

A FRAMEWORK FOR COLLECTING, EXTRACTING
AND MANAGING EVENT IDENTITY INFORMATION
FROM TEXTUAL CONTENT IN SOCIAL MEDIA

A Dissertation Submitted
to the Graduate School
University of Arkansas at Little Rock

in partial fulfillment of requirements
for the degree of

DOCTOR OF PHILOSOPHY

in Integrated Computing

in the Department of Information Science
of the Donaghey College of Engineering and Information
Technology

May 2015

Debanjan Mahata
M.S in Computer Science, Banaras Hindu University, 2010
B.S in Computer Science, Banaras Hindu University, 2008

© Copyright by
Debanjan Mahata
2015

This dissertation, “A Framework for Collecting, Extracting and Managing Event Identity Information from Textual Content in Social Media”, by Debanjan Mahaata, is approved by:

Dissertation Advisor: _____

John R. Talburt
Professor of Information Science

Dissertation Committee: _____

Elizabeth Pierce
Associate Professor of Information Science

Russel Bruhn
Professor of Information Science

Ningning Wu
Professor of Information Science

Mathias Brochhausen
Assistant Professor of the Division of Biomedical
Informatics at UAMS

Program Coordinator: _____

Ningning Wu
Professor of Information Science

Graduate Dean: _____

Paula J. Casey
Professor of Law

Fair Use

This thesis is protected by the Copyright Laws of the United States (Public Law 94-553, revised in 1976). Consistent with fair use as defined in the Copyright Laws, brief quotations from this material are allowed with proper acknowledgment. Use of this material for financial gain without the author's express written permission is not allowed.

Duplication

I authorize the Head of Interlibrary Loan or the Head of Archives at the Ottenheimer Library at the University of Arkansas at Little Rock to arrange for duplication of this thesis for educational or scholarly purposes when so requested by a library user. The duplication will be at the user's expense.

Signature: _____

“A FRAMEWORK FOR COLLECTING, EXTRACTING AND MANAGING EVENT IDENTITY INFORMATION FROM TEXTUAL CONTENT IN SOCIAL MEDIA”, by Debanjan Mahata, May 2015.

Abstract

With the popularity of social media platforms such as Facebook, Twitter and Google Plus, there has been voluminous growth in the digital footprints of real-life events on the Internet. The user-generated colloquial and concise textual content related to different types of real-life events, available in these websites, acts as an extremely useful source for researchers and organizations for extracting valuable and insightful information. There has been significant improvement in natural language processing techniques for mining formal and long textual content commonly found in newspapers. It is still a challenging task to mine textual information from the social media channels producing terse, informal and noisy text with an unusual grammatical structure.

For a real-life event of interest it is necessary to detect and store informative event-specific signals from the noisy social media channels that allows to distinctly identify the event among all others, and characterizes it for extracting actionable insights. These event-specific cues also form 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, data journalism, cyber security, event management, among others. Thus, there is a need of a generic framework that can collect the textual content related to a real-life event, extract event-specific information from it and persistently maintain the information for tracking newly produced content as the event evolves, and provide updated event analytics.

The patent-pending work presented in this dissertation establishes the design and implementation of an extendable framework that enables collection, extraction and persistent management of 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 critical component of the EIIM life 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. The work is concluded by discussing about its novel contributions, practical applications in various other domains and envisaging future directions.

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

Ronald Coase, Economics, Nobel Prize Laureate

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 South Asian 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 (Associate 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 in social media, from 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

| | |
|---|-------------|
| ABSTRACT | iv |
| Acknowledgements | viii |
| Contents | x |
| List of Figures | xiv |
| List of Tables | xvi |
| | |
| 1 Dissertation Overview | 1 |
| 2 Social Media and Real-life Events | 13 |
| 2.1 Social Media | 13 |
| 2.1.1 Blogs | 13 |
| 2.1.2 Microblogs | 14 |
| 2.1.3 Media Sharing | 15 |
| 2.1.4 Social Bookmarking | 15 |
| 2.1.5 Social News | 16 |
| 2.1.6 Social Networking | 16 |
| 2.2 Characteristics of Social Media Websites | 18 |
| 2.3 Events in Social Media | 23 |
| 2.4 Background: Entity Identity Information Management (EIIM) in Master Data Management | 25 |
| 2.5 Problem of Event Identity Information Management (EIIM) in Social Media | 30 |

| | | |
|----------|---|-----------|
| 3 | Literature Review | 39 |
| 3.1 | Identifying High Quality Informative Content in Social Media | 39 |
| 3.2 | Entity Resolution | 44 |
| 3.3 | Event Identification in News Text | 48 |
| 3.4 | Event Identification in Social Media | 51 |
| 4 | Challenges in Mining Event Related Content from Social Media | 55 |
| 4.1 | Information Overload | 56 |
| 4.2 | Veracity of Sources | 58 |
| 4.3 | Multiple Data Sources with Variety of Content | 60 |
| 4.4 | Informal Text | 61 |
| 4.5 | Searching for Information in Long Tail | 63 |
| 4.6 | Sparse Link Structure | 67 |
| 4.7 | Sampling Bias | 67 |
| 4.8 | Lack of Evaluation Datasets | 68 |
| 5 | Event Identity Information Management (EIIM) Life Cycle for Social Media | 70 |
| 5.1 | Identity Integrity | 71 |
| 5.2 | Event Reference Collection | 73 |
| 5.3 | Event Reference Preparation | 74 |
| 5.3.1 | Parts-of-speech tagging | 76 |
| 5.3.2 | Special Character Detection | 77 |
| 5.3.3 | Data Cleansing | 77 |
| 5.3.4 | Duplicate Detection | 77 |
| 5.3.5 | Stop Word Detection and Elimination | 79 |
| 5.3.6 | Slang Word Detection | 79 |
| 5.3.7 | Feeling Word Detection | 80 |
| 5.3.8 | Tokenization | 80 |
| 5.3.9 | Stemming | 81 |
| 5.3.10 | Tweet Meta-data Extraction | 81 |
| 5.3.11 | Named Entity Extraction | 82 |
| 5.4 | Event Information Quality | 82 |
| 5.4.1 | Annotated Dataset | 83 |
| 5.4.2 | Feature Selection and Training | 85 |

| | | |
|----------|--|------------|
| 5.4.3 | Model Evaluation | 88 |
| 5.4.4 | Assignment of Generic Informativeness Score | 88 |
| 5.5 | Event Identity Information Capture | 89 |
| 5.5.1 | Content Analysis of Event Related Tweets | 90 |
| 5.5.2 | Event Identity Information Units | 97 |
| 5.5.3 | Extracting Event Identity Information Units | 99 |
| 5.6 | Event Identity Information Structure (EIIS) | 100 |
| 5.7 | Event Identity Information Processing | 101 |
| 5.7.1 | EventIdentityInfoGraph | 103 |
| 5.7.2 | EventIdentityInfoRank | 109 |
| 5.8 | Event Reference Resolution | 116 |
| 5.9 | Event Analytics | 119 |
| 6 | Evaluations | 127 |
| 6.1 | Evaluation Baselines | 127 |
| 6.2 | Evaluation Setup and Objectives | 133 |
| 6.2.1 | Tweet Annotation | 135 |
| 6.2.2 | Hashtags, Text Units and URL Annotations | 137 |
| 6.2.3 | User Annotations | 138 |
| 6.2.4 | NDCG@n and Precision@n | 138 |
| 7 | Potential Applications of the EIIM Framework | 144 |
| 7.1 | Event Monitoring and Analysis | 145 |
| 7.2 | Event Information Retrieval | 146 |
| 7.3 | Opinion and Review Mining | 147 |
| 7.4 | Event Management and Marketing | 148 |
| 7.5 | Social Media Data Integration | 149 |
| 8 | Conclusion and Future Work | 150 |
| 8.1 | Conclusion | 150 |
| 8.2 | Future Work | 150 |
| 8.2.1 | Summarizing Event Related Content | 150 |
| 8.2.2 | Identifying Insightful Opinionated Content Related to Events | 151 |
| 8.2.3 | Event-specific Recommendations | 153 |
| 8.2.4 | Distributed Processing of EventIdentityInfoGraph | 155 |

| | |
|---|------------|
| 8.2.5 Event Ontology for Social Media | 156 |
| A Appendix Title Here | 158 |
| Bibliography | 159 |

