sunbelt XXXII) organized by the International Network for Social Network Analysis (INSNA), March 12-18, 2012, Redondo Beach, California.
- Vivek Kumar Singh, **Debanjan Mahata**, and Rakesh Adhikari. *Mining the blogosphere from a socio-political perspective.* IEEE International Conference on Computer Information Systems and Industrial Management Applications (CISIM), 2010.
- Vivek Kumar Singh, Rakesh Adhikari, and **Debanjan Mahata**. *A clustering and opinion mining approach to socio-political analysis of the blogosphere.* IEEE International Conference on Computational Intelligence and Computing Research (ICCIC), 2010.

## Related Submitted Publications

- **Debanjan Mahata**, John R. Talburt, Vivek Kumar Singh and Rajesh Piryani; *Chatter that Matter: A Framework for Identifying and Ranking Event-specific Informative Tweets*. 18<sup>th</sup> International Conference on Text, Speech and Dialogue, Plzen, Czech Republic (Notification Due: May 10, 2015)
- **Debanjan Mahata**, John R. Talburt and Vivek Kumar Singh; *A Framework for Collecting, Extracting and Managing Event Identity Information from Twitter*. 20<sup>th</sup> International Conference on Information Quality, M.I.T, Boston (Notification Due: April 30, 2015)
- **Debanjan Mahata**, John R. Talburt and Vivek Kumar Singh; *From Chirps to Whistles : Discovering Event-specific Informative Content from Twitter*. Proceedings of the 7<sup>th</sup> Annual ACM Web Science Conference. ACM, 2015, Oxford, England (Notification Due: April 30, 2015)

The rest of the thesis is organized as follows:

Chapter 2 gives an overview of the different social media websites and challenges in mining information from them. It also looks at the different perspectives of defining an event and gives the definition of events in social media as accepted by the presented work. Finally, it defines the problem of Event Identity Information Management from Social Media whose solution and application is extensively discussed throughout the rest of the thesis.

Chapter 3 reviews the existing literature related to the topic of the thesis and highlights the challenges in applying previously available techniques to the domain of social media. It also discusses the similarities and dissimilarities of our work with the previous ones, and identifies the areas of our novel contributions that makes it different from the available techniques.

Chapter 4 presents a detailed discussion of the *Event Identity Information Management Life Cycle*, that is proposed as a solution to the problem that is solved in this thesis. It goes through all the components of the life cycle and gives a detailed explanation of the design choices, implementation and their working.

Chapter 5 highlights the potential real-life application of the *Event Identity Information Management* framework implemented in this thesis.

Chapter 6 draws conclusions of the work presented in this thesis and points to future directions of the work.

## **Chapter 2**

# **Social Media and Real-life Events**

### **2.1 Social Media**

### **2.2 General Challenges in Social Media Mining**

### **2.3 Events from Different Perspectives**

#### **2.3.1 Topic Detection and Tracking**

#### **2.3.2 Automatic Content Extraction**

#### **2.3.3 Multimedia Event Detection**

### **2.4 Events in Social Media**

### **2.5 Problem of EIIM in Social Media**

# Chapter 3

## Literature Review

### 3.1 Identifying High Quality Informative Content in Social Media

Identifying high quality content from the social media feeds that are related to events, is one of the main objectives of our research. As already discussed in Chapter 2, presence of spams, phishing, farm links, promotion of irrelevant content and development of nepotistic relationships are some of the major concerns of information quality in social media. Several effective solutions has been proposed in combating them by [17–20]. Among the different facets of information quality, credibility and trustworthiness of the references are also important. Due to the popularity and its ability to broadcast information at a tremendous pace, social media is also sometimes used by malicious users to spread misinformation and rumors [21]. In such cases, it becomes necessary to assess the credibility and trustworthiness of the information posted. It was showed by Castillo et al. [22] that selection of different types of features and automated classification based on supervised training can be used for detecting credible information about newsworthy topics in Twitter. In one of their works [23] they also proposed a general classification framework for identifying high quality social media content. They took into account the rich meta data like links between items and explicit quality ratings available in Yahoo! Answers website to train a supervised classification model. Credibility of events in Twitter was studied by Gupta et al. [24]. They used PageRank for propagating credibility scores on a heterogeneous network of events, tweets and users. They further constructed a graph between similar events and propagated the scores of the events from the previous network to estimate the credibility of other events. Ranking of tweets based on their credibility during trending events was proposed by Gupta and Kumaraguru [25]. They showed automated extraction of credible information from Twitter, by adopting

supervised learning combined with relevance feedback approach using different features mined from tweets and the users posting them. Truthy<sup>1</sup>, was developed by Ratkiewicz et al. to study information diffusion on Twitter and compute a trustworthiness score for a public stream of micro-blogging updates related to an event to detect political smears, astroturfing, misinformation, and other forms of social pollution [26].

Several mechanisms for ranking social media content in terms of their informativeness have been proposed. Ranking of microblogs like tweets are of particular interest to us as we consider tweets as a representative of short textual content produced in social media. There are many web hosted applications that supplements the default search provided by Twitter in order to effectively retrieve relevant and high quality tweets from different perspectives<sup>2</sup>. On going through these services we found that the most commonly used criteria for ranking tweets are recency, popularity based on retweets and favorite counts, authority of the users posting the tweets and content relevance. Twitter itself uses the popularity of the tweets and features mined from the profile of the users in order to provide personalized search results ordered by recency<sup>3</sup>. A study of different state-of-the-art features and approaches commonly used for ranking tweets has been documented by [27, 28]. Seen<sup>4</sup> is a new state-of-the-art platform that uses a proprietary algorithm named *SeenRank* for ranking event related tweet content for presenting event highlights and summaries. In this work, we consider *SeenRank* as one of our baselines. As the number of retweets of a tweet is widely used for ranking, we also use it as one of our baselines. In the context of our work we name the ranking scheme as *RTRank*.

Apart from the existing real-world search applications, several adaptations of *PageRank* [29] has been proposed by the scientific community for ranking tweets and users in Twitter [30–32]. TweetRank [32] is one such adaptation that ranks tweets by taking into account the direct relationships between tweets in the form of retweets and replies, as well as indirect follower-friend relationships, and usage of similar hashtags. Various learning to rank approaches have been used for ordering tweets retrieved for a given query in terms of their relevance and quality [33–35]. None of these ranking techniques have been devised for event-specific content. An attempt to solve a similar problem presented in this paper was made by [36]. They represented tweets of an event in a cluster and calculated the similarity of individual tweets with the centroid of the cluster. Then they ranked the tweets based on the decreasing value of their similarity. We use this approach as one of our baselines.

---

<sup>1</sup><http://Truthy.indiana.edu/>

<sup>2</sup><http://mashable.com/2009/04/22/twitter-search-services>

<sup>3</sup><https://blog.twitter.com/2011/engineering-behind-twitter%E2%80%99s-new-search-experience>

<sup>4</sup><http://seen.co>

Recently researchers have shown interest in investigating microblog summarization. Experiments have been conducted using both feature-based and graph-based approaches. However, in the context of our work only graph-based approaches are relevant. A comparison of different Twitter summarization algorithms was performed by [37]. Summarization of tweets for sporting events was performed by [38] using the phrase graph algorithm [39]. The popularly used graph-based summarization algorithms are *LexRank* [40] and *TextRank* [41]. Both the algorithms make use of the PageRank scheme of ranking homogeneous nodes in a graph constructed from the text that needs to be summarized and identify the salient text units for producing the summary. Our algorithm uses a similar technique for heterogeneous nodes. Our proposed framework also defines the semantics of the relationships between the nodes differently in the context of tweets. We use both *LexRank* and *TextRank* as evaluation baselines.

We propose implicit mutually reinforcing relationships between tweets, hashtags, text units, users and URLs forming a heterogeneous graph structure (*TwitterEventInfoGraph*), which is novel and makes our work different from any prior work (refer Chapter 4). Scores are assigned to the association between the nodes representing the semantics of their relationships. We implement an iterative algorithm (*TwitterEventInfoRank*) for ranking the nodes of the graph and propagating the event-specific scores of the nodes to its neighboring nodes based on the measure of their association. To our knowledge, this is the first work that identifies novel relationships between different units of content in Twitter and implements a graph-based algorithm for ranking them simultaneously in the context of an event.

## 3.2 Entity Resolution

Entity resolution has been known for more than five decades as the record linkage or the record matching problem in the statistics community [42–44]. In the database community, the problem is defined as merge-purge [45], data de-duplication [46, 47], and instance identification [48]. In the Artificial intelligence community, this problem is described as database hardening [49], and name matching [50]. The names co-reference resolution, identity uncertainty, and duplicate detection are also commonly used to refer to the same task [51]. The term Entity Resolution (ER) first appeared in publications by researchers at the Stanford InfoLab led by Hector Garcia-Molina and is defined as the process of identifying and merging records judged to represent the same real-world entity [52]. In the context of the work presented in this thesis a pre-defined real-life event is considered as an entity. For detailed definition of an event please refer Chapter 2.

Despite the differences in nomenclature used by these authors, the ER process actually comprises five major sub-tasks or activities [53] which are

1. *Entity reference extraction* – locating entity references in unstructured textual information.
2. *Entity reference preparation* – profiling, standardizing, cleaning, and enhancing reference information in preparation for resolution.
3. *Entity reference resolution* – the process or algorithm for determining when references are equivalent, often through direct matching of attributes.
4. *Entity identity management* - creating and maintaining persistent data structures that represent the identities of external entities, the focus of the proposed research.
5. *Entity relationship analysis* – exploring relationships among distinct entities such as household relationships or shared communication.

The *Event Identity Information Management* Life Cycle (Chapter 4) as proposed in this thesis reflects and implements all of the above activities. Historically the focus of ER research has been on Activity 3, the methods for carrying out the resolution process itself. The majority of published research literature falls into this area. The first formal model for resolution was the Fellegi-Sunter Model of Record Linkage [42], which uses a decision-theoretic approach establishing the validity of principles first used in practice by Newcombe [43]. This was followed by the Stanford Entity Resolution Framework (SERF) developed at the Stanford InfoLab [54]. The SERF Model formalizes the generic ER problem as the interaction of two functions for comparing and merging records as black-boxes and defines the conditions required for these functions to give a unique ER result. It also formulates a family of so called “Swoosh” algorithms (G-Swoosh, R-Swoosh, and F-Swoosh) for carrying out the ER process. With the rise of big data a distributed algorithm D-Swoosh [55], was also proposed that can be implemented in a big data environment. More recently the Talburt-Wang Algebraic Model of ER has been proposed [56] that views ER as a problem of partitioning a given set of references.

In addition to research on Activity 3, there has also been extensive research in the area of information extraction (IE) that is directly related to the ER Activity 1, reference extraction. The task of entity extraction is also more relevant to social media, due to the unstructured nature of the content. One of the main emphases in the realm of unstructured textual content for last two decades has been in the task of extracting named entities and categorizing them into types. Competitions like MUC (Message Understanding Conference), CoNLL (Conference on Computational Natural Language

Learning) and ACE (Automatic Content Extraction) spearheaded the development of new techniques in this domain. This led to the development of sophisticated tools like Stanford NER [57], OpenNLP [58], GATE [59], LingPipe [60] and NLTK [61]. Variety of techniques ranging from hand-coded rules, automatic rules, to statistical machine learning techniques like hidden Markov models, maximum entropy and conditional random fields have been proposed. A comprehensive survey of the techniques could be found in [62, 63]. A study of various efforts in extracting information from micro-blogs could be found in [64] and a survey of named entity recognition and classification could be found in [65]. Efforts have been made by the industry in building crowd sourced knowledge bases like freebase [66] and dbpedia [67] for the purpose of entity extraction. A recent effort from the industry for extracting entities from social media and building scalable knowledge bases for doing so has been documented in [68, 69]. The rise of online social networks, has also motivated new research into the ER Activity 5, entity relationship analysis [70]. With the rise of big data, the modern trend is to perform entity resolution process in humongous volumes of data and scale it horizontally [71, 72]. In spite of the recent efforts in the field of entity extraction and resolution from unstructured text, there is no generic framework that solves the problem of persistently collecting and managing entity identity information from social media. The development of Event Identity Information Management from social media is a pioneering effort in the field of entity resolution and would create new avenues of research.

Traditionally, entity identity resolution and management (Activity 4) has been a subject of system administration and management of user identities in large organizations. For the first time [16], showed the intersection of identity management, master data management and entity resolution could be used for managing identities of real-life entities in information systems, that could further play an important role in data integration and information quality. Entity identity management in social media mainly comprises of resolving and integrating profiles of the same person in social networking websites. The FOAF project has been playing an important role in all such efforts [73–75]. A very nice endeavor has been made by the OKKAM project for integrating and managing the multiple entity identifiers in various knowledge bases across the Internet [76]. To our knowledge, we are the first to propose a framework for collecting and extracting identity information of events from social media and use the concepts of entity identity management and entity resolution for persistently managing their identities with respect to time.

### 3.3 Event Identification in News Text

The event detection task [77] in the TDT program (Topic Detection and Tracking), led to significant advancements in the field of event-based organization of broadcast news. Some of the efforts in the TDT program focused on online event detection from continuous and real-time streams of textual news documents in newswires [14, 15]. While others explored the detection of past events from archived news documents [78].

The textual content in news documents are different from the short informal text common in the realm of social media. Most of these documents contain formal text with well-formed grammatical structures, enabling the researchers to rely on the state-of-the-art natural language processing techniques. Named entity extraction and Parts-of-Speech (POS) tagging are among the widely used techniques. Zhang et al. [79] extracted named entities and POS tags from textual news documents, and used them to reweigh tf-idf representations of these documents for the new event detection task. Filatova and Hatzivassiloglou [80] identified named entities corresponding to participants, locations, and times in text documents, and then used the relationships between certain types of entity pairs to detect event content. Hatzivassiloglou et al. [81] used linguistic features (e.g., noun phrase heads, proper names) and learned a logistic regression model for combining these features into a single similarity value. Makkonen et al. [82] extracted meaningful semantic features such as names, time references, and locations, and learned a similarity function that combines these metrics into a single clustering solution.

Extracting events from text has been the focus of numerous studies as part of the NIST initiative for Automatic Content Extraction (ACE) [83, 84]. The ACE program defines event extraction as a supervised task, given a small set of predefined event categories and entities, with the goal of extracting a unified representation of the event from text via attributes (e.g., type, subtype, modality, polarity) and event roles (e.g., person, place, buyer, seller). Ahn [83] divided the event extraction task into different subtasks, including identification of event keyword triggers, and determination of event coreference, and then used machine learning methods to optimize and evaluate the results of each subtask. Ji and Grishman [84] proposed techniques for extracting event content from multiple topically similar documents, instead of the traditional approach of extracting events from individual documents in isolation. In contrast with the predefined templates outlined by ACE, Filatova et al. [85] presented techniques to automatically create templates for event types, referred to as domains, given a set of domain instances (i.e., documents containing information related to events that belong to the domain).

As already discussed, social media documents are extremely concise, noisy and lacks well-established grammatical structures. Therefore, the techniques used in these works

are not always suitable for identification of events from social media. It has been shown that it is extremely challenging for the state-of-the art information extraction algorithms to perform efficiently and give accurate results for micro-blogs [86]. For example, named entity recognition methods typically show 85-90% accuracy on longer texts, but 30-50% on tweets [87]. Therefore, new approaches had to be taken, leading to new techniques for detecting events in social media, which we discuss next.