List of Figures

| | | |
|-----|--|----|
| 1.1 | Event Identity Information Management (EIIM) Life Cycle for user generated textual content in social media | 5 |
| 2.1 | EIIM components and their interactions as proposed in [1] | 26 |
| 2.2 | Entity Identity Integrity in EIIM process. | 27 |
| 2.3 | Misjudgments made by EIIM process. | 27 |
| 2.4 | Event identity information for real-life event, ‘Egyptian Revolution 2011’. | 30 |
| 2.5 | Relation between elements of Ξ , M and S in Event Identity Information Management from social media. | 33 |
| 4.1 | Top 10 Google search results for “Egyptian Revolution”, “Libyan Revolution”, and “Tunisian Revolution”, visualized using TouchGraph. | 64 |
| 4.2 | Short Head Vs Long Tail media sources. | 64 |
| 4.3 | Top 10 entities from mainstream media and blogs. . . | 65 |
| 5.1 | Event Identity Information Management Life Cycle for Textual Content in Social Media. | 71 |
| 5.2 | Identity Integrity component of the EIIM life cycle. . | 72 |
| 5.3 | Event Reference Collection component of the EIIM life cycle. | 73 |
| 5.4 | Event Reference Preparation component of the EIIM life cycle. | 75 |
| 5.5 | Event Information Quality component of the EIIM life cycle. | 82 |
| 5.6 | Event Identity Information Capture component of the EIIM life cycle. | 89 |
| 5.7 | Distribution of hashtags in event related tweets. . . | 93 |
| 5.8 | Distribution of tokens in event related tweets. . . . | 93 |

| | | |
|------|--|-----|
| 5.9 | Distribution of URLs in event related tweets. | 94 |
| 5.10 | Distribution of users in event related tweets. | 94 |
| 5.11 | Content characteristics of informative and non-informative tweets related to events. | 95 |
| 5.12 | Event Identity Information Structure component of the EIIM life cycle. | 102 |
| 5.13 | Event Identity Information Processing component of the EIIM life cycle. | 103 |
| 5.14 | Mutual Reinforcement Chains in Twitter for an event. | 106 |
| 5.15 | Event Reference Resolution component of the EIIM life cycle. | 117 |
| 5.16 | Event Analytics component of the EIIM life cycle. . . | 120 |
| 6.1 | Performance comparison of ranking techniques using NDCG scores. | 140 |
| 6.2 | Performance comparison of ranking techniques using NDCG scores. . | 141 |
| 6.3 | Performance comparison of ranking techniques using NDCG scores. | 141 |
| 6.4 | Performance comparison of ranking techniques using NDCG scores. | 142 |
| 6.5 | Performance comparison of ranking techniques using precision scores. | 142 |
| 6.6 | Performance comparison of ranking techniques using precision scores. | 143 |

List of Tables

| | | |
|-----|--|-----|
| 2.1 | Popular social media websites belonging to different categories. | 17 |
| 2.2 | Number of active users for the top five social media websites used by American adults. | 18 |
| 4.1 | Examples of different event related tweets. | 58 |
| 4.2 | Presence of different meta-data in popular social media websites. | 61 |
| 5.1 | Details of data collected for analyzing event related tweet content. | 75 |
| 5.2 | Evaluation measures for logistic regression model. . . | 88 |
| 6.1 | Avg IIC scores and total avg scores of annotations for Millions March NYC event. | 134 |
| 6.2 | Avg IIC scores and total avg scores of annotations for Sydney Siege event. | 135 |

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

Chapter 1

Dissertation Overview

Social media is a paradigm shift in the way people communicate with each other. It has grown from being just a medium to a global medium of communication between people. Different types of social media platforms provide multiple venues for people to share first-hand experiences and exchange information about real-life events. It has become an indispensable means for disseminating news and real-time information about current events, using websites such as Twitter, Facebook, Instagram, Flickr, Youtube, Google Plus and Vine. These applications allow users to post short textual messages accompanied by images and videos. At the same time users also share their detailed journalistic experiences in the form of diaries through blogging platforms such as Blogger, Wordpress and Medium. Studies have shown the importance of social media platforms as a news circulation service [2], and a source for gauging public interest and opinions [3–6]. Its efficacy as a real-time, citizen-journalistic source

of information has been recently harnessed in the detection, extraction and analysis of real-life events [7–9]. The activities of users producing content in social media has also been studied for gaining deep insights about how users form communities around topics related to real-life events [10–12] that lead to collective action [13, 14].

With the popularity of social media there has been proliferation of unstructured textual content on the Internet about different real-life events. Tracking social media for useful content such as live reporting of an event and recent updates can be harnessed to identify important nuggets of information. It can lead to identification and analysis of insightful user-generated information about named entities (people, place, organization, etc). Event summaries can be generated from the identified event-specific informative content. Actionable event-specific insights can be extracted such as, “what is happening”, “who is involved”, “where is it happening”, “when did it happen”, and so on. There are tremendous applications in the areas of real-life event analysis, data journalism, event management, opinion mining, online targeted marketing, cyber security, among others. Thus, there is a need of a generic framework that has the following capabilities:

- Collecting different types of textual content produced in social media related to an event.

- Extracting information that acts as an identity of the event used for characterizing it.
- Maintaining the extracted event identity information persistently for resolving constantly produced new content and discovering important event-specific information.