### 3.4 Event Identification in Social Media

Identification of events and event related content from social media is still in its infancy and needs to be studied more. Several related papers explored the unknown event identification scenario in social media. Weng and Lee [88] proposed wavelet-based signal detection techniques for identifying real-life events from Twitter. These techniques can detect significant bursts or trends in a Twitter data stream. Sankaranarayanan et al. [89] identified late breaking news events on Twitter using clustering, along with a text-based classifier and a set of handpicked news seeders. But they do not take into account the filtering of non-event content, which results in poor performance. Segregating the messages that have high likelihood of containing event related informative content from the ones with chances of having non-informative content, or content that are not at all related to an event are at the core of the work presented in this thesis. Petrovic et al. [90] used locality-sensitive hashing to detect the first tweet associated with an event in a stream of Twitter messages. Rattenbury et al. [91] analyzed the temporal usage distribution of tags to identify tags that correspond to events. Chen and Roy [92] used the time and location associated with Flickr image tags to discover event-related tags with significant distribution patterns (e.g.bursts) in both of these dimensions. Becker et al. [93] defined multi-feature similarity metrics based on the textual and non-textual features associated with the social media documents in order to automatically identify events and their related content. They use the general text-based classifier suggested in [89] and a method for identifying top events suggested by [90] as baseline approaches in their evaluations and achieved better precision scores.

New techniques have been proposed recently for identification of known events in social media. Many of these techniques rely on a set of manually selected terms to retrieve event-related documents from a single social media site [94, 95]. Sakaki et al. [94] developed techniques for identifying earthquake events on Twitter by monitoring keyword triggers (e.g., earthquake or shaking). In their setting, the type of event must be known *a priori*, and should be easily represented using simple keyword queries. Benson et al.

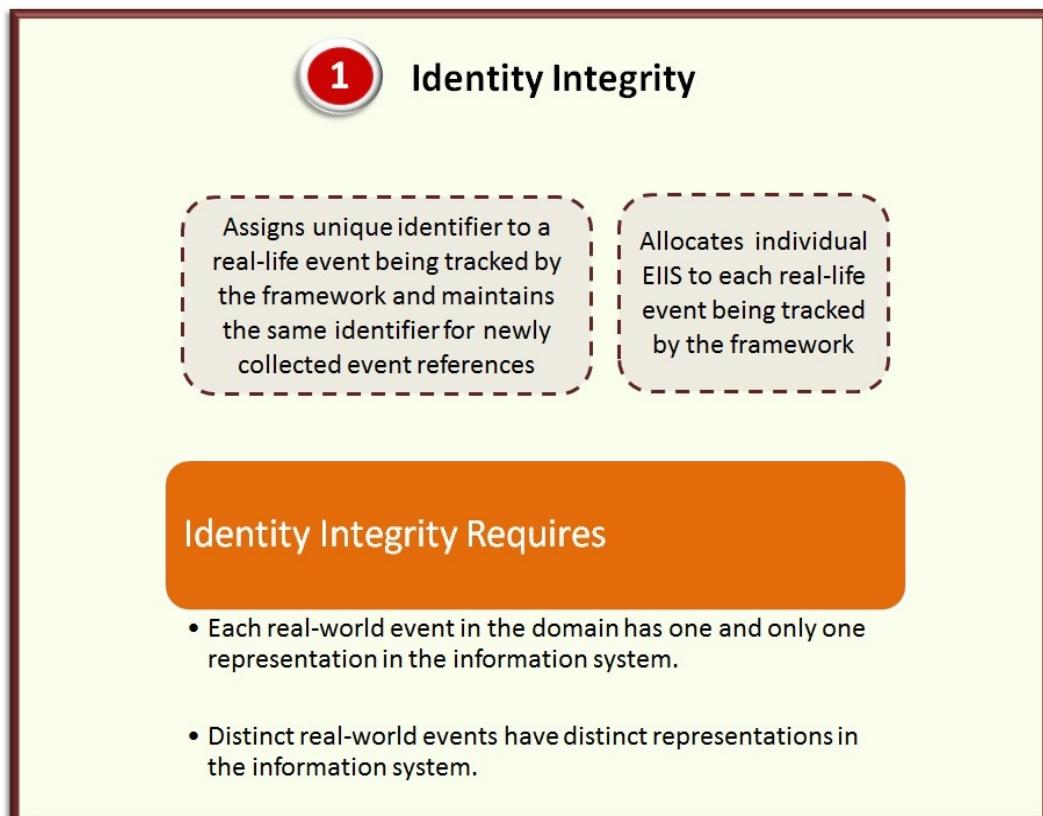
[96] identified Twitter messages for concert events using statistical models to automatically tag artist and venue terms in Twitter messages. Their approach is novel and fully automatic, but it limits the set of identified messages for concert events to those with explicit artist and venue mentions. Most of these approaches are tailored towards one specific social media site. Becker et al. [97] extracts event features, that are often noisy and missing and use them to develop query formulation strategies for retrieving content associated with a planned event from Twitter [98] as well as different social media websites [97].

Our method of tracking events is similar to the idea of identification of known events. We also use predefined hashtags and query words to bootstrap the process of collecting data related to a known set of events. However, we introduce and implement the concept of Event Identity Information Structures that are mapped in a one-to-one mapping with the events that we track. The Event Identity Information Structures persistently stores information that acts as identity of an event as the event evolves with time. This identity information is further processed and ranked in order to identify the top event-specific informative units that is further used for tracking new event related content being generated in different social media channels. Also the emphasis of our research is more on information quality, which is absent in most of the previous research in social media. Instead of just identifying event related content, we identify event-specific informative content. Also, the technique that we develop for identifying event-specific informative content from microblogs (Twitter) leverages hashtags, text units, users, posts and URLs. All these metadata are available in most of the social media websites producing short textual content. Therefore our technique should be applicable to other such platforms. We plan to explore it in the future.

## Chapter 4

# Event Identity Information Management (EIIM) Life Cycle for Social Media

FIGURE 4.1: Identity Integrity component of the EIIM life cycle.

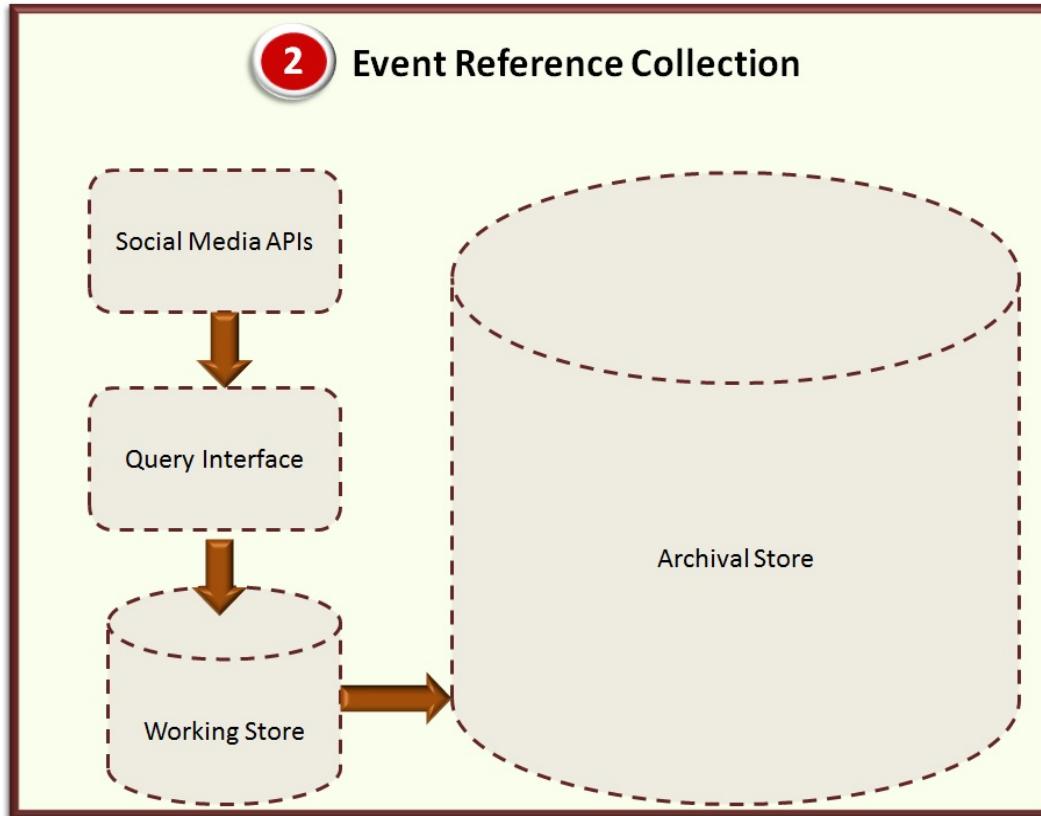


## 4.1 Identity Integrity

One of the fundamental goals of the proposed framework is to maintain a one-to-one correspondence between real-world events being monitored and the Event Identity Information Structure (EIIS) of the corresponding events for ensuring identity integrity. Therefore, a separate EIIS is maintained corresponding to each event. As new events are introduced to the framework, a unique identifier is assigned to them along with the allocation of individual EIIS structures. The framework is expected to maintain the integrity throughout the EIIM life cycle, by consistently assigning the same identifier to the references of a tracked event. Modules of this component assigns 12 byte unique integers known as ObjectId to each event, and is also responsible for maintaining the same ObjectId for event ids of collected references and related EIIS. It is also the functionality of this component to assign the right identifier to the references resolved for an event by the Event Reference Resolution component.

## 4.2 Event Reference Collection

FIGURE 4.2: Event Reference Collection component of the EIIM life cycle.



This component allows the framework to collect event references from different social media websites using its publicly available APIs (Application Programming Interface), and store them in the database after processing them using the next two components of the EIIM life cycle. Due to the semi-structured nature of the collected data, a NOSQL document oriented database management system (MongoDb ) is used for storage. The choice of MongoDb was also driven by its ability to scale horizontally and perform operations on large volumes of data.

For performing the experiments, data was collected from two types of social media websites,

- Blogs - representing the longer genre of textual references.
- Microblogs - representing the shorter genre of textual references.

#### 4.2.1 Blog Reference Collection

TABLE 4.1: Details of Data Collected.

Service Used	Event	Number of Blog Posts
GlobalVoices	Egyptian Revolution	234
	Libyan Revolution	86
	Tunisian Revolution	77
Google Blogger	Egyptian Revolution	579
	Libyan Revolution	600
	Tunisian Revolution	484
Icerocket Blog Search	Egyptian Revolution	5900
	Libyan Revolution	2198
	Tunisian Revolution	1220

Blog posts from GlobalVoices, Blogger<sup>1</sup> and Icerocket Blog Search<sup>2</sup> respectively, were collected for the study. The details of the dataset used is given in Table 4.1. The dataset includes 11,378 blog posts from various blogging platforms like blogspot.com, wordpress.com, livejournal.com, typepad.com, etc. We also filter out the non-english blogs. The data from GlobalVoices is used for constructing event dictionaries, as explained later in Section 4.5. We collect blog posts related to the three events from Blogger using Google Search, and from other blogging platforms using Icerocket blog search.

---

<sup>1</sup><http://blogger.com>

<sup>2</sup><http://icerocket.com>

#### 4.2.2 Microblog Reference Collection

Due to extreme popularity of Twitter, data from it was collected for representing the microblog genre. Four million tweets (approx) related to five different events were collected. Details of the collected event references are provided in Table 2. The tweets were collected over the given period of time, by providing a popular hashtag to the Twitter streaming API (for details about Twitter Data Collection please refer Appendix A).

TABLE 4.2: Details of data collected for analyzing event related tweet content.

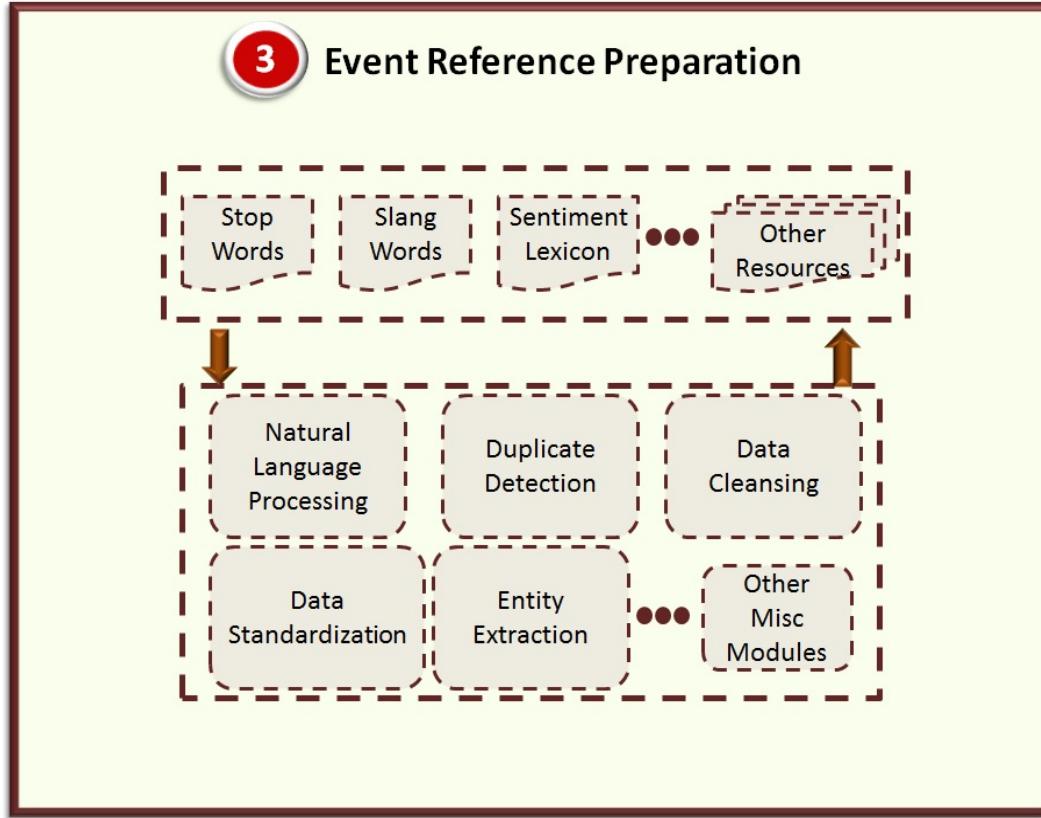
Event	Query Hashtag	No. of Tweets	Time Period
Sochi Winter Games 2014 ( <a href="http://goo.gl/sG4Rqd">http://goo.gl/sG4Rqd</a> )	#sochi2014	1958220	11th Feb, 2014 to 3rd March, 2014
SXSW 2014 ( <a href="http://goo.gl/b6Nd6X">http://goo.gl/b6Nd6X</a> )	sxsw2014	1880557	8th March, 2014 to 16th March, 2014
CPAC 2014 ( <a href="http://goo.gl/9o1KUx">http://goo.gl/9o1KUx</a> )	#cpac2014	18104	7th March, 2014 to 16th March, 2014
Millions March NYC ( <a href="http://goo.gl/I8WR4B">http://goo.gl/I8WR4B</a> )	#millionsmarchnyc	56927	13th Dec, 2014 20:25:43 to 14th Dec, 2014 03:30:41
Sydney Siege ( <a href="http://goo.gl/qLguvG">http://goo.gl/qLguvG</a> )	#sydneysiege	398204	15th Dec, 2014 07:21:16 to 15th Dec, 2014 22:46:45

#### 4.3 Event Reference Preparation

Preprocessing the raw references is an important stage of any data intensive application. This component performs a series of data preparation steps on the collected event references in order to make them suitable for further processing by the other components of the EIIM life cycle.

It performs deduplication of tweets using md5 hashing scheme. Redundant copies of a tweet are filtered out keeping a single copy in the database. Parts-of-speech tagging is done using the default POS tagger available in the NLTK module. A standard list of English stop words is used for eliminating the stop words from the tweet text. All the characters of a tweet are converted into lower case and special characters are removed.

FIGURE 4.3: Event Reference Preparation component of the EIIM life cycle.



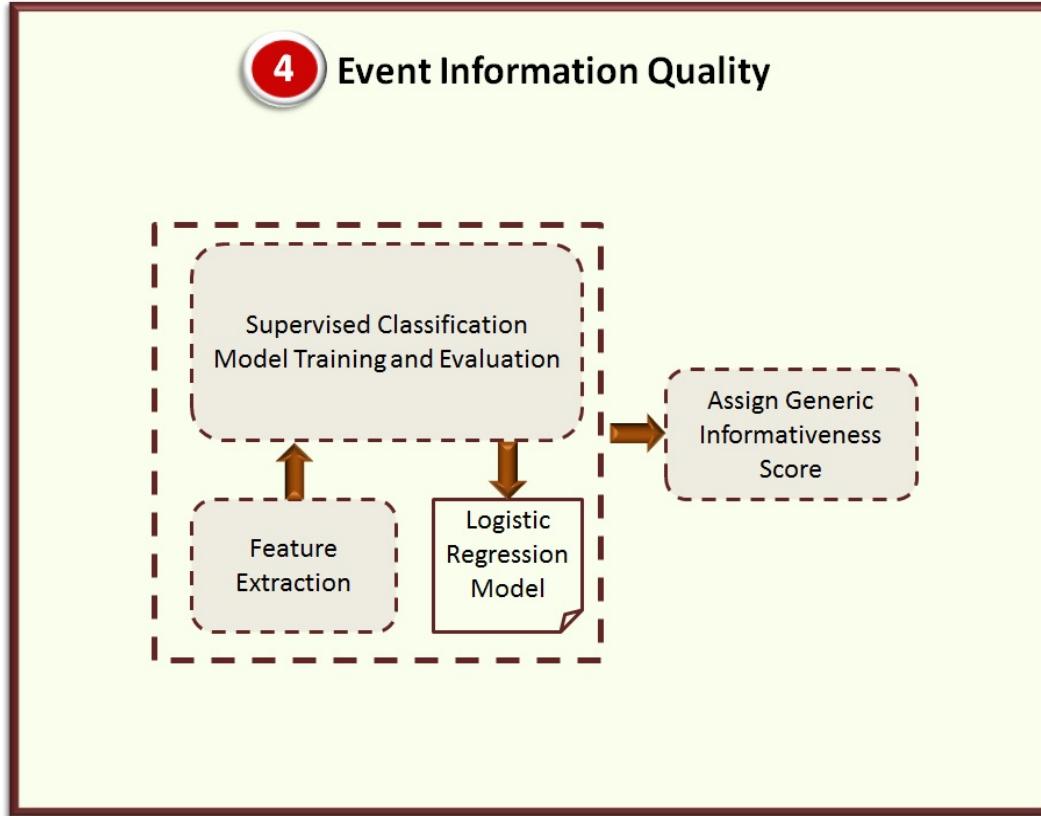
The tweets are tokenized into unigram tokens. User mentions, retweet symbol and URLs are removed during tokenization and are not considered as tokens.

A list of words expressing feelings in the internet, obtained from wefeelfine.org is used for detecting and extracting the feeling words from a tweet. Slang words commonly used in the internet and twitter specific slang publicly shared by FBI is combined together for compiling a list of English slang words. The modules use this list for detecting and extracting the slang words from the tweets, hashtags and text units. Retweet counts, favorite counts, verification information, user follower count, time information and expanded form of the URLs shared in the tweets are extracted from the metadata associated with each tweet, as retrieved using the Twitter API.

#### 4.4 Event Information Quality

This component examines the quality of information present in the tweets collected for the events. It segregates the references having high likelihood of containing good quality event related information from the ones that are less likely to contain or point to good quality information. In order to make a generic module for identifying high quality event

FIGURE 4.4: Event Information Quality component of the EIIM life cycle.



related informative references we implemented a logistic regression classifier trained on a publicly available annotated dataset provided by [28]. The tweets labeled as ‘related and informative’ were assigned a score of 1 and all the other tweets labeled as ‘related-but not informative’, and ‘not related’ were assigned a score of 0. Table 3 lists the features extracted from each tweet. The choice of features was governed by previous works related to identifying high quality information from Twitter as already pointed in the Related Work section. 10-fold cross validation was performed resulting in a model with an accuracy of 76.64The trained model is used for assigning a score between 0 (least informative) and 1 (most informative) to the tweets in real-time. Both the ‘Event Reference Preparation’ and the ‘Event Information Quality’ components work in collaboration with the ‘Event Reference Collection’ component in order to collect, prepare, assign quality score and store the tweets related to an event, obtained from Twitter streaming API, in real-time.

FIGURE 4.5: Content characteristics of informative and non-informative tweets related to events.

		Average No. of Tokens	Average No. of Slang Words	Average Length	Average No. of Top Hashtags	Average No. of Top Nouns	Percentage of URLs
Sochi Winter Games 2014	<i>Informative</i>	8.55	0.47	115.55	0.44	5.14	96.32%
	<i>Non-informative</i>	3.55	0.77	69.92	1.23	1.78	1.04%
SXSW 2014	<i>Informative</i>	7.24	0.62	114.01	0.81	4.36	92.21%
	<i>Non-informative</i>	3.08	0.91	62.64	0.94	1.52	0.34%
CPAC 2014	<i>Informative</i>	6.81	0.53	126.83	1.84	2.42	76.01%
	<i>Non-informative</i>	3.55	0.9	88.65	2.04	2.04	0.68%

FIGURE 4.6: Event Identity Information Capture component of the EIIM life cycle.

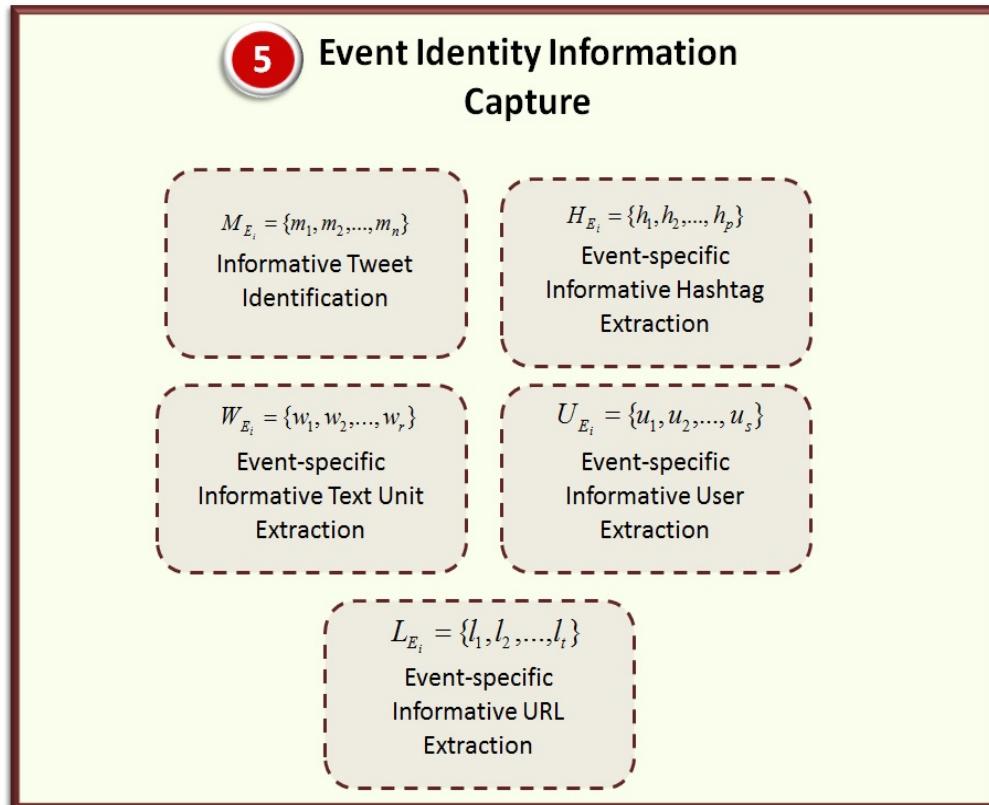


TABLE 4.3: Tweet features for content informativeness.

Has Url, No. of words, No. of stopwords, No. of feeling words, No. of slang words, No. of hashtags, No. of user mentions, Tweet length (No. of characters), No. of unique characters, No. of special characters, Favorite count, Retweet count, Formality, Is tweet verified, No. of nouns, No. of adjectives, No. of verbs, No. of adverbs, No. of pronouns, No. of interjections, No. of articles, No. of prepositions.

TABLE 4.4: Evaluation measures for logistic regression model.

	Precision	Recall	F1-score
<b>Non-informative (0)</b>	0.70	0.49	0.57
<b>Informative (1)</b>	0.78	0.90	0.84
<b>Avg/Total</b>	0.76	0.77	0.75
<b>Accuracy</b>	= 76.64%		

## 4.5 Event Identity Information Capture

It is the component that aids in extracting event identity information units (explained later) from the already processed tweets and build the Event Identity Information Structure (EIIS) for an event. It also enables the framework to set a threshold between 0.0-1.0 for differentiating between high quality informative tweets from low quality non-informative ones related to an event. The event identity information units are then extracted from the high quality informative tweets. In order to understand what might consist of the event identity information units that would represent the EIIS, we conducted a detailed analysis of 3.8 million tweets collected for three events. Details of the data collected are provided in Table 6. The data collection task was accomplished by Event Reference Collection component and was then preprocessed by the Event Reference Preparation component.

The logistic regression model developed for the Event Information Quality component was used for assigning scores to all the 3.8 million tweets in the dataset. The tweets getting a score greater than 0.7 were considered as instances of high quality informative tweets. Those getting a score lesser than 0.3 were considered as instances of low quality non-informative tweets. Average values of different content characteristics of the tweets were calculated. Top ten percent of the frequently occurring hashtags and nouns were considered as top hashtags and top nouns respectively, for the analysis. Some of the characteristics that were prominently different for informative and non-informative tweets are listed in Table 5. As presented in the table, for all the three events, on an average the informative tweets are marked by a higher number of tokens per tweet and greater occurrence of top nouns. The average length of informative tweets is also more than the non-informative ones. The percentage of informative tweets having URLs is strikingly

high. A greater use of slang words is observed in non-informative tweets. However, greater occurrence of top hashtags in non-informative tweets intrigued us to look into the content and obtain a detailed view of it. We observed that a lot of non-informative tweets have used popular hashtags with unrelated content and URLs directing to irrelevant information. This is typical of spam tweets as already reported by [30]. Although not shown due to space constraints, the average number of follower counts for users posting informative tweets was also observed to be higher than the ones posting non-informative ones. The average number of feeling words used in informative tweets were also relatively higher than the feeling words used in the non-informative tweets.

The above observations gave us an idea of how high quality informative content related to events is produced in Twitter and the characteristics that differentiate them from low quality non- informative content. It is now intuitive that the informative tweets are more expressive, formal and lengthier, marked by higher presence of nouns. The high presence of nouns indicates that these tweets also contain information about people, places, organizations, etc, associated with the events, which is vital information about any event and is ideal for representing its identity. Due to the limitations imposed by Twitter on the number of characters in a tweet, the users tend to share URLs along with the textual content that might lead to more information about the event. Also, users with high follower counts tend to post informative tweets. This can also be concluded by the fact that as they have more followers they are encouraged to share informative content. Conversely, since they share informative content they are followed by a large number of other users interested in the content shared by them. Based on the above analysis we decided to build the EIIS for an event composed of the following event identity information units:

## 4.6 Event Identity Information Structure

## 4.7 Event Identity Information Processing

For an event  $E_i$

- a *tweet is an event-specific informative tweet* if it is strongly associated with:
  - (a) *event-specific informative hashtags,*
  - (b) *event-specific informative text units,*
  - (c) *event-specific informative users,*
  - (d) *event-specific informative URLs.*

FIGURE 4.7: Event Identity Information Structure component of the EIIM life cycle.

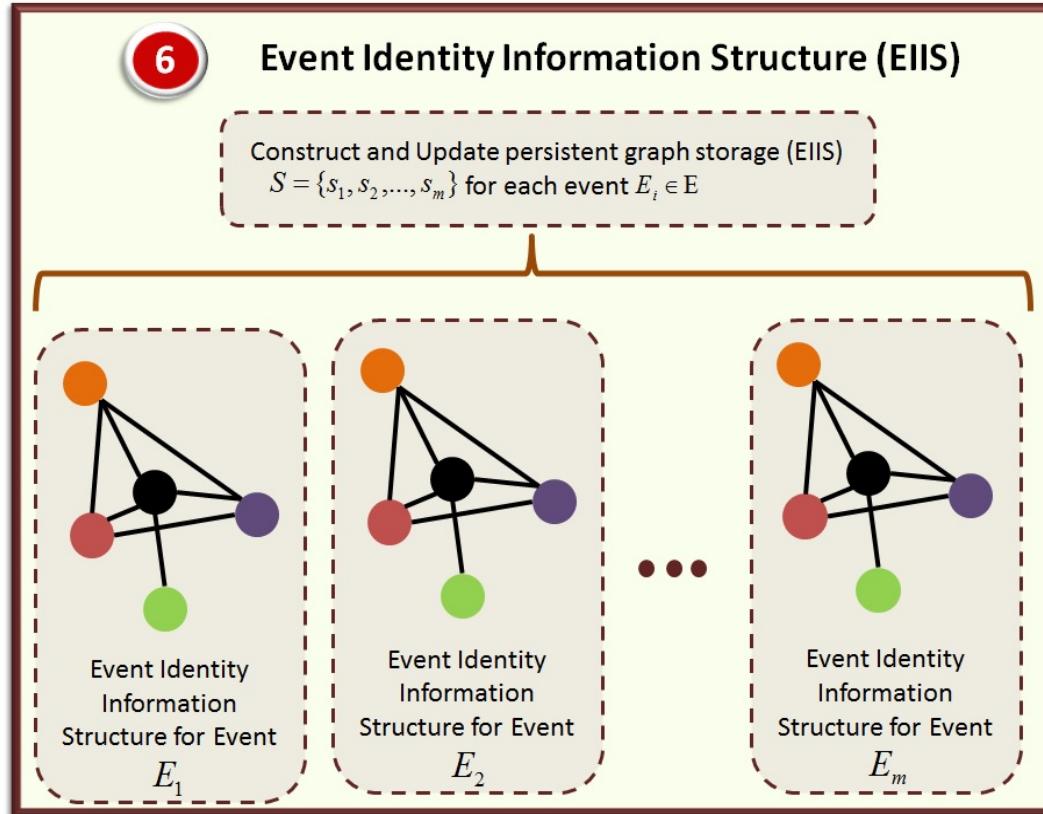


FIGURE 4.8: Event Identity Information Processing component of the EIIM life cycle.