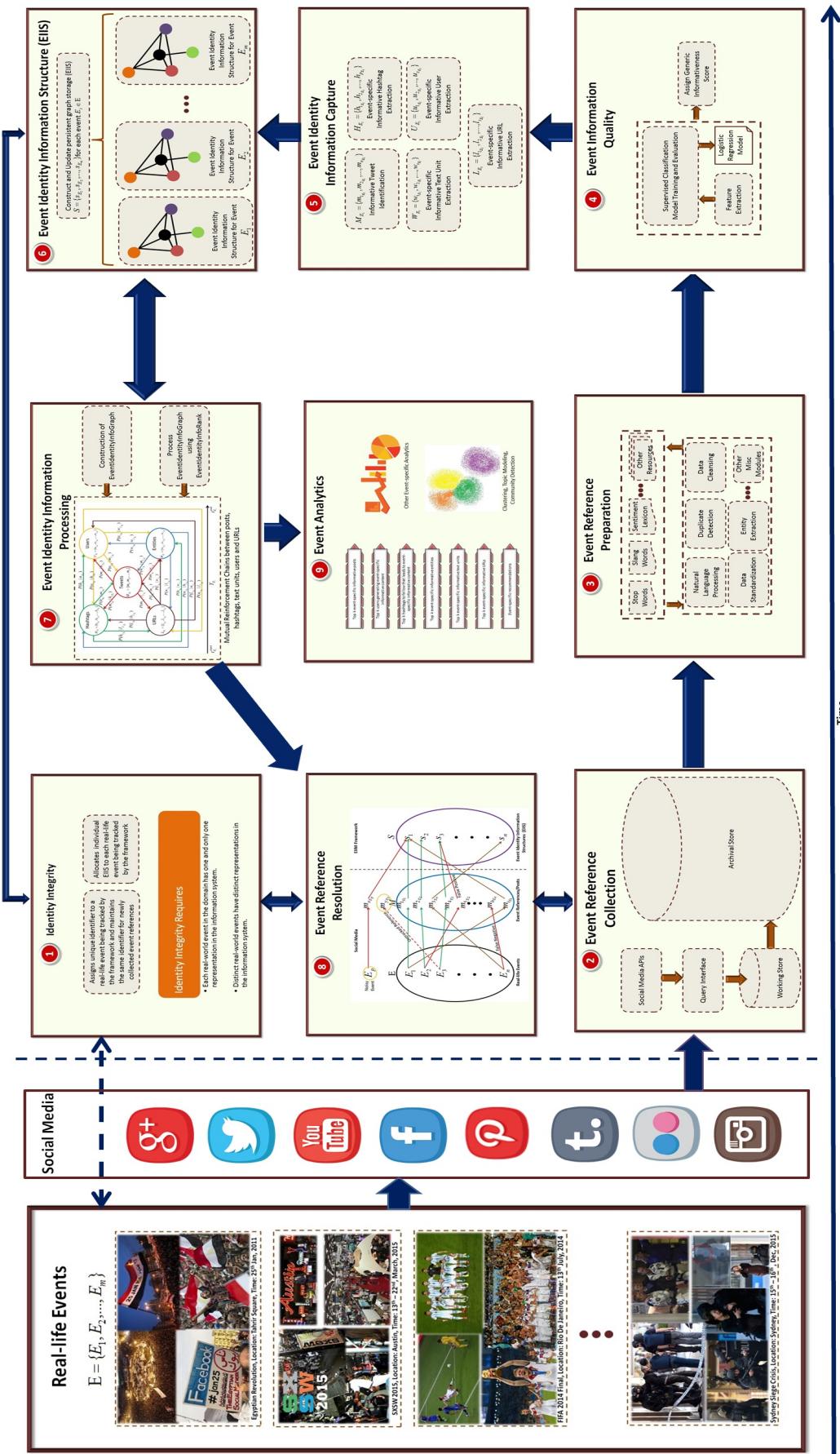
The problem of collecting and extracting event identity information from social media is similar to the task of event detection and tracking from newswires [15, 16]. However, in this dissertation, we add new components that creates persistent identity structures of an event and update it with new information over time. In order to make our task well-defined we focus on the problem of tracking 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 lacks proper grammatical structure. This makes it a challenging task for even the state-of-the-art natural language processing techniques to extract useful information and to 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 dissertation establishes the conceptual design and implementation of a framework capable of collecting, extracting and persistently managing event identity information from user-generated short 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) [1], 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 the challenges of identifying event-specific informative content is explained. This approach is a critical component of the framework. The dissertation further explores the applications of the research and concludes by discussing future directions of the work.

Some of the main contributions of the work are:

1. Extending the *Entity Identity Information Management* (EIIM) model [1] from the closed world domain of *Master Data Management* (MDM) to the open and unstructured domain of social media.
2. The design and implementation of an *Event Identity Information Management* framework that is capable of tracking and identifying event-specific information from short user-generated textual content in social media. Towards this objective a data

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



processing pipeline named *Event Identity Information Management Life Cycle* is developed, 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 references to 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 and ranking event-specific informative content produced in social media.
3. 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.
4. Analysis of informative and non-informative event related content from more than 3.8 million short textual social media messages.

5. A novel model based on the 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*.
6. A scalable graph processing iterative algorithm - *EventIdentityInfoRank*, 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.
7. Evaluation of the proposed technique against popularly used baseline techniques using large scale datasets.

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

Related Filed Patent

1. 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

1. **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

1. **Debanjan Mahata**, John R. Talburt and Vivek Kumar Singh; *A Framework for Collecting, Extracting and Managing Event Identity Information from Twitter.* 20th International Conference on Information Quality, M.I.T, Boston.
2. **Debanjan Mahata**, John R. Talburt and Vivek Kumar Singh; *From Chirps to Whistles : Discovering Event-specific Informative Content from Twitter.* Proceedings of the 7th Annual ACM Web Science Conference. ACM, 2015, Oxford, England.
3. **Debanjan Mahata**, John R. Talburt and Vivek Kumar Singh; *Identifying and Ranking of Event-specific Entity-centric Informative Content from Twitter.* 20th International Conference On Applications Of Natural Language To Information Systems (NLDB 2015), Passau, Germany. 17th – 19th June, 2015.

4. **Debanjan Mahata** and John R. Talburt; *A Framework for Collecting and Managing Entity Identity Information from Social Media.* 19th International Conference on Information Quality, Xi'An, China.
5. **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.
6. 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.
7. **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.
8. Nitin Agarwal, and **Debanjan Mahata**. *Grouping the Similar among the Disconnected Bloggers.* Social Media Mining and Social Network Analysis: Emerging Research (2013), 54.
9. **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.

10. **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.
11. 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.
12. 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.

The rest of the dissertation is organized as follows:

Chapter 2 gives an overview of the different social media websites and their common characteristics. It also gives the definition of events in social media as accepted by the presented work. Finally, it

introduces the problem of Event Identity Information Management from Social Media.

Chapter 3 reviews the existing literature related to the topic of the dissertation and highlights the challenges in applying previously available techniques to the domain of social media. It also discusses the similarities and dissimilarities of the work presented in this dissertation with the previous ones, and identifies the novel contributions that make it different from the state-of-the-art techniques.

Chapter 4 gives a detailed discussion of the challenges that are relevant to the problem of this dissertation. It discusses different scenarios in which these problems occur and point to their solutions as proposed in this dissertation.

Chapter 5 presents a detailed explanation of the *Event Identity Information Management Life Cycle*, that is proposed as a solution to the problem that is posed in this dissertation. It goes through all the components of the life cycle and gives a detailed explanation of the design choices, implementation and their working.

Chapter 6 discusses the baselines selected for evaluating the approach presented in this dissertation and compares its performance with them.

Chapter 7 highlights the potential real-life application of the *Event Identity Information Management* framework implemented in this dissertation.

Chapter 8 draws conclusions of the work presented in this dissertation and points to future directions of the work.

Chapter 2

Social Media and Real-life Events

2.1 Social Media

Social media is defined as a group of Internet-based applications that build on the ideological and technological foundations of Web 2.0, and that allow the creation and exchange of user-generated content [17]. These Internet-based applications broadly ranges from blogs, microblogs, media sharing webistes, social bookmarking websites, social news to social networking websites. Brief descriptions of most popular types of social media websites is given below.

2.1.1 Blogs

A blog can be defined as a website that displays, in a reverse chronological order, the entries by one or more individuals and usually has links to comments on specific postings. Blogs often provide opinions,

commentaries, or news on a particular subject, such as food, politics, or local news. Some of them also function as personal online diaries. Most of the time the entries of a blog is archived and is accessible at a later time. For the purpose of constant syndication, RSS or XML feeds for the blogs are made available. An individual entry in a blog is known as a blog post. A typical blog post can combine text, images and links to other blogs, web pages and other media related to its topic. The universe of all the blogs on the Internet is known as blogosphere [11].