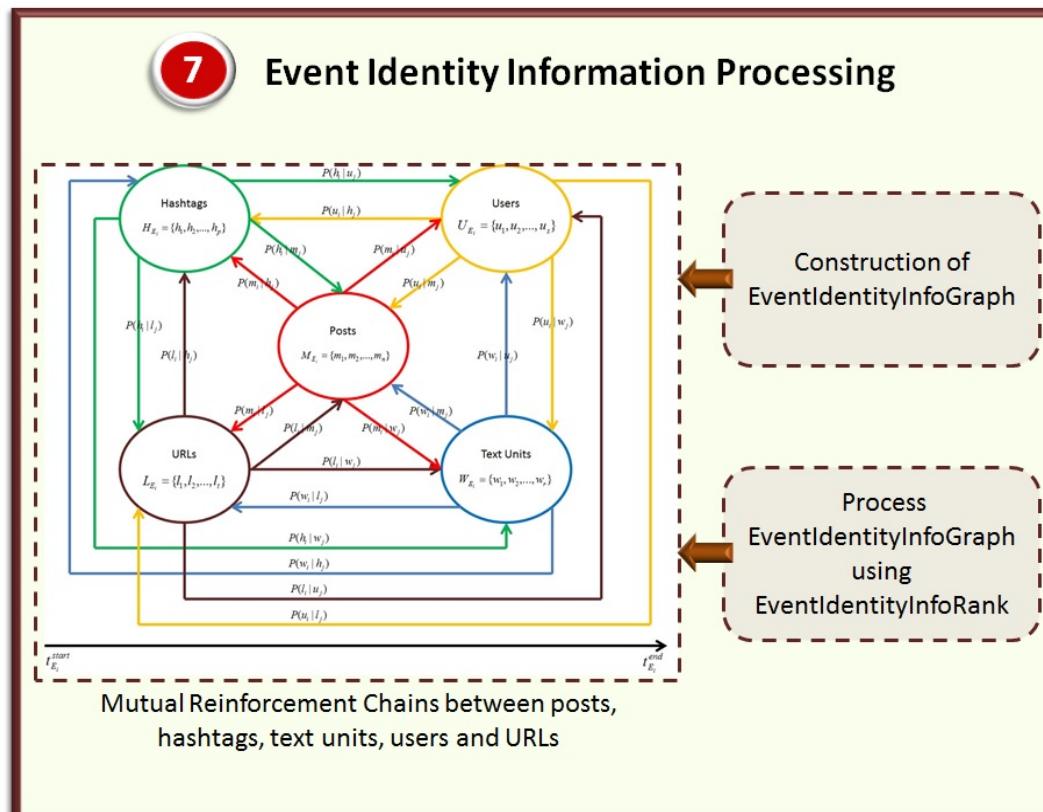


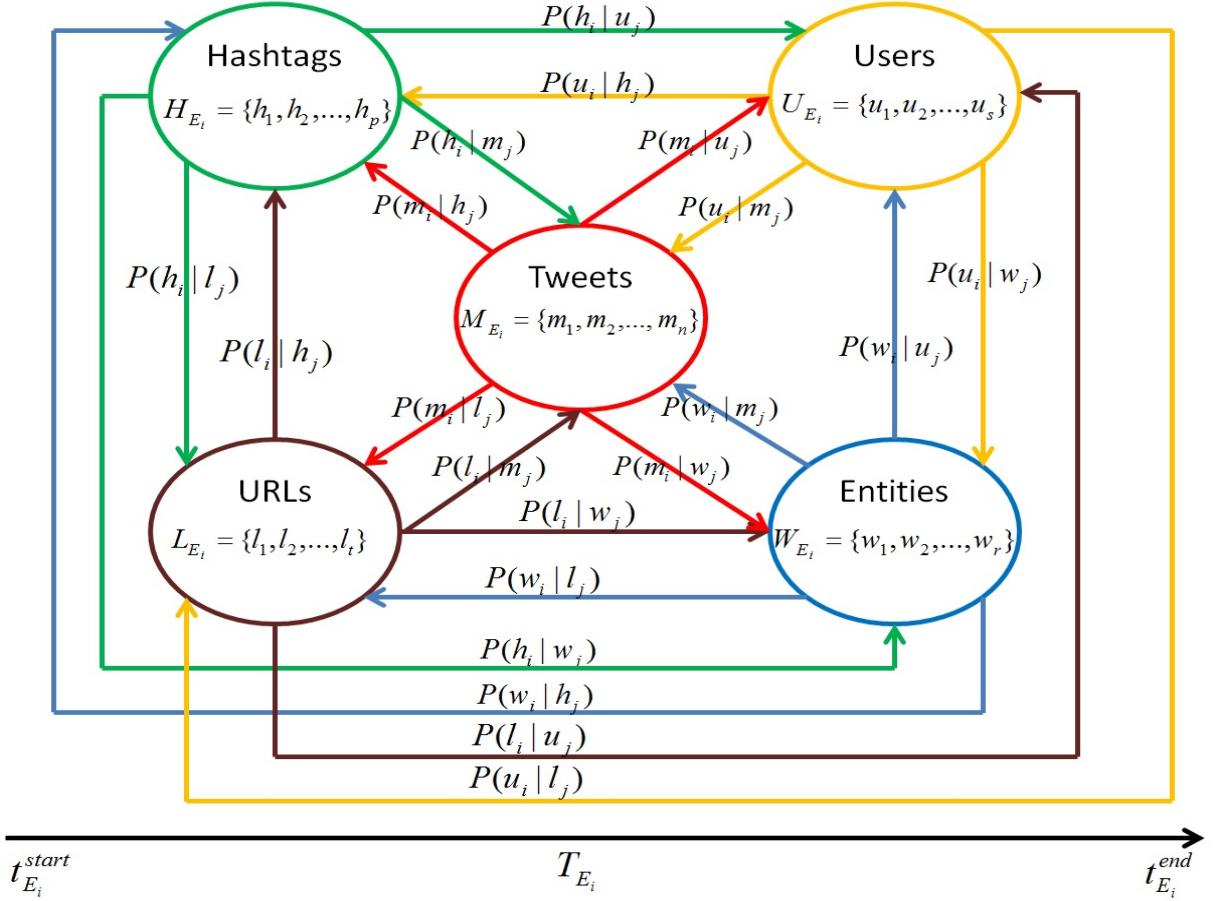
TABLE 4.5: Affinity scores of edges between vertices of TwitterEventInfoGraph

<b>Affinity scores (edge weights) between different vertices <math>\in M_{E_i}, H_{E_i}, W_{E_i}, U_{E_i}, L_{E_i}</math>:</b>
$P(h_i   w_j) = \frac{\text{No. of tweets } h_i \text{ and } w_j \text{ occur together}}{\text{No. of tweets } w_j \text{ occurs}}, P(w_i   h_j) = \frac{\text{No. of tweets } w_i \text{ and } h_j \text{ occur together}}{\text{No. of tweets } h_j \text{ occurs}},$
$P(h_i   l_j) = \frac{\text{No. of tweets } h_i \text{ and } l_j \text{ occur together}}{\text{No. of tweets } l_j \text{ occurs}}, P(l_i   h_j) = \frac{\text{No. of tweets } l_i \text{ and } h_j \text{ occur together}}{\text{No. of tweets } h_j \text{ occurs}},$
$P(h_i   u_j) = \frac{\text{No. of tweets } h_i \text{ and } u_j \text{ occur together}}{\text{No. of tweets } u_j \text{ occurs}}, P(u_i   h_j) = \frac{\text{No. of tweets } u_i \text{ and } h_j \text{ occur together}}{\text{No. of tweets } h_j \text{ occurs}},$
$P(w_i   l_j) = \frac{\text{No. of tweets } w_i \text{ and } l_j \text{ occur together}}{\text{No. of tweets } l_j \text{ occurs}}, P(l_i   w_j) = \frac{\text{No. of tweets } l_i \text{ and } w_j \text{ occur together}}{\text{No. of tweets } w_j \text{ occurs}},$
$P(w_i   u_j) = \frac{\text{No. of tweets } w_i \text{ and } u_j \text{ occur together}}{\text{No. of tweets } u_j \text{ occurs}}, P(u_i   w_j) = \frac{\text{No. of tweets } u_i \text{ and } w_j \text{ occur together}}{\text{No. of tweets } w_j \text{ occurs}},$
$P(u_i   l_j) = \frac{\text{No. of tweets } u_i \text{ and } l_j \text{ occur together}}{\text{No. of tweets } l_j \text{ occurs}}, P(l_i   u_j) = \frac{\text{No. of tweets } l_i \text{ and } u_j \text{ occur together}}{\text{No. of tweets } u_j \text{ occurs}},$
$P(h_i   m_j) = P(m_i   h_j) = P(w_i   m_j) = P(m_i   w_j) = P(u_i   m_j) = P(m_i   u_j) = P(l_i   m_j) = P(m_i   l_j) = 1.0$

**Note:**  $P(h_i | w_j)$  should be read as the probability of occurrence of hashtag  $h_i$  given the occurrence of the text unit  $w_j$  in the stream of tweets  $M_{E_i}$  related to event  $E_i$  collected over the time period  $T_{E_i}$ . Similarly, for others.

- a *hashtag is an event-specific informative hashtag* if it is strongly associated with:
  - (a) *event-specific informative tweets,*
  - (b) *event-specific informative text units,*
  - (c) *event-specific informative users,*
  - (d) *event-specific informative URLs.*
- a *text unit is an event-specific informative text unit* if it is strongly associated with:
  - (a) *event-specific informative tweets,*
  - (b) *event-specific informative hashtags,*
  - (c) *event-specific informative users,*
  - (d) *event-specific informative URLs.*
- a *user is an event-specific informative user* if it is strongly associated with:
  - (a) *event-specific informative tweets,*
  - (b) *event-specific informative hashtags,*
  - (c) *event-specific informative text units,*
  - (d) *event-specific informative URLs.*
- a *URL is an event-specific informative URL* if it is strongly associated with:

FIGURE 4.9: Mutual Reinforcement Chains in Twitter for an event.



- (a) event-specific informative tweets,
- (b) event-specific informative hashtags,
- (c) event-specific informative text units,
- (d) event-specific informative users.

The relationships for an event  $E_i$  as stated above, forms a *Mutual Reinforcement Chain* [99] for the event  $E_i$  as shown in Figure 4.9. We represent this relationship in a graph  $\mathbf{G} = (\mathbf{V}, \mathbf{D})$ , which we call as *TwitterEventInfoGraph*, where  $\mathbf{V} = \mathbf{M}_{\mathbf{E}_i} \cup \mathbf{H}_{\mathbf{E}_i} \cup \mathbf{W}_{\mathbf{E}_i} \cup \mathbf{U}_{\mathbf{E}_i} \cup \mathbf{L}_{\mathbf{E}_i}$ , is the set of vertices and  $\mathbf{D}$  is the set of directed edges between different vertices.

Whenever two vertices are associated, there are two edges between them that are oppositely directed. Each directed edge is assigned a weight, which determines the degree of association of one vertex with the other. The weights for each edge is calculated according to the conditional probabilities given in Table 4.5.

We do not consider an edge between two vertices of same type. That is, we don't connect a tweet with another tweet. Similarly, for hashtags, text units, users and URLs. This

constraint was imposed in order to deal with the nepotistic relationships between high quality content and low quality content introduced by the malicious users for promoting the low quality content. We observe these malicious side effects in the results obtained for *TextRank* explained in Section 6.5.

Next, we explain *TwitterEventInfoRank*.

#### 4.7.1 TwitterEventInfoRank

In this section, we introduce an iterative algorithm that takes into account the mutually reinforcing relationships between the vertices of *TwitterEventInfoGraph* as explained in the previous section and propagates event-specific scores of each vertex to connected vertices across the graph for ranking its vertices ( $\in V$ ) in terms of event-specific informativeness.

We first assign a event-specific score to all the vertices of the graph. Event-specific scores for vertices ( $\in H_{E_i}, W_{E_i}, U_{E_i}, L_{E_i}$ ) are calculated using equations (1-4) as presented in Table 4.5. The tweets ( $\in M_{E_i}$ ) are assigned an initial informativeness score as obtained from the logistic regression model explained in Section 3. The event-specific scores for vertices ( $\in H_{E_i}, W_{E_i}, U_{E_i}, L_{E_i}$ ) and informativeness score for vertices ( $\in M_{E_i}$ ) gives an initial ranking of all the vertices of *TwitterEventInfoGraph*. We aim to refine the initial scores and assign a final score for ranking the vertices by leveraging the mutually reinforcing relationships between them.

$$\text{Score}(h_i) = \frac{\text{freq}(h_i)}{\max\{\text{freq}(h_1), \text{freq}(h_2), \dots, \text{freq}(h_p)\}} \quad (4.1)$$

$$\text{Score}(w_i) = \frac{\text{freq}(w_i)}{\max\{\text{freq}(w_1), \text{freq}(w_2), \dots, \text{freq}(w_r)\}} \quad (4.2)$$

$$\text{Score}(u_i) = \frac{\text{followers}(u_i)}{\max\{\text{followers}(u_1), \dots, \text{followers}(u_r)\}} \quad (4.3)$$

$$\text{Score}(l_i) = \frac{\text{freq}(l_i)}{\max\{\text{freq}(l_1), \text{freq}(l_2), \dots, \text{freq}(l_r)\}} \quad (4.4)$$

The relationships between two different subsets of vertices in graph  $\mathbf{G}$  is denoted by an affinity matrix. For e.g.,  $\mathbf{A}_{\mathbf{E}_i}^{\mathbf{MH}}$  denotes the  $\mathbf{M}_{\mathbf{E}_i} - \mathbf{H}_{\mathbf{E}_i}$  affinity matrix for event  $E_i$ , where  $(\mathbf{i}, \mathbf{j})^{\mathbf{th}}$  entry is the edge weight quantifying the association between  $i^{th}$  tweet ( $\in M_{E_i}$ ) and  $j^{th}$  hashtag ( $\in H_{E_i}$ ), calculated using Table 4.5. Similarly,  $\mathbf{A}_{\mathbf{E}_i}^{\mathbf{WH}}$  denotes

the  $\mathbf{W}_{\mathbf{E}_i} - \mathbf{H}_{\mathbf{E}_i}$  affinity matrix between set of text units  $W_{E_i}$  and set of hashtags  $H_{E_i}$  for event  $E_i$ , and so on.

The rankings of *tweets*, *hashtags*, *text units*, *users* and *URLs* in terms of event-specific informativeness, can be iteratively derived from the Mutual Reinforcement Chain for the event. Let  $R_{E_i}^M$ ,  $R_{E_i}^H$ ,  $R_{E_i}^W$ ,  $R_{E_i}^U$  and  $R_{E_i}^L$  denote the ranking scores for the set of tweets ( $\in M_E$ ), set of hashtags ( $\in H_E$ ), set of text units ( $\in W_E$ ), set of users ( $\in U_E$ ), and set of URLs ( $\in L_E$ ), respectively. Therefore, the Mutual Reinforcement Chain ranking for the  $k^{th}$  iteration can be formulated as follows:

$$R_{E_i}^{M(k+1)} = A_{E_i}^{MM(k)} R_{E_i}^{M(k)} + A_{E_i}^{MH(k)} R_{E_i}^{H(k)} + A_{E_i}^{MW(k)} R_{E_i}^{W(k)} + A_{E_i}^{MU(k)} R_{E_i}^{U(k)} + A_{E_i}^{ML(k)} R_{E_i}^{L(k)} \quad (4.5)$$

$$R_{E_i}^{H(k+1)} = A_{E_i}^{HM(k)} R_{E_i}^{M(k)} + A_{E_i}^{HH(k)} R_{E_i}^{H(k)} + A_{E_i}^{HW(k)} R_{E_i}^{W(k)} + A_{E_i}^{HU(k)} R_{E_i}^{U(k)} + A_{E_i}^{HL(k)} R_{E_i}^{L(k)} \quad (4.6)$$

$$R_{E_i}^{W(k+1)} = A_{E_i}^{WM(k)} R_{E_i}^{M(k)} + A_{E_i}^{WH(k)} R_{E_i}^{H(k)} + A_{E_i}^{WW(k)} R_{E_i}^{W(k)} + A_{E_i}^{WU(k)} R_{E_i}^{U(k)} + A_{E_i}^{WL(k)} R_{E_i}^{L(k)} \quad (4.7)$$

$$R_{E_i}^{U(k+1)} = A_{E_i}^{UM(k)} R_{E_i}^{M(k)} + A_{E_i}^{UH(k)} R_{E_i}^{H(k)} + A_{E_i}^{UW(k)} R_{E_i}^{W(k)} + A_{E_i}^{UU(k)} R_{E_i}^{U(k)} + A_{E_i}^{UL(k)} R_{E_i}^{L(k)} \quad (4.8)$$

$$R_{E_i}^{L(k+1)} = A_{E_i}^{LM(k)} R_{E_i}^{M(k)} + A_{E_i}^{LH(k)} R_{E_i}^{H(k)} + A_{E_i}^{LW(k)} R_{E_i}^{W(k)} + A_{E_i}^{LU(k)} R_{E_i}^{U(k)} + A_{E_i}^{LL(k)} R_{E_i}^{L(k)} \quad (4.9)$$

The equations 5-9 can be represented in the form of a block matrix  $\Delta_{E_i}$ , where,

$$\Delta_{E_i} = \begin{pmatrix} A_{E_i}^{MM} & A_{E_i}^{MH} & A_{E_i}^{MW} & A_{E_i}^{MU} & A_{E_i}^{ML} \\ A_{E_i}^{HM} & A_{E_i}^{HH} & A_{E_i}^{HW} & A_{E_i}^{HU} & A_{E_i}^{HL} \\ A_{E_i}^{WM} & A_{E_i}^{WH} & A_{E_i}^{WW} & A_{E_i}^{WU} & A_{E_i}^{WL} \\ A_{E_i}^{UM} & A_{E_i}^{UH} & A_{E_i}^{UW} & A_{E_i}^{UU} & A_{E_i}^{UL} \\ A_{E_i}^{LM} & A_{E_i}^{LH} & A_{E_i}^{LW} & A_{E_i}^{LU} & A_{E_i}^{LL} \end{pmatrix}$$

Let

$$R_{E_i} = \begin{pmatrix} R_{E_i}^M \\ R_{E_i}^H \\ R_{E_i}^W \\ R_{E_i}^U \\ R_{E_i}^L \end{pmatrix}$$

then,  $R_{E_i}$  can be computed as the dominant eigenvector of  $\Delta_{E_i}$ .

$$\Delta_{E_i} \cdot R_{E_i} = \lambda \cdot R_{E_i} \quad (4.10)$$

In order to guarantee a unique  $R_{E_i}$ ,  $\Delta_{E_i}$  must be forced to be stochastic and irreducible.