2.1.2 Microblogs

Microblogs are similar to blogs, but a shorter version of it. Most of the microblogging websites pose limitations on the length of an individual post. Twitter¹, one of the most popular microblogging website has a limitation of 140 characters. This makes the textual posts in these platforms, extremely concise. Users often associate URLs that lead to external sources of information related to the posts. A post may also contain attached image or video. The microblogging services mostly focus on short updates that are pushed out to anyone subscribed to receive the updates. This is made possible by enabling the users to form directed networks of *friends* and *followers*. The

¹<http://twitter.com>

followers of an user are entitled to get all the updates posted by him.

Mostly these updates are public.

2.1.3 Media Sharing

Media sharing services allow its users to upload and share various multimedia content such as pictures and videos. Most services have additional social features such as profiles, commenting, etc. The most popular are Instagram², Pinterest³, YouTube⁴ and Flickr⁵. The media elements are often enriched with geographical and topical “tags” by the users who create them and the consumers who browse them. These tags acts as very useful meta data and allow automated programs to leverage them for efficient organization and retrieval of the videos and images that otherwise have very less textual content.

2.1.4 Social Bookmarking

These are the genre of social media services that allow its users to save, organize and manage links to various websites and resources on the Internet. Most allow to “tag” URLs for making them easy to search and share. The most popular are Delicious⁶ and StumbleUpon⁷. Some of the services like StumbleUpon also allow their users

²<http://instagram.com>

³<http://pinterest.com>

⁴<http://youtube.com>

⁵<http://flickr.com>

⁶<http://delicious.com>

⁷<http://stumbleupon.com>

to form friendship networks. These websites often provide different browsing experiences through interfaces that help the users to search for most recent tags, most popular ones, and so on.

2.1.5 Social News

Social news websites allow people to post various news items or links to articles that are external to the website, and allow its users to cast their “vote” on the items. The voting is the core social aspect as the items that get the most votes are displayed most prominently. This makes it an ideal crowdsourced news platform. It is up to the community of users to decide which news items gets seen by more people. Users can also “tag” the news stories and comment on them. The most popular are Digg⁸ and Reddit⁹.

2.1.6 Social Networking

Social networking websites are the ones that allow its users to connect with each other and form networks. The connections are generally non-directional and reciprocal. Two users who are connected to each other are considered as *friends*. Usually the users in these webistes have a profile that presents the personal information of the user as provided by him. The users have various ways to interact with other

⁸<http://digg.com>

⁹<http://reddit.com>

users, and also sometimes have the ability to set up groups. These social networks may be based on a certain theme such as interests, location, and profession. Facebook¹⁰ is the most popular personal social network and LinkedIn¹¹ is the most popular professional network.

Some of the other types of websites that can also be categorized as social media services are, social messaging services, collaboration tools, rating or review sites, personal broadcasting tools, virtual worlds, and group buying. Table 2.1, lists popular social media websites in different categories. Some of the websites may overlap and fall into multiple categories due to the broad range of services provided by them. For example, Facebook is not only a popular social networking website, but also a widely used social messaging service.

TABLE 2.1: Popular social media websites belonging to different categories.

| Category | Popular Social Media Websites |
|------------------------------|--|
| <i>Blogs</i> | Blogger, Medium, Wordpress, Squarespace |
| <i>Microblogs</i> | Twitter, Tumblr, Posterous |
| <i>Media Sharing</i> | Flickr, Instagram, YouTube, Vimeo, Dailymotion, Metacafe, Viddler, Pinterest |
| <i>Social Bookmarking</i> | Delicious, StumbleUpon, Scoop, Slashdot |
| <i>Social News</i> | Digg, Reddit, Newsvine, Propeller |
| <i>Social Networking</i> | Facebook, Google Plus, LinkedIn, Ello, CafeMom, Gather, Fitsugar |
| <i>Virtual Worlds</i> | Second Life, World of Warcraft, Farmville |
| <i>Group Buying</i> | Groupon, Living Social, CrowdSavings |
| <i>Personal Broadcasting</i> | Blog Talk radio, Ustream, Livestream |
| <i>Review/Rating</i> | Amazon ratings, Angie's List |
| <i>Collaboration Tools</i> | Wikipedia, WikiTravel, WikiBooks |
| <i>Social Messaging</i> | WhatsApp, Viber |

¹⁰<http://facebook.com>

¹¹<http://linkedin.com>

According to Pew Research Center, Facebook, LinkedIn, Pinterest, Instagram and Twitter are the top five most popular social media websites used by American adult Internet users¹². The number of active users world-wide for all the five social media sites is shown in Table 2.2. The numbers are obtained from the official pages of the respective websites.

TABLE 2.2: Number of active users for the top five social media websites used by American adults.

| Social Media Website | Number of Active Users |
|----------------------|------------------------|
| <i>Facebook</i> | 1.31 billion |
| <i>LinkedIn</i> | 347 million |
| <i>Twitter</i> | 289 million |
| <i>Instagram</i> | 100 million |
| Pinterest | 70 million |

2.2 Characteristics of Social Media Websites

All the above social media websites exhibit certain common characteristics that is also responsible for their wide usage and huge popularity. We revisit some of the characteristics as already suggested by Agarwal et al. [18] in the context of this dissertation.

1. **Accessibility:** Social media websites are freely available to whoever has an Internet connection. This makes these websites

¹²<http://www.pewinternet.org/2015/01/09/social-media-update-2014/>

easily accessible all over the world. One of the latest initiatives by Facebook and Google is to make social media accessible even to the most remote corner of the world through their Internet.org¹³ and “Loon for All”¹⁴ projects, respectively. With the popularity of hand held devices and increase in the Internet bandwidth, social media is accessible to anyone who has a smart phone and can use it. This is unlike the mainstream media or the print media, to which people subscribe and buy in the form of magazines, newspapers, journals, etc. Also, the mainstream media can be controlled by the government that may lead to propagation of biased information. For example, during the “Egyptian Revolution of 2011”, the mainstream media was biased, regulated by the government, and did not portray the true picture of the situation in Egypt. On the other hand, it was social media through which people discussed about the actual atrocities of the government and grouped together to incite the entire revolution [19].

2. **Permanence:** Social media websites show dynamic nature and the content can be altered any time. Users can easily edit the content shared by them. On the other hand the traditional print media and television media is not at all dynamic. Once

¹³<http://internet.org/>

¹⁴<http://www.google.com/loon/>

an article is printed in a magazine/newspaper, or a television show is recorded and broadcasted, it cannot be changed.

3. **Reach:** As already presented in Table 2.2, some of the popular social media websites have a reach of billions. Moreover, the ability of an individual user to simultaneously network with many other users makes this reach more effective than any other means of communication. These connections acts as networks of information flow, which helps in spreading any kind of information at a lightning speed [20]. Also it provides equal opportunity to everyone for reaching their intended audience, unlike the traditional mainstream media. This characteristic is now regularly used by politicians for launching election campaigns and reaching out to people in social media [21]. The marketers also leverage social media to a great extent. Event managers also take advantage of it, due to which social media has become an integral part of event management for getting connected with the event audience [22]. Social media is regularly used during planned events for making announcements, building and tracking audience, building focused communities, developing public relations [23], and targeted marketing [24]. Due to easy reachability in social media, a focused group of people can also get together very quickly and organize events such as flash mobs and protest movements [12].

4. **Recency:** The time lag at which communication can take place and information can flow is almost zero for the social media websites. Content is produced and communicated in real-time. Once this content is consumed, the users discuss about it instantaneously. Due to the reach and recency of social media it can make people aware of newsworthy events at a faster pace than traditional mainstream media [2, 25]. It might take hours, days, or sometimes months to present news or an event through mainstream media like newspapers, magazines and television. For example, the death of Bin Laden and the entire covert operation was reported in Twitter even before the US president made an official announcement in the mainstream media [26]. The users in social media were not only aware of the event but were also sharing and discussing it with great zeal. Another example is of the theater shootings in Colorado [27]. The shooting incident was reported, covered and analyzed in real-time, with traditional news media lagging behind social media by several minutes.

5. **Usability:** Social media sites are extremely easy to use and are user-friendly. An user does not require special training or skills to create content in a social media website. Whoever, can use a device connected to the Internet can share information in social media. Therefore, the operational cost of any social

media application is mostly negligible. This is not the case for mainstream media. In order to report an event one needs to be skilled and specially trained. Also, the printing and telecasting of any event has to go through many other processes and has to be finally approved by the editor. This makes traditional media unusable by common people. The use of social media at the time of Internet blackout during the “Egyptian Revolution of 2011” is a great example of its usability [28]. The Egyptian government had throttled the Internet connection and there was an Internet blackout for hours in order to stop the spreading of messages in social media, which was the major media of communication that led to the protest movements and finally to the revolution. In order to counterattack the government and to allow the Egyptians to use social media for giving latest updates, Google and Twitter launched a service that enabled them to leave a voicemail on a specific number. This voicemail was then posted in Twitter as a text message. Thus, people who didn’t have a smart phone or didn’t know how to use one could also post messages in social media.

2.3 Events in Social Media

Events have been discussed from different perspectives and have gained attention in various areas of academics. Right from philosophy [29] to psychology [30] events have been defined in various ways. In the context of this dissertation we define events as defined by Becker et al. [31].

Event: *An event is defined as a real-world occurrence (E_i) with an associated time period T_{E_i} ($t_{E_i}^{start}$ - $t_{E_i}^{end}$) and a time ordered stream of social media references $M_{E_i} = \{m_{1_{E_i}}, m_{2_{E_i}}, \dots, m_{n_{E_i}}, \dots, m_{z_{E_i}}\}$, of substantial volume, discussing about the event and posted in time T_{E_i} .*

Although, this definition of event was originally proposed for the domain of Twitter, it can be applicable to all types of social media platforms as all of them fulfill the basic requirements of the above definition. Becker et al. [31], further categorizes events in social media into two categories, planned and unplanned, depending on whether the context of an event (e.g., topic or hashtags on Twitter, time, location) is available or not.

Planned event: A planned event PE_i is an event in Twitter with preknown corresponding event context information consisting both

- topic or hashtags of the event.
- time, at which PE_i is planned to occur.

We only track planned events in this dissertation whose topics and hashtags are already known to us. For all the data collection efforts, an already known and event related popular hashtag is provided for bootstrapping the process of collecting event related messages. We consider all the messages in this stream as related to the event for which the query is placed and concentrate all our efforts to segregate high quality event-specific informative content from the non-informative ones.

Next, in contrast to the planned event, there is no information about an unplanned event. To characterize and define such events, we have to rely on other signals that could indicate their presence. Such events are often known as trending events and is defined as follows for the domain of Twitter.

Unplanned trending event: An unplanned trending event is an event in Twitter with one or more features (e.g., terms) of the corresponding Twitter messages exhibiting bursty patterns during the event's time period.

We don't consider the trending events for this dissertation as don't intend to solve the task of identifying a trending event from streams of social media messages.

2.4 Background: Entity Identity Information Management (EIIM) in Master Data Management

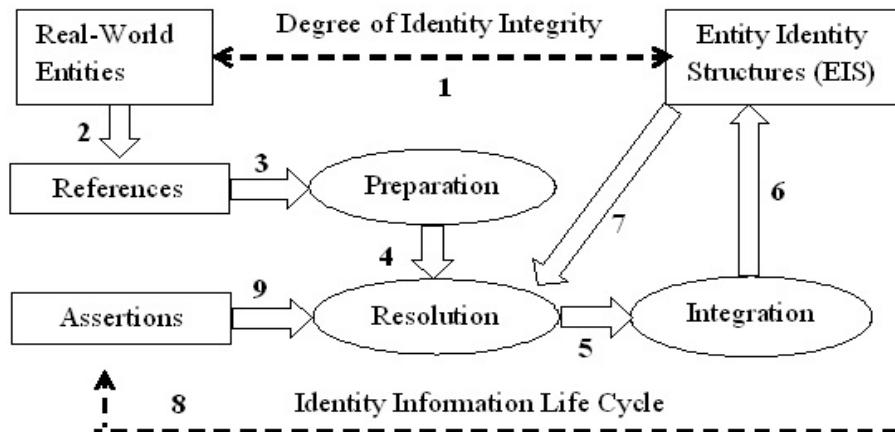
The idea of Entity Identity Information Management (EIIM) as defined by Zhou et al. [1], is the collection and management of identity information of real-world entities with the goal of sustaining *entity identity integrity*. *Entity identity integrity* is one of the basic tenets of data quality that applies to the representation of a given domain of real-world entities in an information system [32]. In order to maintain the property of *entity identity integrity* following conditions should be satisfied:

1. Each real-world entity in the domain has one and only one representation in the information system.
2. Distinct real-world entities have distinct representations in the information system.

Their model of EIIM was motivated by the problem of entity resolution in information systems, particularly in the domain of MDM

(Master Data Management). They define entity resolution as the process of determining whether two references to real-world objects in an information system are referring to the same object, or to different object [32]. The EIIM life cycle as proposed by them is an iterative process that combines entity resolution and data structures representing entity identity into specific operational configurations (EIIM configurations, as shown in Figure 2.1), that when executed in concert, work to maintain the entity identity integrity of master data over time. The EIIM framework is implemented by developing open source software known as OYSTER¹⁵.

FIGURE 2.1: EIIM components and their interactions as proposed in [1]



Some of the definitions as specified by the current EIIM model are:

- **Definition 1.** An *entity* (e_i) is defined as a real-life object that has a distinct identity.

¹⁵<http://sourceforge.net/projects/oysterer/>

- **Definition 2** *Entity Identity Information* is defined as a set of attributes of a given entity that distinctly characterizes it and allows that entity to be distinguished from all the other entities maintained by the framework.
- **Definition 3.** An *Entity Identity Information Structure (EIIS)* (s_i), is defined as a data structure that can persistently and efficiently store, retrieve, and manipulate entity identity information.

FIGURE 2.2: Entity Identity Integrity in EIIM process.