To make  $\Delta_{E_i}$  stochastic we divide the value of each element in a column of  $\Delta_{E_i}$  by the sum of the values of all the elements in that column. This finally makes  $\Delta_{E_i}$  column stochastic. We now denote it by  $\hat{\Delta}_{E_i}$ .

Next, we make  $\hat{\Delta}_{E_i}$  irreducible. This is done by making the graph  $G$  strongly connected by adding links from one node to any other node with a probability vector  $p$ . Now,  $\hat{\Delta}_{E_i}$  is transformed to

$$\bar{\Delta}_{E_i} = \alpha \hat{\Delta}_{E_i} + (1 - \alpha) E \quad (4.11)$$

$$E = p \times [1]_{1 \times k} \quad (4.12)$$

where  $0 \leq \alpha \leq 1$  is set to 0.85 according to *PageRank*, and  $k$  is the order of  $\hat{\Delta}_{E_i}$ . We set  $p = [1/k]_{k \times 1}$  by assuming a uniform distribution over all elements. Now,  $\bar{\Delta}_{E_i}$  is stochastic and irreducible and it can be shown that it is also primitive by checking  $\bar{\Delta}_{E_i}^2$  is greater than 0.

Following steps are taken next,

1. We initialize the rank vectors  $(R_{E_i}^{M(0)}, R_{E_i}^{H(0)}, R_{E_i}^{W(0)}, R_{E_i}^{U(0)}, R_{E_i}^{L(0)})$  for each subset of vertices  $(M_{E_i}, H_{E_i}, W_{E_i}, U_{E_i}, L_{E_i})$ . We use the event-specific scores calculated for the set of hashtags, text units, users and urls as their initial scores. All the scores lie between 0 and 1. For the tweets we use the logistic regression model and assign each one of them an initial informativeness score between 0 and 1.

**2.** Then we assign

$$R_{E_i}^0 = \begin{pmatrix} R_{E_i}^{M(0)} \\ R_{E_i}^{H(0)} \\ R_{E_i}^{W(0)} \\ R_{E_i}^{U(0)} \\ R_{E_i}^{L(0)} \end{pmatrix}$$

and normalize  $R_{E_i}^0$  such that  $\| R_{E_i}^0 \|_1 = 1$

- 3.** Apply power iteration method using the same parameters as used in PageRank with the convergence tolerance set at  $1e-08$  and  $\lambda = 0.85$ .
- 4.** We get the final rank vectors for each subset of the vertices  $(R_{E_i}^M, R_{E_i}^H, R_{E_i}^W, R_{E_i}^U, R_{E_i}^L)$  after convergence.
- 5.** We finally obtain the subsets  $\hat{M}_{E_i}, \hat{H}_{E_i}, \hat{W}_{E_i}, \hat{L}_{E_i}, \hat{U}_{E_i}$  consisting of the *tweets*, *hashtags*, *text units*, *URLs* and *users*, respectively arranged in descending order of their final scores.

The final ordered subsets  $\hat{\mathbf{M}}_{\mathbf{E}_i}, \hat{\mathbf{H}}_{\mathbf{E}_i}, \hat{\mathbf{W}}_{\mathbf{E}_i}, \hat{\mathbf{L}}_{\mathbf{E}_i}, \hat{\mathbf{U}}_{\mathbf{E}_i}$ , thus obtained are the tweets, hashtags, text units, URLs and users, ranked in terms of their event-specific informativeness.

**Input** : Sets of vertices  $M_{E_i}, H_{E_i}, W_{E_i}, U_{E_i}, L_{E_i}$  of graph G,  $\alpha = 0.85$ ,  $\varepsilon = 1e - 08$ .

**Output:** Ordered set of vertices  $\hat{M}_{E_i}$ , containing tweets ranked in order of event-specific informative content sharing information about event related entities.

**Steps:**

Initialize rank vectors  $[R_{E_i}^{M(0)}, R_{E_i}^{H(0)}, R_{E_i}^{W(0)}, R_{E_i}^{U(0)}, R_{E_i}^{L(0)}]$ ;

Assign  $R_{E_i}^0 = [R_{E_i}^{M(0)}, R_{E_i}^{H(0)}, R_{E_i}^{W(0)}, R_{E_i}^{U(0)}, R_{E_i}^{L(0)}]^T$ ;

Normalize  $R_{E_i}^0$  such that  $\| R_{E_i}^0 \|_1 = 1$  ;

Construct matrix  $\Delta_{E_i}$ ;

Make matrix  $\Delta_{E_i}$  stochastic and irreducible converting it to  $\overline{\Delta}_{E_i}$ ;

$k \leftarrow 1$

**repeat**

$R_{E_i}^k \leftarrow \overline{\Delta}_{E_i} R_{E_i}^{k-1}$ ;  
 $k \leftarrow k + 1$ ;

**until**  $\| R_{E_i}^k - R_{E_i}^{k-1} \|_1 < \varepsilon$  OR  $k \geq 100$ ;

$R_{E_i}^M \leftarrow R_{E_i}^{M(k)}, R_{E_i}^H \leftarrow R_{E_i}^{H(k)}, R_{E_i}^W \leftarrow R_{E_i}^{W(k)}, R_{E_i}^U \leftarrow R_{E_i}^{U(k)}, R_{E_i}^L \leftarrow R_{E_i}^{L(k)}$ ;

$\hat{M}_{E_i} \leftarrow R_{E_i}^M, \hat{H}_{E_i} \leftarrow R_{E_i}^H, \hat{W}_{E_i} \leftarrow R_{E_i}^W, \hat{U}_{E_i} \leftarrow R_{E_i}^U, \hat{L}_{E_i} \leftarrow R_{E_i}^L$ ;

return  $\hat{M}_{E_i}, \hat{H}_{E_i}, \hat{W}_{E_i}, \hat{U}_{E_i}, \hat{L}_{E_i}$ ;

During the implementation of the *TwitterEventInfoRank* algorithm the slang hashtags were removed. We only considered nouns as the text units and removed the slang words. We already reported in our analysis that non-informative tweets have higher slang content. Therefore, removal of slang hashtags and text units was done in order to obtain high quality results. We also showed higher occurrence of nouns in informative tweets. Also, the occurrence of a noun in a tweet intuitively suggests that the tweet has information about a person, place, or thing. Thus, we only considered the set of nouns extracted from the tweets as the set of text units.

The text units are generic units in the framework and can be changed according to specific requirements. Entities extracted from the textual content of tweets could be experimented, in place of nouns. Since the algorithm uses power iteration method for ranking the vertices of the graph, it could be easily made scalable using mapreduce paradigm [100]. We plan to work on it in the future and implement our framework using hadoop and mapreduce environment.

Since, our proposed framework takes a hybrid approach by using both supervised and unsupervised component, it is easily applicable in situations where an event needs to be tracked over time. The supervised portion assigns an initial generic informativeness score to the tweets for bootstrapping an unsupervised process that finally assigns event-specific informativeness scores. When applied over a time period the method for assigning the initial supervised scores might remain the same and the unsupervised process can change the rankings of the tweet contents as the event evolves.

TABLE 4.6: Avg IIC scores and total avg scores of annotations for Millions March NYC event.

Millions March NYC	IIC	Total Avg Score (1-3)
<b>Top 50 event-specific informative Hashtags</b>	0.786	1.980
<b>Top 50 event-specific informative Text Units</b>	0.880	1.320
<b>Top 50 event-specific informative URLs</b>	0.926	2.560
<b>Top 50 event-specific informative Users</b>	0.700	2.386
<b>Top 100 event-specific informative Tweets</b>	0.760	2.59

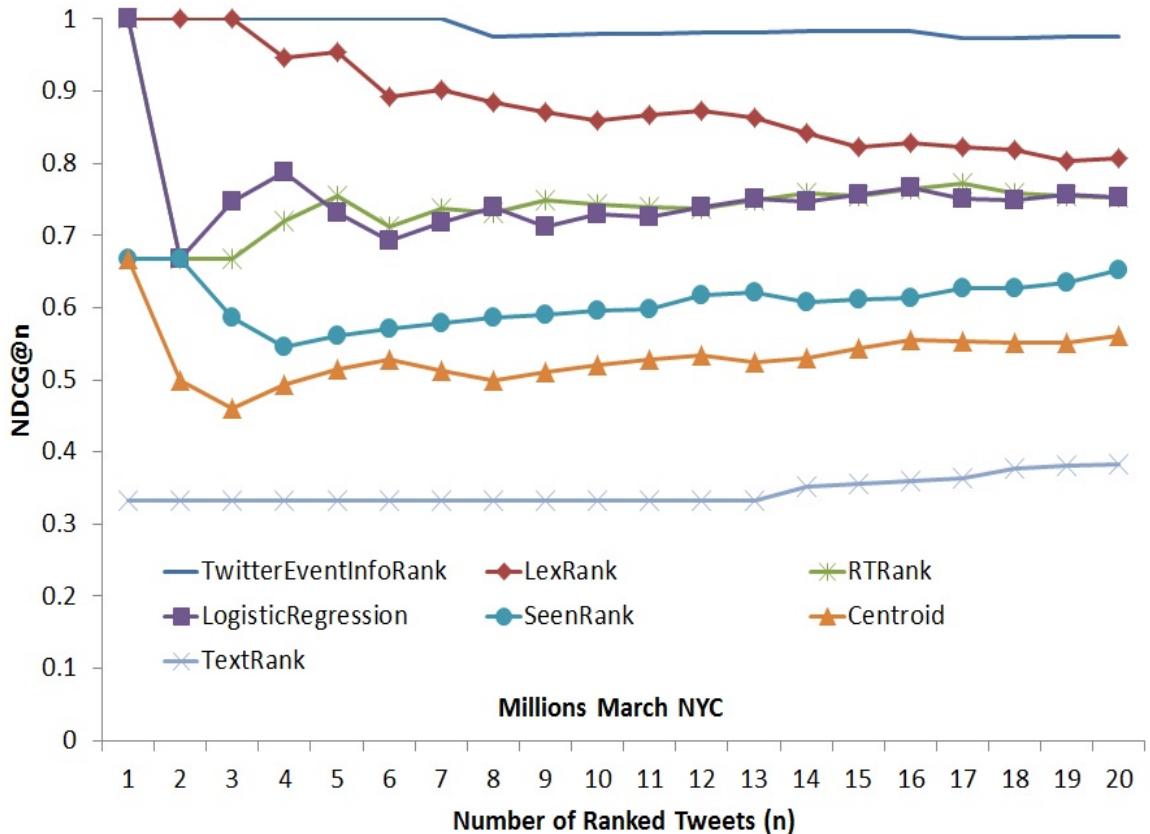


FIGURE 4.10: Performance comparison of ranking techniques using NDCG scores.

TABLE 4.7: Avg IIC scores and total avg scores of annotations for Sydney Siege event.

Sydney Siege	IIC	Total Avg Score (1-3)
Top 50 event-specific informative Hashtags	0.880	2.027
Top 50 event-specific informative Text Units	0.986	1.487
Top 50 event-specific informative URLs	0.893	2.413
Top 50 event-specific informative Users	0.646	2.353
Top 100 event-specific informative Tweets	0.83	2.62

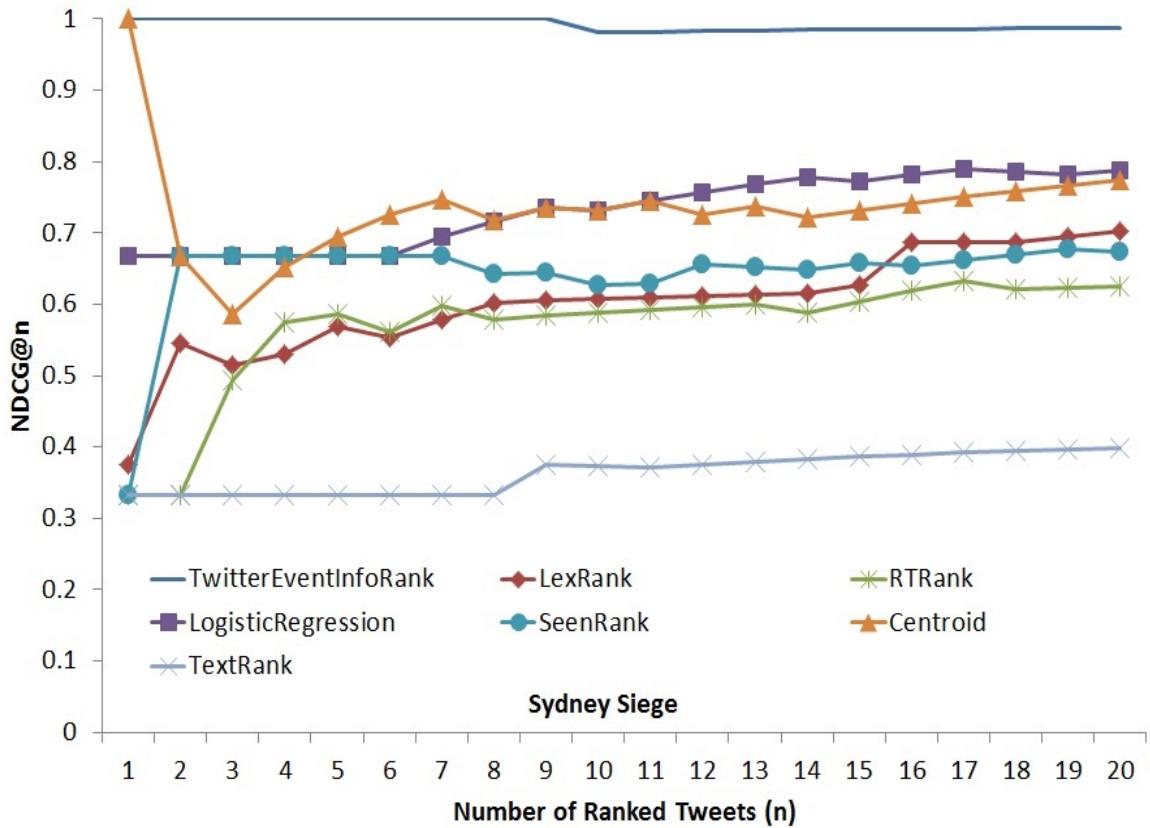


FIGURE 4.11: Performance comparison of ranking techniques using NDCG scores.

Technique	@ 10	@ 20	@ 30	@ 40	@ 50	@ 60	@ 70	@ 80	@ 90	@ 100
TwitterEventInfoRank	0.979	0.975	0.966	0.966	0.957	0.936	0.951	0.960	0.967	0.989
LexRank	0.859	0.807	0.830	0.813	0.822	0.825	0.834	0.878	0.922	0.944
RTRank	0.744	0.752	0.749	0.765	0.792	0.822	0.861	0.870	0.884	0.922
Logistic Regression	0.729	0.753	0.757	0.752	0.757	0.776	0.792	0.839	0.878	0.915
SeenRank	0.595	0.652	0.708	0.733	0.745	0.759	0.801	0.828	0.859	0.884
Centroid	0.519	0.560	0.623	0.658	0.690	0.727	0.747	0.788	0.835	0.857
TextRank	0.333	0.383	0.418	0.468	0.499	0.564	0.633	0.681	0.729	0.782

FIGURE 4.12: Performance comparison of ranking techniques using NDCG scores.

Technique	@ 10	@ 20	@ 30	@ 40	@ 50	@ 60	@ 70	@ 80	@ 90	@ 100
TwitterEventInfoRank	0.980	0.987	0.968	0.957	0.954	0.941	0.946	0.952	0.960	0.990
LexRank	0.607	0.701	0.684	0.707	0.737	0.768	0.764	0.806	0.838	0.868
RTRank	0.588	0.624	0.677	0.716	0.729	0.751	0.769	0.821	0.863	0.880
Logistic Regression	0.730	0.787	0.790	0.791	0.794	0.821	0.855	0.883	0.896	0.927
SeenRank	0.626	0.673	0.728	0.751	0.746	0.779	0.806	0.839	0.869	0.892
Centroid	0.731	0.773	0.779	0.810	0.800	0.779	0.787	0.839	0.880	0.918
TextRank	0.373	0.398	0.485	0.540	0.624	0.664	0.714	0.728	0.764	0.783

FIGURE 4.13: Performance comparison of ranking techniques using NDCG scores.

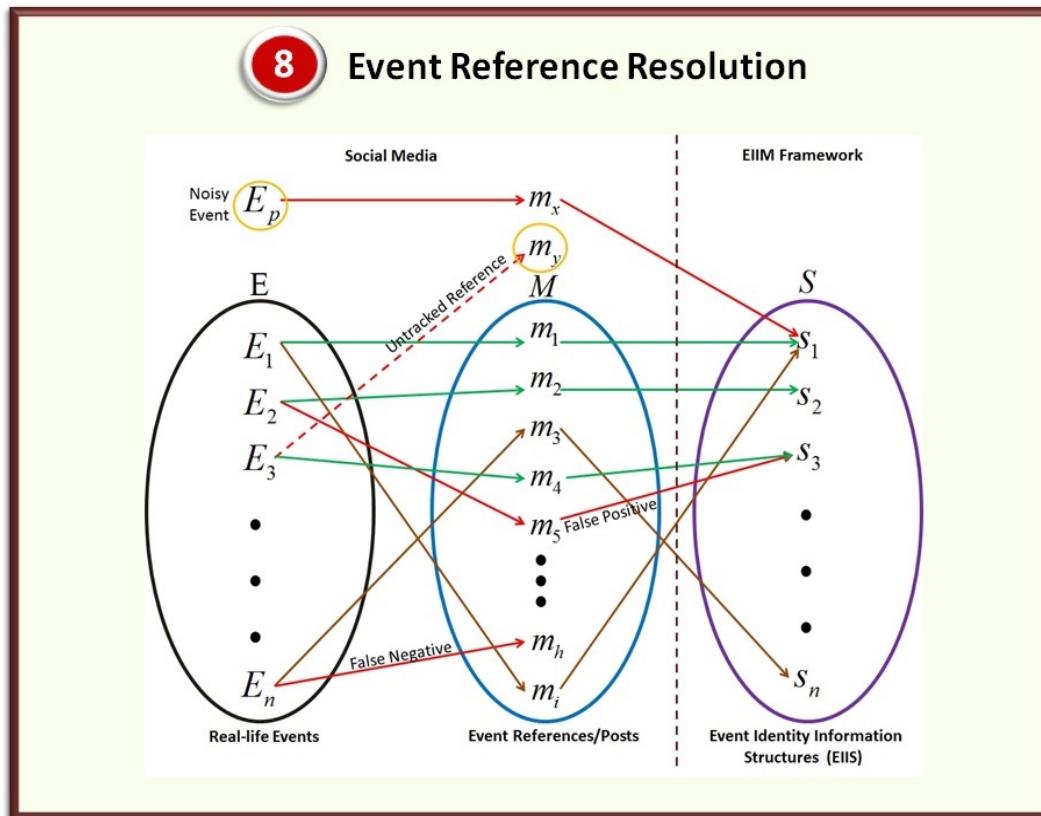
Technique	@ 10	@ 20	@ 30	@ 40	@ 50	@ 60	@ 70	@ 80	@ 90	@ 100
TwitterEventInfoRank	100%	100%	100%	100%	100%	100%	100%	97.5%	96.6%	96.0%
LexRank	90.0%	80.0%	76.6%	65.0%	64.0%	63.3%	60.0%	62.5%	64.4%	64.0%
RTRank	80.0%	85.0%	86.6%	85.0%	86.0%	88.3%	90.0%	91.3%	92.2%	90.0%
Logistic Regression	60.0%	75.0%	76.6%	72.5%	74.0%	71.6%	68.5%	71.3%	71.1%	73.0%
SeenRank	80.0%	85.0%	80.0%	75.0%	72.0%	68.3%	70.0%	67.5%	65.5%	64.0%
Centroid	60.0%	60.0%	60.0%	62.5%	64.0%	66.6%	67.1%	67.5%	70.0%	68.0%
TextRank	0.00%	10.0%	13.3%	25.0%	28.0%	35.0%	42.8%	45.0%	47.8%	51.0%

FIGURE 4.14: Performance comparison of ranking techniques using precision scores.

Technique	@ 10	@ 20	@ 30	@ 40	@ 50	@ 60	@ 70	@ 80	@ 90	@ 100
TwitterEventInfoRank	100%	100%	100%	97.5%	98%	96.7%	95.7%	95.0%	95.5%	96.0%
LexRank	80.0%	85.0%	76.6%	72.5%	76.0%	78.3%	72.8%	73.7%	73.3%	74.0%
RTRank	60.0%	70.0%	76.6%	75.0%	70.0%	71.6%	71.4%	75.0%	73.3%	69.0%
Logistic Regression	100%	100%	100%	97.5%	96.0%	91.6%	92.8%	93.7%	93.3%	92.0%
SeenRank	70.0%	65.0%	70.0%	67.5%	62.0%	61.6%	57.1%	57.5%	55.5%	55.0%
Centroid	70.0%	75.0%	76.7%	82.5%	78.0%	71.6%	65.7%	66.3%	66.7%	66.0%
TextRank	10.0%	5.00%	13.3%	15.0%	22.0%	21.6%	24.3%	21.3%	22.2%	21.0%

FIGURE 4.15: Performance comparison of ranking techniques using precision scores.

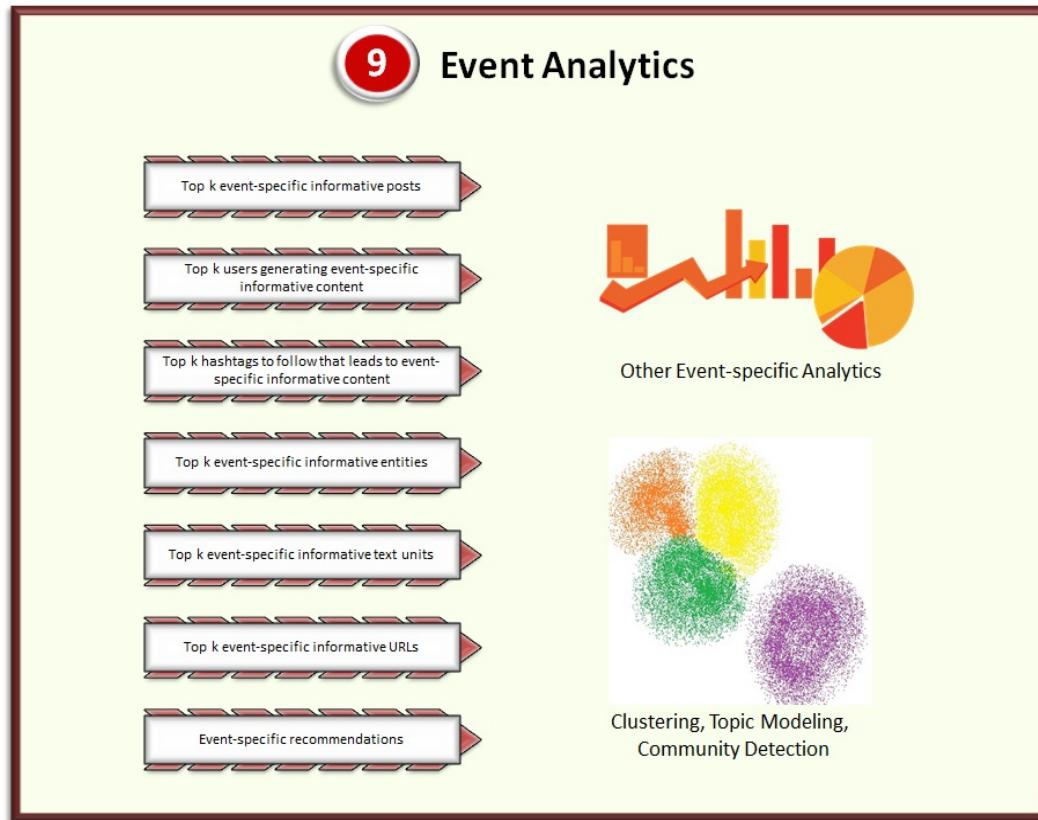
FIGURE 4.16: Event Reference Resolution component of the EIIM life cycle.



## 4.8 Event Reference Resolution

## 4.9 Event Analytics

FIGURE 4.17: Event Analytics component of the EIIM life cycle.



### Top Five Event-specific Informative Hashtags for Sydney Siege Event

1. #sydneysiege
2. #SydneySiege
3. #Sydneysiege
4. #MartinPlace
5. #9News

### Top Five Event-specific Informative Text Units for Sydney Siege Event

1. police

2. sydney
3. reporter
4. lindt
5. isis

### **Top Five Event-specific Informative URLs for Sydney Siege Event**

1. <http://www.cnn.com/2014/12/15/world/asia/australia-sydney-hostage-situation/index.html>
2. <http://www.bbc.co.uk/news/world-australia-30474089>
3. <http://edition.cnn.com/2014/12/15/world/asia/australia-sydney-siege-scene/index.html>
4. <http://rt.com/news/214399-sydney-hostages-islamists-updates/>
5. <http://www.newsroompost.com/138766/sydney-cafe-siege-ends-gunner-among-two-killed>

### **Top Five Event-specific Informative Tweet Excerpts for Sydney Siege Event**

1. RT @faithcnn: Hostage taker in Sydney cafe has demanded 2 things: ISIS flag and; phone call with Australia PM Tony Abbott #SydneySiege <http://t.co/a2vgrn30Xh>
2. Aussie grand mufti and; Imam Council condemn #SydneySiege hostage capture <http://t.co/ED98YKMxqM> - LIVE UPDATES <http://t.co/ED98YKMxqM>
3. RT @PatDollard: #SydneySiege: Hostages Held By Jihadis In Australian Cafe - WATCH LIVE VIDEO COVERAGE <http://t.co/uGxmd7zLpc> #tcot #pjnet <http://t.co/uGxmd7zLpc> sydney-siege-scene/index.html
4. RT @FoxNews: MORE: Police confirm 3 hostages escape Sydney cafe, unknown number remain inside <http://t.co/pcAt91LIdS> #SydneySiege
5. Watch #sydneySiege police conference live as hostages are still being held inside a central Sydney cafe <http://t.co/OjulBqM7w2> #c4news

**Three Randomly Selected Tweets for Top Three Event-specific Informative Users posting about Sydney Siege Event.**

**1. User 1**

- (a) RT @cnni: Hostage taker in Sydney cafe demands ISIS flag and call with Australian PM, Sky News reports. <http://t.co/a2vgrn30Xh> #sydneyseige
- (b) RT @DR\_SHAHID: Hostage taker demands delivery of an #ISIS flag and a conversation with Prime Minister Tony Abbott <http://t.co/xTSDMKCPcD>
- (c) RT @SkyNewsBreak: Update - New South Wales police commissioner confirms five hostages have escaped from the Lindt cafe in Sydney #sydneyseige

**2. User 2**

- (a) RT @smh: NSW Police Deputy Commissioner Catherine Burn will hold a press conference to update on the #SydneySiege at 6.30pm.
- (b) RT @Y7News: Helpful travel advice for commuters heading out of #Sydney's CBD this evening - <http://t.co/aQx2lvSosm> #sydneyseige
- (c) RT @hughwhitfeld: British PM David Cameron informed of #sydneyseige ..UK Foreign Office is in touch with Aus authorities

**3. User 3**

- (a) RT @RT\_com: #SYDNEY: Gunman tall man in late 40s, dressed in black – eyewitness <http://t.co/m51P8dUPhB> #SydneySiege <http://t.co/NvJzFsGrFN>
- (b) RT @NewsAustralia: 2GB's Ray Hadley claims hostage takers in #SydneySiege "wants to speak to Prime Minister Abbott live on radio."
- (c) RT @BBCWorld: "Profoundly shocking" -Australia PM Tony Abbott delivers second #sydneyseige statement. MORE: <http://t.co/VaKt3ZpRZR>

**Top Five Event-specific Informative Hashtags for Millions March NYC Event**

1. #MillionsMarchNYC
2. #BlackLivesMatter
3. #ICantBreathe
4. #ShutItDown
5. #millionsmarchnyc

**Top Five Event-specific Informative Text Units for Millions March NYC Event**

1. police
2. nyc
3. eric
4. protesters
5. nypd

**Top Five Event-specific Informative URLs for Millions March NYC Event**

1. <http://rt.com/usa/214203-protests-police-brutality-nationwide/index.html>
2. [http://mashable.com/2014/12/13/time-lapse-new-york-protest-march/?utm\\_cid=mash-com-Tw-main-link](http://mashable.com/2014/12/13/time-lapse-new-york-protest-march/?utm_cid=mash-com-Tw-main-link)
3. <http://www.cbsnews.com/news/eric-garner-ferguson-missouri-protesters-converge-on-washington/>
4. [http://www.huffingtonpost.com/2014/12/13/millions-march-nyc\\_n\\_6320348.html?ncid=tweetlnk](http://www.huffingtonpost.com/2014/12/13/millions-march-nyc_n_6320348.html?ncid=tweetlnk)
5. <https://www.youtube.com/watch?v=Iz7hkfNmTY&feature=youtu.be>

**Top Five Event-specific Informative Tweet Excerpts for Millions March NYC Event**

1. RT @rightnowio\_feed: Timelapse video reveals massive size of New York City prot... <http://t.co/oHtIhEK969> #Soho #Millionsmarchnyc #NEWYorkC..
2. ”@Breaking911: BREAKING NOW: #NYPD OFFICER INJURED ON THE BROOKLYN BRIDGE BY PROTESTERS THROWING ITEMS AT OFFICERS #MillionsMarchNYC” Great
3. RT @mohkeit: MT @WSJ: march to NYPD headquarters to protest police brutality #MillionsMarchNYC <http://t.co/zhNSngjbkN> <http://t.co/YLMJ8uJnJ>
4. RT @NaomiCampbell: Peaceful March Saturday Dec 13th Washington Square Park NYC 2:00pm march Tell everyone U know #MillionsMarchNYC
5. RT @anregarret: Incredible day! #MillionsMarchNYC On NYPD Headquarters To Protest Police Killings <http://t.co/P2QHvxl9xb> via @blackvoices

**Three Randomly Selected Tweets for Top Three Event-specific Informative Users posting about Millions March NYC Event for a particular hour.**

**1. User 1**

- (a) RT @mashable: Timelapse video reveals massive size of New York City protests <http://t.co/zhqHpkDLk1> #MillionsMarchNYC <http://t.co/WktxssAfDp>
- (b) RT @DahmPublishing: RT@wendycarrillo: Real thugs wear flag pics and Eric Garner's eyes are haunting image #MillionsMarchNYC <http://t.co/7wY...>
- (c) RT @TheRoot: RT @mfmartinez: Protesters continue gathering in Washington Square Park #MillionsMarchNYC #TheRootMOW <http://t.co/IwkQG1KjFg>

**2. User 2**

- (a) RT @roqchams: Thousands march on NYPD headquarters to protest police terrorism <http://t.co/yVyUVYkd9X> <http://t.co/X4QZrfOISh> #MillionsMarchNYC
- (b) RT @NYjusticeleague: Hundreds killed. Ten Demands. One Continued Fight. Sign our petition at: <http://t.co/KETNo6bS0V> #MillionsMarchNYC <http://t.co/...>
- (c) RT @cobismith: Union Square now with NYPD in foreground, #MillionsMarchNYC protesters at right and; US national debt ticker on the left <http://t.co/...>

**3. User 3**

- (a) RT @mashable: Timelapse video reveals massive size of New York City protests <http://t.co/zhqHpkDLk1> #MillionsMarchNYC <http://t.co/WktxssAfDp>
- (b) RT @KeeganNYC: LOTS of NYPD waiting for protesters on the BK side of the Brooklyn Bridge #MillionsMarchNYC #ShutItDown #ICantBreathe <http://t.co/...>
- (c) RT @Zegota42: . @KeeganNYC Protesters on Brooklyn Bridge leaving Manhattan Skyline behind. #MillionsMarchNYC #ICantBreathe <http://t.co/UPvN...>

# Chapter 5

## Potential Applications of the EIIM Framework

### 5.1 Event Monitoring and Analysis

References related to real-life events are extremely abundant in social media. Right from natural disasters such as the ‘Haiti Earthquake’ [101] to international sporting events like the ‘Winter Olympics’ [102] to socio-political [4] and socio-economical [103] events that shook the world such as presidential elections [104], ‘Egyptian Revolution’ [105], and recessions were covered, analyzed, extrapolated and informed by social media. This prolific event-specific content in social media makes it a promising ground for performing event analytics. Platforms like Geofeedia<sup>1</sup>, TwitterStand<sup>2</sup>, Twitris<sup>3</sup>,Truthy<sup>4</sup>, and TweetTracker<sup>5</sup> have developed techniques to provide analytics related to different local and global real-life events.