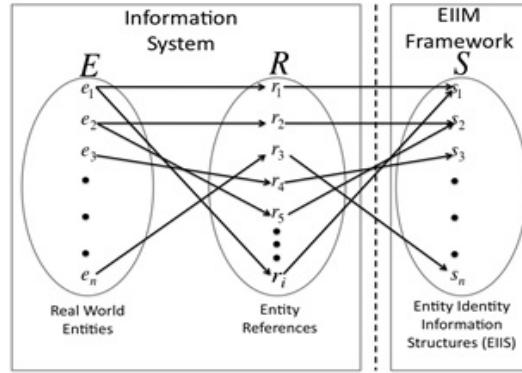
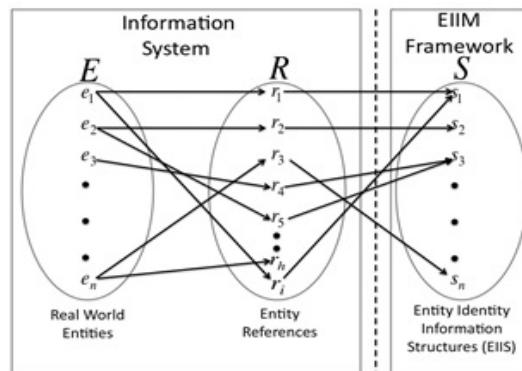


FIGURE 2.3: Misjudgments made by EIIM process.



Therefore, ideally in an information system, if $E = \{e_1, e_2, \dots, e_n\}$ represents a finite set of entities, $R = \{r_1, r_2, \dots, r_m\}$ represents a finite set of references to the entities, and $S = \{s_1, s_2, \dots, s_n\}$ represents a finite set of EIIS maintaining identity information of the entities then there should be one-to-one correspondence between the real-life entities ($\in E$) and the EIIS ($\in S$) representing their identity information. Also, the references ($\in R$) of a particular entity ($\in E$) should always map to one and only one EIIS ($\in S$) maintaining its identity information. This is shown in Figure 2.2. Such a situation ensures that the condition of entity identity integrity is satisfied by the information system. One of the main aims of EIIM is to satisfy the conditions of entity identity integrity along with persistently maintaining the entity identity information.

The current EIIM model deals with a closed environment of an information system where there is fixed number of entities along with fixed number of references to them. In an ideal situation the EIIM process should always satisfy the conditions of entity identity integrity as shown in Figure 2.2, and previously explained. However, in practice, all the references to an entity in the information system might not get mapped to the EIIS maintained for that particular entity due to misjudgments made by the automated processes as shown in Figure 2.3. This might result in *false negative* and *false positive* errors. A *false negative* error arises when the system fails to

map a reference of an entity to its corresponding EIIS. This is shown in Figure 2.3, where the system fails to map the reference r_h ($\in R$) of entity e_n ($\in E$) to an EIIS ($\in S$). A *false positive* error arises when the system maps two references of different entities to a single EIIS. This is shown in Figure 5, where the system wrongly maps reference r_5 ($\in R$) of entity e_2 ($\in E$), to the EIIS s_3 ($\in S$) being maintained for entity e_3 ($\in E$). Such a situation creates dissonance between the actual identity of the real-world entities being stored in the information system and their identities interpreted by the automated processes, resulting in low entity identity integrity of the system. Asserted resolutions are introduced in order to deal with such problems (shown in Figure 2.1.).

The EIIM processes and life cycle is a step ahead of the basic record linking process that identifies references to same entities for a given dataset. The goal of EIIM is to consistently label references to the same entity with the same identifier across different datasets processed at different times. Through the management of persistent entity identity structures, EIIM provides an added functionality for an entity resolution system to create and assign persistent entity identifiers that do not change from process to process. The current EIIM can also be thought of as forming a nexus between ER and MDM by adding an explicit longitudinal dimension to the management of identity information. The EIIM model proposed in the

presented research expands the current model into the unstructured domain of social media, bringing in new challenges and devising new techniques for solving them. The next section gives a detailed discussion and definition of the problem of extending the EIIM model to social media.

2.5 Problem of Event Identity Information Management (EIIM) in Social Media

FIGURE 2.4: Event identity information for real-life event, ‘Egyptian Revolution 2011’.



The definitions as given in the previous section also hold true for the work presented in this dissertation. The only differences are:

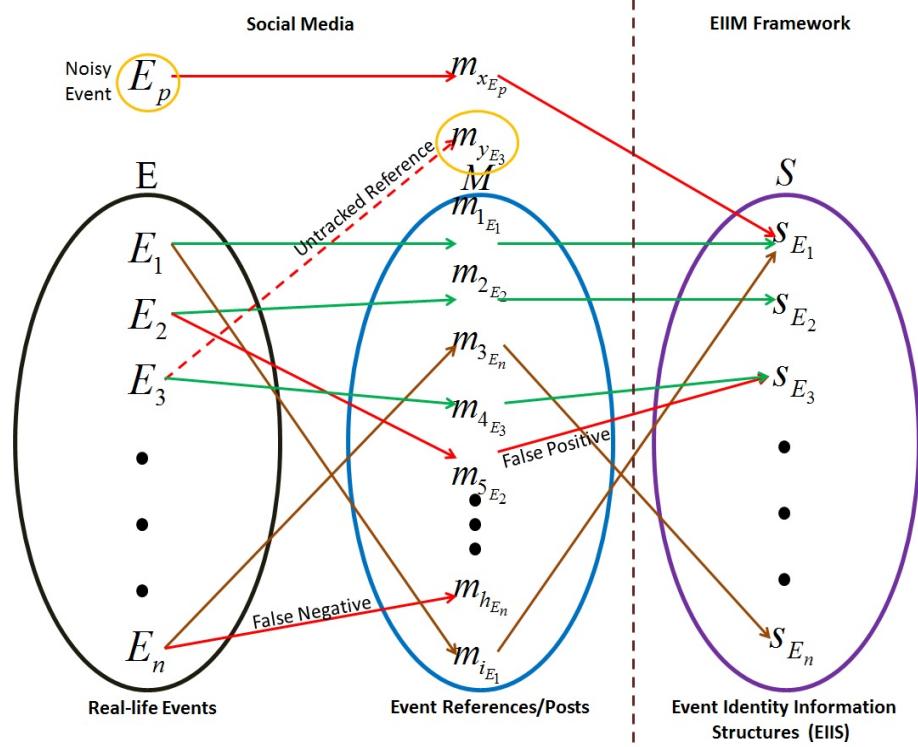
1. We consider a **real-life event** E_i as an individual **entity**, and assume that every **real-life event** also has a distinct identity. Only a pre-specified set of events $\Xi = \{E_1, E_2, \dots, E_n\}$ are considered, with their ordered stream of references $M = \{m_{1_{E_1}}, m_{2_{E_1}}, \dots, m_{i_{E_i}}, \dots, m_{n_{E_n}}\}$ collected from social media instead of a fixed set of references already present in a closed information system.
2. Instead of **Entity Identity Information**, we define **Event Identity Information** as a set of attributes that distinctly characterizes it and allows that event to be distinguished from all other events maintained by the framework. For example, names of the main active participants (people, organizations, etc), popular keywords, timespan of occurrence and place of occurrence could be identity information for a real-life event as shown in Figure 2.4.
3. Instead of **Entity Identity Information Structure (EIIS)**, we define **Event Identity Information Structures (EIIS)**, for real-life events as a data structure that can persistently and efficiently store, retrieve, and manipulate entity identity information. A set of EIIS, $S = \{s_{E_1}, s_{E_2}, \dots, s_{E_n}\}$ is maintained, corresponding to each event $E_i \in \Xi$

4. Lastly, this dissertation solves the problem of *Event Identity Information Management* in the open and unstructured domain of social media, that gives rise to new challenges. It added complexity to the original problem of EIIM, for which following steps are taken (explained in Chapter 5):

- A completely new life cycle of data processing components had to be introduced that is more suitable for processing unstructured references of events in social media.
- What consists of **Event Identity Information** became loosely defined due to the unstructured and highly dynamic nature of the environment. This led to development of new techniques that can identify identity information of an event w.r.t time, extract it from the references and store them in the **Event Identity Information Structures (EIIS)**.
- A complete new representation of **Event Identity Information Structures (EIIS)** is proposed, which is more suitable for managing the identity information of events in a dynamic and unstructured environment.