Monitoring social media has become one of the essential activities of national security agencies for predicting potential threats and mass protests [106]. Social media is being used for tracking terrorism activities [107], collective actions [12], and countering cyber-attack threats<sup>6</sup>. One of the main components of each of these applications is tracking references related to the events. The proposed EIIM model could be an essential component of such systems. It would help in identifying, tracking and analyzing events and its related references in an organized manner over time.

---

<sup>1</sup><http://geofeedia.com/>

<sup>2</sup><http://twitterstand.umiacs.umd.edu/>

<sup>3</sup><http://twitris.knoesis.org/>

<sup>4</sup><http://truthy.indiana.edu/>

<sup>5</sup><http://tweettracker.fulton.asu.edu/>

<sup>6</sup><https://www.recordedfuture.com/>

## 5.2 Event Information Retrieval

Retrieving informative content related to real-life events shared in social media and presenting them in an organized way to the interested users has led to web based services like Seen<sup>7</sup>. It allows users to follow live updates of the events and also aids in witnessing and re-living the events at a later stage from the archives. Showing useful and interesting content to users by filtering out the pointless babbles from social media streams is an important component of such services. Additionally, such systems could get immensely benefitted by identification of event-specific informative hashtags, text units, users and URLs over time as the event proceeds. This would further enable efficient indexing of event-specific terms and hashtags that leads to high quality information, and effective processing of information. It would enhance the user experience, allowing better consumption and summarization of information related to the events, and positively impact triggering of event-specific recommendations. Thus, the proposed EIIM model in this thesis can act as the core component of information retrieval systems retrieving and organizing information related to real-life events from social media.

## 5.3 Opinion and Review Mining

Every day millions of people express their opinions in social media about products and companies they like and dislike. Their communications often include thoughts about good and bad experiences with the products and services. This provides a great opportunity for companies to understand its customers and to get unbiased valuable feedback from them about their product offerings without asking them to fill out time consuming outdated surveys. The EIIM framework when used for monitoring references of products/services from social media during product launch events could be useful in mining insightful and informative opinionated content. Combined with sentiment analysis, the invention could be a powerful tool for review analysis. One of the important contributions of the system could be to identify the sources having high chances of containing insightful information and filter them out for further processing. This would make a review mining system more efficient and increase its overall quality. Mining opinions related to entities related to an event could be used in many other contexts like political campaigns, socio-political studies, market behavior analysis, e-commerce applications, etc. Steps are being taken for adding this capability to the EIIM framework. On considering a mix of named entities and unigram opinionated words as text units in the *EventIdentityInfoGraph* we obtained some preliminary encouraging results. A glimpse

---

<sup>7</sup><http://seen.co>

of the results obtained for a basketball game "Miami Heats VS Cleveland Cavaliers", played on 25th December, 2014 is as follows:

Top 10 insightful and opinionated tweets for an hour related to the game

1. Good win for the Heat tonight against Cavs and Lebron. Great game for Wade and Deng. Just imagine if Bosh were healthy. #HeatvsCavs
2. Good work Dwayne Wade. Good work Miami Heat. LeBron is embarrassed. It's all over his face. #NBA #heatvscavs
3. Great game on Christmas Heat Showed up and spoiled Lebron Return to MIA! #Wade County #HeatvsCavs #NBAChristmas
4. Lebron leaves Miami high and dry and they cheer his return. Some even cheering cavs. Embarrassing bandwagon fan base. #heatv...
5. I totally understand LBJ move to Cleveland and like it. But if I'm a #Miami fan, I would boo LeBron like crazy today. #heatvscavs #CLEvsMIA
6. Stay classy #Miami. Good game vs. Lebron and; Cavs. #NBA #MIAvsCLE #HeatvsCavs #Heat #HeatNation
7. Loul Deng playing both ends of the floor. He's playing good D to LBJ #heatvscavs
8. Heat fans ; Cavs fans. Class vs no class. No burning a jersey in Miami #heatvscavs #HeatNation
9. WE FUCKING WON!!!!!! LETS GO HEAT #HEATgame #HeatNation #Heatvs-  
Cavs Wade with 31 points 5 assist 5 rebounds! Good shit MIAMI
10. Kevin Love is overrated. Big fish, small pond in MN and injury prone. #Heatvs-  
Cavs #NBAXmas

The above tweets point to the reactions of the viewers on the game as well as the players participating in the event.

## 5.4 Recommender Systems

The EIIM framework can be used for developing event related recommender systems. The ranked list of event identity information can be used for giving useful recommendations. For example following is a refined tweet recommendation for an event obtained

from a snapshot of the *EventIdentityInfoGraph* created for the event: “BlackLivesMatter”: Protest movement against the killing of Eric Garner.

#### Original Tweet:

- #BREAKING #NEWS — New York City Mayor Says, #BlackLivesMatter  
<http://t.co/qYvp8L8gDh> — #BLACK HCP520

#### Recommended Tweets:

- New York: What’s the plan? Where are the protests happening tonight? #EricGarner #BlackLivesMatter #MichaelBrown #ICantBreathe
- Brooklyn District Attorney to Convene Grand Jury in Case of #AkaiGurley NBC New York <http://t.co/mLiYPy39Pa> #BlackLivesMatter
- New York Today! #ShutItDown #economicshutdown #BlackLivesMatter #ICantBreathe #EricGarner #nojusticenoprofits <http://t.co/F0TrZtx2Y5>

Similarly an user can get other recommended users who are talking on the same topic. Hashtags and topics can also be recommended. It can further lead to clustering of similar content and discovery of communities around different topics related to the event. We wish to work on this in the future.

## 5.5 Event Management and Marketing

Social media is increasingly being used by event management practitioners while organizing conferences, seminars, music festivals, fashion shows, fundraisers and various other types of planned events. Tracking and producing useful and informative content before, during and after the events in social media from the perspective of event management has proved to be extremely beneficial <sup>8</sup>. Right from promoting the events, collecting RSVPs, creating communities around topics, announcing important information, getting real-time unbiased feedbacks, to marketing right content to the users creating buzz about the events, social media plays an important role. It also helps in building long term relationships with the communities of users interested in an event and track their related activities. In such a scenario the EIIM life cycle can constantly track and persistently store salient information related to events right from its inception. The *EventIdentityInfoGraph* can aid in identifying event-specific informative content and users producing

---

<sup>8</sup><http://oursocialtimes.com/using-social-media-to-make-your-event-a-dazzling-success-infographic/>

them, which could further lead to effective targeting of user communities, generating event summaries, mining opinions, broadcasting interesting information, among other things related to an event.

## 5.6 Social Media Data Integration

Organizations have increasingly started integrating the data available in social media with the enterprise data<sup>9</sup>. Social media data is most powerful when it is combined with daily transactional data and the master data to give a comprehensive view of customers, products and business conditions. Customers often openly talk about the products in social media and build communities around hashtags [? ] related to different topics. The EIIM framework could go a long way in collecting right information about the entities of concern maintained in the enterprise databases and integrate the collected information with the already existing ones. The entity resolution aspect would further help in managing the data quality issues related to data integration. In such conditions the EIIM model proposed could be used for integrating entity information from two distinct domains of enterprise system and social media in order to gain strategic intelligence related to business of an organization. This would further help an organization in marketing, corporate communications, public relations, customer support, product development, advertising, market research, product recommendations and gaining competitive intelligence.

---

<sup>9</sup><http://www.altimetergroup.com/research/reports/social-data-intelligence>

# **Chapter 6**

## **Conclusion and Future Work**

### **6.1 Conclusion**

### **6.2 Future Work**

#### **6.2.1 Summarizing Event Related Content**

#### **6.2.2 Identifying Insightful Opinionated Content Related to Events**

#### **6.2.3 Event Topic Modeling**

#### **6.2.4 Event-specific Recommendations**

#### **6.2.5 Distributed Processing of EventIdentityInfoGraph**

#### **6.2.6 Event Ontology for Social Media**

## **Appendix A**

### **Appendix Title Here**

Write your Appendix content here.

# Bibliography

- [1] Owen Phelan, Kevin McCarthy, and Barry Smyth. Using twitter to recommend real-time topical news. In *Proceedings of the third ACM conference on Recommender systems*, pages 385–388. ACM, 2009.
- [2] Brendan O’Connor, Ramnath Balasubramanyan, Bryan R Routledge, and Noah A Smith. From tweets to polls: Linking text sentiment to public opinion time series. *ICWSM*, 11:122–129, 2010.
- [3] Vivek Kumar Singh, Rakesh Adhikari, and Debanjan Mahata. A clustering and opinion mining approach to socio-political analysis of the blogosphere. In *Computational Intelligence and Computing Research (ICCIC), 2010 IEEE International Conference on*, pages 1–4. IEEE, 2010.
- [4] Vivek Kumar Singh, Debanjan Mahata, and Rakesh Adhikari. Mining the blogosphere from a socio-political perspective. In *Computer Information Systems and Industrial Management Applications (CISIM), 2010 International Conference on*, pages 365–370. IEEE, 2010.
- [5] Nitin Agarwal, Merlyna Lim, and Rolf T Wigand. Online collective action and the role of social media in mobilizing opinions: A case study on women’s right-to-drive campaigns in saudi arabia. In *Web 2.0 Technologies and Democratic Governance*, pages 99–123. Springer, 2012.
- [6] Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. Tweet analysis for real-time event detection and earthquake reporting system development. *Knowledge and Data Engineering, IEEE Transactions on*, 25(4):919–931, 2013.
- [7] Ana-Maria Popescu, Marco Pennacchiotti, and Deepa Paranjpe. Extracting events and event descriptions from twitter. In *Proceedings of the 20th international conference companion on World wide web*, pages 105–106. ACM, 2011.
- [8] Hemant Purohit and Amit P Sheth. Twitris v3: From citizen sensing to analysis, coordination and action. In *ICWSM*, 2013.

- [9] Nitin Agarwal and Debanjan Mahata. Grouping the similar among the disconnected bloggers. In *Social Media Mining and Social Network Analysis: Emerging Research: Emerging Research*, page 54, 2013.
- [10] Nitin Agarwal, Debanjan Mahata, and Huan Liu. Time-and event-driven modeling of blogger influence. In *Encyclopedia of Social Network Analysis and Mining*, pages 2154–2165. Springer, 2014.
- [11] Fatih Sen, Rolf T Wigand, Nitin Agarwal, Debanjan Mahata, and Halil Bisgin. Identifying focal patterns in social networks. In *CASoN*, pages 105–108, 2012.
- [12] Nitin Agarwal, Merlyna Lim, and Rolf T Wigand. *Online collective action: Dynamics of the crowd in social media*. Springer, 2014.
- [13] Nitin Agarwal, Merlyna Lim, and Rolf T Wigand. Raising and rising voices in social media. *Business & Information Systems Engineering*, 4(3):113–126, 2012.
- [14] James Allan, Ron Papka, and Victor Lavrenko. On-line new event detection and tracking. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 37–45. ACM, 1998.
- [15] Giridhar Kumaran and James Allan. Text classification and named entities for new event detection. In *Proceedings of the 27th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 297–304. ACM, 2004.
- [16] Yinle Zhou and John Talburt. Entity identity information management (eiim). In *International Conference on Information Quality (ICIQ-11), Adelaide, Australia*, pages 327–341, 2011.
- [17] Fabricio Benevenuto, Gabriel Magno, Tiago Rodrigues, and Virgilio Almeida. Detecting spammers on twitter. In *Collaboration, electronic messaging, anti-abuse and spam conference (CEAS)*, volume 6, page 12, 2010.
- [18] Sidharth Chhabra, Anupama Aggarwal, Fabricio Benevenuto, and Ponnurangam Kumaraguru. Phishy social: the phishing landscape through short urls. In *Proceedings of the 8th Annual Collaboration, Electronic messaging, Anti-Abuse and Spam Conference*, pages 92–101. ACM, 2011.
- [19] Chris Grier, Kurt Thomas, Vern Paxson, and Michael Zhang. @ spam: the underground on 140 characters or less. In *Proceedings of the 17th ACM conference on Computer and communications security*, pages 27–37. ACM, 2010.
- [20] Sarita Yardi, Daniel Romero, Grant Schoenebeck, et al. Detecting spam in a twitter network. *First Monday*, 15(1), 2009.

- [21] Emma Tonkin, Heather D Pfeiffer, and Greg Tourte. Twitter, information sharing and the london riots? *Bulletin of the American Society for Information Science and Technology*, 38(2):49–57, 2012.
- [22] Carlos Castillo, Marcelo Mendoza, and Barbara Poblete. Information credibility on twitter. In *Proceedings of the 20th international conference on World wide web*, pages 675–684. ACM, 2011.
- [23] Eugene Agichtein, Carlos Castillo, Debora Donato, Aristides Gionis, and Gilad Mishne. Finding high-quality content in social media. In *Proceedings of the 2008 International Conference on Web Search and Data Mining*, pages 183–194. ACM, 2008.
- [24] Manish Gupta, Peixiang Zhao, and Jiawei Han. Evaluating event credibility on twitter. In *SDM*, pages 153–164. SIAM, 2012.
- [25] Aditi Gupta and Ponnurangam Kumaraguru. Credibility ranking of tweets during high impact events. In *Proceedings of the 1st Workshop on Privacy and Security in Online Social Media*, page 2. ACM, 2012.
- [26] Jacob Ratkiewicz, Michael Conover, Mark Meiss, Bruno Gonçalves, Snehal Patil, Alessandro Flammini, and Filippo Menczer. Truthy: mapping the spread of astroturf in microblog streams. In *Proceedings of the 20th international conference companion on World wide web*, pages 249–252. ACM, 2011.
- [27] Firas Damak, Karen Pinel-Sauvagnat, Mohand Boughanem, and Guillaume Cabanac. Effectiveness of state-of-the-art features for microblog search. In *Proceedings of the 28th Annual ACM Symposium on Applied Computing*, pages 914–919. ACM, 2013.
- [28] Rinkesh Nagmoti, Ankur Teredesai, and Martine De Cock. Ranking approaches for microblog search. In *2010 IEEE/WIC/ACM International joint conference on Web Intelligence-Intelligent Agent Technology (WI-IAT)*, volume 1, pages 153–157. IEEE Computer Society, 2010.
- [29] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. 1999.
- [30] Jianshu Weng, Ee-Peng Lim, Jing Jiang, and Qi He. Twitterrank: finding topic-sensitive influential twitterers. In *Proceedings of the third ACM international conference on Web search and data mining*, pages 261–270. ACM, 2010.
- [31] D Tunkelang. A twitter analog to pagerank. *The Noisy Channel*, 2009.