Events have been defined from various perspectives and in different contexts as already discussed in the previous section. In the context of the work presented in this dissertation we adopt a definition similar to [31].

FIGURE 2.5: Relation between elements of Ξ , M and S in Event Identity Information Management from social media.



Event: An event is defined as a real-world occurrence (E_i) with an associated time period T_{E_i} ($t_{E_i}^{start}$ - $t_{E_i}^{end}$) and a time ordered stream of social media references $M_{E_i} = \{m_{1_{E_i}}, m_{2_{E_i}}, \dots, m_{n_{E_i}}, \dots, m_{z_{E_i}}\}$, of substantial volume, discussing about the event and posted in time T_{E_i} .

Users in social media post about real-life events in huge volumes and with high velocity. This results in multiple footprints of an event in different social media channels. The footprints could be in the form of textual posts, images, videos or other types of multimedia documents. For example, following are three different tweets retrieved for the same event “Sydney Siege Crisis”:

1. RT @cnni: Hostage taker in Sydney cafe demands ISIS flag and call with Australian PM, Sky News reports. <http://t.co/a2vgrn30Xh> #sydneysiege
2. “@carlyeinfeld: Whats with the barking dogs in the background of the briefing on the #sydneysiege right now?!” Police dogs.
3. @Channel4News The least you could do is to pixelate the faces in the window #sydneysiege

We consider these footprints as references to the event. One of the main aims of the proposed framework would be to consolidate all the references of a particular event and map it to its corresponding EIIS maintained by the framework. However, due to the noisy nature of social media references an additional problem arise. We call it the problem of noisy reference.

- **Noisy Reference:** The three tweets in the example above are different in terms of information content related to the event. While the first one is very informative, the other two are surely related to the event, but not informative. The non-informative tweets might be interesting to the users having a personal conversation related to the event, but it has no value in conveying useful information about the event, that can act as an identity of the event. We consider these event related non-informative

references as noise and do not want to map these tweets to the EIIS of the corresponding event.

Also, there might be a mismatch between the way an informative event reference is posted by an user and the way the proposed framework interprets the same in an automated fashion. This would result in loss of data integrity and increase the chances of erroneous results. Such problems are prevalent in the vanilla flavor of the EIIM as discussed in the previous section, giving rise to *false negative* and *false positive* errors.

Since the Event Identity Information Management framework operates in an open environment of social media as discussed earlier, two new scenarios leading to erroneous conditions are observed:

- **Noisy Event** - The first situation occurs when a reference gets mapped to an existing EIIS in the framework, although it refers to an event, which is out of the pre-specified list of events that are being tracked. Such a scenario is shown in Figure 2.5, where the event $E_p \in \Xi$ generates the reference $m_{x_{E_i}} \in M_{E_i}$, yet the external reference gets mapped to s_{E_1} .

Spamming activity in social media channels can give rise to noisy events. For example, if we consider the following tweet:

- RT @BFDealz: <http://t.co/TSJAigrVJI> WHEELS SUPER TREASURE HUNT SUPERIZED HARLEY DAVIDSON FAT BOY LONG CARD 2014 #cpac2014 #sxsw

The above tweet might have been posted in the context of the event SXSW 2014 in order to market a Harley Davidson bike, but it also appears in the set of references for an unrelated event CPAC 2014. This is probably due to the fact that both the events took place parallelly and the user posting this reference wanted to make it more visible by using the hashtags for both the trending events. This is often the case of spamming activities as discussed in Chapter 4. The tweet might not be related to both the events, yet it can get tracked by the framework.

- **Untracked reference** - The second situation occurs when there is a reference, which the framework is unable to track and map it to a pre-specified list of events although it refers to it. Such a scenario is shown in Figure 2.5, where the framework should have tracked the reference $m_{y_{E_3}} \in M_{E_3}$ and associate it with event $E_3 \in \Xi$, yet it is unable to do so and lose track of the reference in the process. This can happen due to the sampling bias as discussed in Chapter 4, which is one of the main challenges in collecting content from social media. The other reason could be errors in data collection process.

Therefore, the main broader problem that this dissertation solves is stated below.

Problem: *Given a pre-specified finite set of real-life events, $\Xi = \{E_1, E_2, \dots, E_n\}$ generating a finite ordered stream of references $M = \{m_{1_{E_1}}, m_{2_{E_1}}, \dots, m_{i_{E_i}}, \dots, m_{n_{E_n}}\}$ in social media, and a finite set of EIIS structures $S = \{s_{E_1}, s_{E_2}, \dots, s_{E_n}\}$ corresponding to each real-life event ($E_i \in \Xi$), the problem is to resolve references of an event (E_i), and to persistently extract, store and manage identity information of the event in its corresponding EIIS s_{E_i} .*

Towards the above objective, the dissertation seeks answers to the following questions:

- How to create the set of EIIS (S) for the pre-specified set of events of interest ?
- How to collect reference to the events containing identity information of the event from different social media platforms ?
- How to extract event identity information from the social media references ?
- How to integrate the event identity information from the unstructured references into the appropriate EIIS ?
- How to determine that new references from social media relates to an EIIS ?

- How to continuously integrate new information into the existing EIIS structures ?

Chapter 5 presents the design and implementation of the entire Event Identity Information Management framework that explains the strategies and novel techniques contributed by this dissertation in order to find answers to the above questions. A detailed survey of the existing literature related to the problem is presented in the next Chapter.

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 [33–36]. 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 [37]. In

such cases, it becomes necessary to assess the credibility and trustworthiness of the information posted. It was showed by Castillo et al. [38] 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 [39] 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. [40]. 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 [41]. 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. Truthy¹, 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

¹<http://Truthy.indiana.edu/>

to an event to detect political smears, astroturfing, misinformation, and other forms of social pollution [42].

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 perspectives². 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 recency³. A study of different state-of-the-art features and approaches commonly used for ranking tweets has been documented by [43, 44]. Seen⁴ 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

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

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

⁴ <http://seen.co>

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* [45] has been proposed by the scientific community for ranking tweets and users in Twitter [46–48]. Tweet-Rank [48] 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 [49–51]. 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 [52]. 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.

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 [53]. Summarization of tweets for sporting events was performed by [54] using the phrase graph algorithm [55]. The popularly used graph-based summarization algorithms are *LexRank* [56] and *TextRank* [57]. 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 5). 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 [58–60]. In the database community, the problem is defined as merge-purge [61], data de-duplication [62, 63], and instance identification [64]. In the Artificial intelligence community, this problem is described as database hardening [65], and name matching [66]. The names co-reference resolution, identity uncertainty, and duplicate detection are also commonly used to refer to the same task [67]. 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 [68]. 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 [32] 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 5) 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 [58], which uses a decision-theoretic approach establishing the validity of principles first used in practice by Newcombe [59]. This was followed by the Stanford Entity Resolution Framework (SERF) developed at the Stanford InfoLab [69]. 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 [70], was also proposed that can be implemented in a big data environment. More recently the Talburt-Wang Algebraic Model of ER has been proposed [71] 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 [72], OpenNLP

[73], GATE [74], LingPipe [75] and NLTK [76]. 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 [77, 78]. A study of various efforts in extracting information from micro-blogs could be found in [79] and a survey of named entity recognition and classification could be found in [80]. Efforts have been made by the industry in building crowd sourced knowledge bases like freebase [81] and dbpedia [82] 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 [83, 84]. The rise of online social networks, has also motivated new research into the ER Activity 5, entity relationship analysis [85]. With the rise of big data, the modern trend is to perform entity resolution process in humongous volumes of data and scale it horizontally [86, 87]. 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 [1], 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 [88–90]. 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 [91]. 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 [92] 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 [15, 16]. While others explored the detection of past events from archived news documents [93].