- [32] V Hallberg, A Hjalmarsson, J Puigcerver, C Rydberg, and J Stjernberg. An adaptation of the pagerank algorithm to twitter world. 2012.
- [33] Yajuan Duan, Long Jiang, Tao Qin, Ming Zhou, and Heung-Yeung Shum. An empirical study on learning to rank of tweets. In *Proceedings of the 23rd International Conference on Computational Linguistics*, pages 295–303. Association for Computational Linguistics, 2010.
- [34] Richard McCreadie and Craig Macdonald. Relevance in microblogs: Enhancing tweet retrieval using hyperlinked documents. In *Proceedings of the 10th conference on open research areas in information retrieval*, pages 189–196. LE CENTRE DE HAUTES ETUDES INTERNATIONALES D'INFORMATIQUE DOCUMENTAIRE, 2013.
- [35] Jan Vosecky, Kenneth Wai-Ting Leung, and Wilfred Ng. Searching for quality microblog posts: Filtering and ranking based on content analysis and implicit links. In *Database Systems for Advanced Applications*, pages 397–413. Springer, 2012.
- [36] Hila Becker, Mor Naaman, and Luis Gravano. Selecting quality twitter content for events. *ICWSM*, 11, 2011.
- [37] David Inouye and Jugal K Kalita. Comparing twitter summarization algorithms for multiple post summaries. In *Privacy, Security, Risk and Trust (PASSAT) and 2011 IEEE Third International Conference on Social Computing (SocialCom), 2011 IEEE Third International Conference on*, pages 298–306. IEEE, 2011.
- [38] Jeffrey Nichols, Jalal Mahmud, and Clemens Drews. Summarizing sporting events using twitter. In *Proceedings of the 2012 ACM international conference on Intelligent User Interfaces*, pages 189–198. ACM, 2012.
- [39] Beaux Sharifi, M-A Hutton, and Jugal K Kalita. Experiments in microblog summarization. In *Social Computing (SocialCom), 2010 IEEE Second International Conference on*, pages 49–56. IEEE, 2010.
- [40] Günes Erkan and Dragomir R Radev. Lexrank: graph-based lexical centrality as salience in text summarization. *Journal of Artificial Intelligence Research*, pages 457–479, 2004.
- [41] Rada Mihalcea and Paul Tarau. Textrank: Bringing order into texts. Association for Computational Linguistics, 2004.
- [42] Ivan P Fellegi and Alan B Sunter. A theory for record linkage. *Journal of the American Statistical Association*, 64(328):1183–1210, 1969.

- [43] Howard B Newcombe, James M Kennedy, SJ Axford, and Allison P James. Automatic linkage of vital records computers can be used to extract” follow-up” statistics of families from files of routine records. *Science*, 130(3381):954–959, 1959.
- [44] Thomas N Herzog, Fritz J Scheuren, and William E Winkler. *Data quality and record linkage techniques*. Springer Science & Business Media, 2007.
- [45] Mauricio A Hernández and Salvatore J Stolfo. Real-world data is dirty: Data cleansing and the merge/purge problem. *Data mining and knowledge discovery*, 2(1):9–37, 1998.
- [46] Sunita Sarawagi and Anuradha Bhamidipaty. Interactive deduplication using active learning. In *Proceedings of the eighth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 269–278. ACM, 2002.
- [47] Rohit Ananthakrishna, Surajit Chaudhuri, and Venkatesh Ganti. Eliminating fuzzy duplicates in data warehouses. In *Proceedings of the 28th international conference on Very Large Data Bases*, pages 586–597. VLDB Endowment, 2002.
- [48] Y Richard Wang and Stuart E Madnick. The inter-database instance identification problem in integrating autonomous systems. In *Data Engineering, 1989. Proceedings. Fifth International Conference on*, pages 46–55. IEEE, 1989.
- [49] William W Cohen, Henry Kautz, and David McAllester. Hardening soft information sources. In *Proceedings of the sixth ACM SIGKDD international conference on Knowledge discovery and data mining*, pages 255–259. ACM, 2000.
- [50] Mikhail Bilenko, Raymond Mooney, William Cohen, Pradeep Ravikumar, and Stephen Fienberg. Adaptive name matching in information integration. *IEEE Intelligent Systems*, 18(5):16–23, 2003.
- [51] Ahmed K Elmagarmid, Panagiotis G Ipeirotis, and Vassilios S Verykios. Duplicate record detection: A survey. *Knowledge and Data Engineering, IEEE Transactions on*, 19(1):1–16, 2007.
- [52] Hector Garcia-Molina. Pair-wise entity resolution: overview and challenges. In *Proceedings of the 15th ACM international conference on Information and knowledge management*, pages 1–1. ACM, 2006.
- [53] John R Talburt. *Entity resolution and information quality*. Elsevier, 2011.
- [54] Omar Benjelloun, Hector Garcia-Molina, Hideki Kawai, Tait Elliott Larson, David Menestrina, Qi Su, Suttipong Thavisomboon, and Jennifer Widom. Generic

- entity resolution in the serf project. *IEEE Data Engineering Bulletin, June 2006 Issue*, 2006.
- [55] Omar Benjelloun, Hector Garcia-Molina, Heng Gong, Hideki Kawai, Tait Elliott Larson, David Menestrina, and Suttipong Thavisomboon. D-swoosh: A family of algorithms for generic, distributed entity resolution. In *Distributed Computing Systems, 2007. ICDCS'07. 27th International Conference on*, pages 37–37. IEEE, 2007.
  - [56] John Talburt, Richard Wang, Kimberly Hess, and Emily Kuo. An algebraic approach to data quality metrics for entity resolution over large datasets. *Information quality management: Theory and applications*, pages 1–22, 2007.
  - [57] Jenny Rose Finkel. Named entity recognition and the stanford ner software, 2007.
  - [58] Jason Baldridge. The opennlp project. *URL: <http://opennlp.apache.org/index.html>, (accessed 2 February 2012)*, 2005.
  - [59] Hamish Cunningham. Gate, a general architecture for text engineering. *Computers and the Humanities*, 36(2):223–254, 2002.
  - [60] Breck Baldwin and Bob Carpenter. Lingpipe. *Available from World Wide Web: <http://alias-i.com/lingpipe>*, 2003.
  - [61] Steven Bird. Nltk: the natural language toolkit. In *Proceedings of the COLING/ACL on Interactive presentation sessions*, pages 69–72. Association for Computational Linguistics, 2006.
  - [62] Jakub Piskorski and Roman Yangarber. Information extraction: Past, present and future. In *Multi-source, multilingual information extraction and summarization*, pages 23–49. Springer, 2013.
  - [63] Sunita Sarawagi. Information extraction. *Foundations and trends in databases*, 1(3):261–377, 2008.
  - [64] Wen Hua, Dat T Huynh, Saeid Hosseini, Jiaheng Lu, and Xiaofang Zhou. Information extraction from microblogs: A survey. *Int. J. Soft. and Informatics*, 6(4):495–522, 2012.
  - [65] David Nadeau and Satoshi Sekine. A survey of named entity recognition and classification. *Linguisticae Investigationes*, 30(1):3–26, 2007.
  - [66] Kurt Bollacker, Colin Evans, Praveen Paritosh, Tim Sturge, and Jamie Taylor. Freebase: a collaboratively created graph database for structuring human knowledge. In *Proceedings of the 2008 ACM SIGMOD international conference on Management of data*, pages 1247–1250. ACM, 2008.

- [67] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. *Dbpedia: A nucleus for a web of open data*. Springer, 2007.
- [68] Omkar Deshpande, Digvijay S Lamba, Michel Tourn, Sanjib Das, Sri Subramaniam, Anand Rajaraman, Venky Harinarayan, and AnHai Doan. Building, maintaining, and using knowledge bases: A report from the trenches. In *Proceedings of the 2013 ACM SIGMOD International Conference on Management of Data*, pages 1209–1220. ACM, 2013.
- [69] Abhishek Gattani, Digvijay S Lamba, Nikesh Garera, Mitul Tiwari, Xiaoyong Chai, Sanjib Das, Sri Subramaniam, Anand Rajaraman, Venky Harinarayan, and AnHai Doan. Entity extraction, linking, classification, and tagging for social media: a wikipedia-based approach. *Proceedings of the VLDB Endowment*, 6(11):1126–1137, 2013.
- [70] Mustafa Bilgic, Louis Licamele, Lise Getoor, and Ben Shneiderman. D-dupe: An interactive tool for entity resolution in social networks. In *Visual Analytics Science And Technology, 2006 IEEE Symposium On*, pages 43–50. IEEE, 2006.
- [71] Lars Kolb, Andreas Thor, and Erhard Rahm. Dedoop: efficient deduplication with hadoop. *Proceedings of the VLDB Endowment*, 5(12):1878–1881, 2012.
- [72] John R Talburt and Yinle Zhou. *Entity Information Life Cycle for Big Data: Master Data Management and Information Integration*. Morgan Kaufmann, 2015.
- [73] Paolo Bouquet and Stefano Bortoli. Entity-centric social profile integration. In *Proceedings of the International Workshop on Linking of User Profiles and Applications in the Social Semantic Web (LUPAS 2010)*, pages 52–57, 2010.
- [74] Stefano Bortoli, Heiko Stoermer, Paolo Bouquet, and Holger Wache. Foaf-o-matic—solving the identity problem in the foaf network. In *SWAP*, 2007.
- [75] Elie Raad, Richard Chbeir, and Albert Dipanda. User profile matching in social networks. In *Network-Based Information Systems (NBiS), 2010 13th International Conference on*, pages 297–304. IEEE, 2010.
- [76] Paolo Bouquet, Heiko Stoermer, Michele Mancioppi, and Daniel Giacomuzzi. Okkam: Towards a solution to the “identity crisis” on the semantic web. In *SWAP*, volume 201, 2006.
- [77] James Allan. *Topic detection and tracking: event-based information organization*, volume 12. Springer Science & Business Media, 2002.

- [78] Yiming Yang, Tom Pierce, and Jaime Carbonell. A study of retrospective and online event detection. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 28–36. ACM, 1998.
- [79] Kuo Zhang, Juan Zi, and Li Gang Wu. New event detection based on indexing-tree and named entity. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 215–222. ACM, 2007.
- [80] Vasileios Hatzivassiloglou and Elena Filatova. Domain-independent detection, extraction, and labeling of atomic events. 2003.
- [81] Vasileios Hatzivassiloglou, Luis Gravano, and Ankineedu Maganti. An investigation of linguistic features and clustering algorithms for topical document clustering. In *Proceedings of the 23rd annual international ACM SIGIR conference on Research and development in information retrieval*, pages 224–231. ACM, 2000.
- [82] Juha Makkonen, Helena Ahonen-Myka, and Marko Salmenkivi. Simple semantics in topic detection and tracking. *Information Retrieval*, 7(3-4):347–368, 2004.
- [83] David Ahn. The stages of event extraction. In *Proceedings of the Workshop on Annotating and Reasoning about Time and Events*, pages 1–8. Association for Computational Linguistics, 2006.
- [84] Heng Ji and Ralph Grishman. Refining event extraction through cross-document inference. In *ACL*, pages 254–262, 2008.
- [85] Elena Filatova, Vasileios Hatzivassiloglou, and Kathleen McKeown. Automatic creation of domain templates. In *Proceedings of the COLING/ACL on Main conference poster sessions*, pages 207–214. Association for Computational Linguistics, 2006.
- [86] Leon Derczynski, Diana Maynard, Niraj Aswani, and Kalina Bontcheva. Microblog-genre noise and impact on semantic annotation accuracy. In *Proceedings of the 24th ACM Conference on Hypertext and Social Media*, pages 21–30. ACM, 2013.
- [87] Alan Ritter, Sam Clark, Oren Etzioni, et al. Named entity recognition in tweets: an experimental study. In *Proceedings of the Conference on Empirical Methods in Natural Language Processing*, pages 1524–1534. Association for Computational Linguistics, 2011.

- [88] Jianshu Weng and Bu-Sung Lee. Event detection in twitter. *ICWSM*, 11:401–408, 2011.
- [89] Jagan Sankaranarayanan, Hanan Samet, Benjamin E Teitler, Michael D Lieberman, and Jon Sperling. Twitterstand: news in tweets. In *Proceedings of the 17th ACM SIGSPATIAL international conference on advances in geographic information systems*, pages 42–51. ACM, 2009.
- [90] Saša Petrović, Miles Osborne, and Victor Lavrenko. Streaming first story detection with application to twitter. In *Human Language Technologies: The 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics*, pages 181–189. Association for Computational Linguistics, 2010.
- [91] T. Rattenbury et al. Towards automatic extraction of event and place semantics from flickr tags. In *Proceedings of the 30th annual international ACM SIGIR conference on Research and development in information retrieval*, pages 103–110. ACM, 2007.
- [92] Ling Chen and Abhishek Roy. Event detection from flickr data through wavelet-based spatial analysis. In *Proceedings of the 18th ACM conference on Information and knowledge management*, pages 523–532. ACM, 2009.
- [93] Hila Becker, Mor Naaman, and Luis Gravano. Learning similarity metrics for event identification in social media. In *Proceedings of the third ACM international conference on Web search and data mining*, pages 291–300. ACM, 2010.
- [94] Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. Earthquake shakes twitter users: real-time event detection by social sensors. In *Proceedings of the 19th international conference on World wide web*, pages 851–860. ACM, 2010.
- [95] Sarita Yardi and Danah Boyd. Tweeting from the town square: Measuring geographic local networks. In *ICWSM*, 2010.
- [96] Edward Benson, Aria Haghghi, and Regina Barzilay. Event discovery in social media feeds. In *Proceedings of the 49th Annual Meeting of the Association for Computational Linguistics: Human Language Technologies-Volume 1*, pages 389–398. Association for Computational Linguistics, 2011.
- [97] Hila Becker, Dan Iter, Mor Naaman, and Luis Gravano. Identifying content for planned events across social media sites. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 533–542. ACM, 2012.
- [98] Hila Becker, Feiyang Chen, Dan Iter, Mor Naaman, and Luis Gravano. Automatic identification and presentation of twitter content for planned events. In *ICWSM*, 2011.

- [99] Furu Wei, Wenjie Li, Qin Lu, and Yanxiang He. Query-sensitive mutual reinforcement chain and its application in query-oriented multi-document summarization. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 283–290. ACM, 2008.
- [100] Jimmy Lin and Michael Schatz. Design patterns for efficient graph algorithms in mapreduce. In *Proceedings of the Eighth Workshop on Mining and Learning with Graphs*, pages 78–85. ACM, 2010.
- [101] Huiji Gao, Geoffrey Barbier, Rebecca Goolsby, and Daniel Zeng. Harnessing the crowdsourcing power of social media for disaster relief. Technical report, DTIC Document, 2011.
- [102] Shaun Walker. Russia to monitor 'all communications' at winter olympics in sochi. *The Guardian*, October, 6, 2013.
- [103] Johan Bollen, Alberto Pepe, and Huina Mao. Modeling public mood and emotion: Twitter sentiment and socio-economic phenomena. *arXiv preprint arXiv:0911.1583*, 2009.
- [104] Emily Metzgar and Albert Maruggi. Social media and the 2008 us presidential election. *Journal of New Communications Research*, 4(1):141–165, 2009.
- [105] Alok Choudhary, William Hendrix, Kathy Lee, Diana Palsetia, and Wei-Keng Liao. Social media evolution of the egyptian revolution. *Communications of the ACM*, 55(5):74–80, 2012.
- [106] Jeffrey Ghannam. Social media in the arab world: Leading up to the uprisings of 2011. *Center for International Media Assistance*, 3, 2011.
- [107] Onook Oh, Manish Agrawal, and H Raghav Rao. Information control and terrorism: Tracking the mumbai terrorist attack through twitter. *Information Systems Frontiers*, 13(1):33–43, 2011.