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. [94] 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 [95] 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. [96] 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. [97] 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) [98, 99]. 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 [98] 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 [99] 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. [100] 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 [101]. For example, named entity recognition methods typically show 85-90% accuracy on longer texts, but 30-50% on tweets [102]. 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 [103] 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. [104] 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. [105] used locality-sensitive hashing to detect the first tweet associated with an event in a stream of Twitter messages. Rattenbury et al. [106] analyzed the temporal usage distribution of tags to identify tags that correspond to events. Chen and Roy [107] 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. [108] 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 [104] and a method for identifying top events suggested by [105] 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 [109, 110]. Sakaki et al. [109] 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. [111]

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. [112] 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 [113] as well as different social media websites [112].

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

Challenges in Mining Event Related Content from Social Media

The characteristics of social media websites as discussed in Section 2.2 of Chapter 2, makes the domain extremely unstructured and uncontrolled. This gives rise to different challenges in mining information from all types of social media platforms. Some of the major challenges applicable to the context of this dissertation are discussed below. The solutions to these challenges as proposed in this dissertation is also referred while explaining them. Social media posts from the datasets collected for the experiments in Chapter 5 (Section 5.2) are presented as examples while discussing the problems.

4.1 Information Overload

A daily average of 58 million tweets is posted in Twitter¹. On an average 60 million photos are shared in Instagram daily². Facebook stores 300 petabytes of data related to its users from all over the world³. These are some compelling statistics that makes social media not only rich in volume of data, but also variety, and the velocity at which data is being generated. Due to the great pace at which data is produced in social media, the search engines and content filtering algorithms often face the problem of information overload [114]. They suffer from the dilemma of assessing the accuracy and quality of information content in the sources being produced over their freshness. Thus, collecting different types of references of events from social media, assessing their quality, resolving and extracting identity information of the events poses great challenges in such a situation.

For example, 284 million monthly users of Twitter posting 500 million tweets per day produces a variety of content⁴. A significant proportion of it are related to different real-life events (e.g., football matches, conferences, music shows, etc). Majority of this content are personal updates (e.g. *Thanks for the memories Sochi! I've had the*

¹<http://www.statisticbrain.com/twitter-statistics/>

²<http://instagram.com/press/>

³<http://expandedramblings.com/index.php/by-the-numbers-17-amazing-facebook-stats/>

⁴<http://about.twitter.com/company>

*time of my life #Sochi2014 #sochiselfie <http://t.co/DqkLEaAMpo>), pointless babbles (e.g. *Ted Cruz is a dangerous man. Crazy and gaining support. Megalomaniac leaders are bad, mkay.* #CPAC #politics #joke) and spams (e.g *New post: Sochi Was For Suckers - Laugh Studios/* <http://t.co/cWQJCBp3Ow> #lol #funny #rofl #funnypic #wtf).*

Personal views and conversations might be of interest to a specific group of people. However, they are meaningless and provides no information to the general audience. On the other hand there are tweets that presents newsworthy content, recent updates and real-time coverage of on-going events (e.g. *In #Sochi, the Dutch are dominating the overall Olympic medal count* <http://t.co/jMR1WUqEK4> (*Reuters*) <http://t.co/dAfDhEgTGA>). These tweets provide event-specific informative content and are more useful for the users interested to know about the event. In this dissertation, we call them as event-specific informative tweets. One of the main problems due to information overload is to identify these tweets among millions of tweets being produced during the event. Table 4.1 presents some examples of different types of tweets shared during real-life events.

Techniques are developed in this dissertation that overcomes this challenge. In Chapter 5 (Section 5.4), a generic supervised classifier for Twitter is developed in order to identify event references

in real-time that has greater chances of containing high quality information. This results in segregation of informative references from the non-informative ones, and identification of the informative references for further processing. The *EventIdentityInfoRank* algorithm implemented in Chapter 5 (Section 5.7.2), processes the content in the filtered tweets and results in a ranked list of tweets that has high quality event-specific informative content.

TABLE 4.1: Examples of different event related tweets.

| |
|---|
| Ted Cruz is a dangerous man. Crazy and gaining support. Megalomaniac leaders are bad, mkay. #CPAC #politics #joke [personal/uninformative] Event: ‘CPAC 2014’ |
| Thanks for the memories Sochi! I’ve had the time of my life #Sochi2014 #sochiselfie http://t.co/DqkLEaAMpo. [personal/uninformative] Event: ‘Sochi Games’ |
| #SXSW14 #SXSW #sxswinteractive #CPAC2014 #CPAC #CPACPick-upLines #CPACPanels Be squared away perky TOP TWEETED of http://t.co/h0igdOVNW0. [spam/uninformative] Event: ‘CPAC 2014’ |
| In #Sochi, the Dutch are dominating the overall Olympic medal count http://t.co/jMR1WUqEK4 (Reuters) http://t.co/dAfDhEgTGA. [event-specific informative] Event: ‘Sochi Games’ |
| New post: Sochi Was For Suckers - Laugh Studios/ http://t.co/cWQJCBp3Ow #lol #funny #rofl #funnypic #fail #wtf. [spam/uninformative] Event: ‘Sochi Games’ |
| It’s tedcruz vs. SenJohnMcCain in a #CPAC spat. What did they say? Find out on #AC360 8p on CNN. [event-specific informative] Event: ‘CPAC 2014’ |

4.2 Veracity of Sources

Judging the accuracy of the information and detecting relevant, event-specific informative content from social media constitutes another challenging situation due to the malevolent practices of spam users. For trending topics, the search engines have started showing

real-time feeds from social media websites in their search results. This has attracted spammers who post trending hashtags or keywords along with their spam content in order to attract people to their websites offering products or services [33]. For example, the tweet (*RT @BFDealz: http://t.co/TSJAigrVJI WHEELS SUPER TREASURE HUNT SUPERIZED HARLEY DAVIDSON FAT BOY LONG CARD 2014 #cpac2014 #sxsw*) was posted during the parallel occurrence of CPAC 2014 (a political conference) and SXSW 2014, but has nothing to do with the events. Instead it leads to a deal related to Harley Davidson bike promoted using popular event related hashtags #cpac2014 and #sxsw.

An alarming 355% growth of spam in social media has been reported in 2013⁵. Social media has also been instrumental in spreading misinformation and rumors. Spread of misinformation not only results in pandemonium among the users⁶ but also result in extraction of completely wrong information about events. The users in the social media websites also develop nepotistic relationships in order to get higher scores in the ranking techniques with malicious intentions [115]. This also helps them to spread spam and other malicious content. Such behavior can also lead to cyber attacks.

Some examples of spam tweets in the collected dataset are shown

⁵<http://www.likeable.com/blog/2013/11/10-surprising-social-media-statistics/>

⁶<http://www.theguardian.com/uk/interactive/2011/dec/07/london-riots-twitter>

in Table 4.1. Most of the existing techniques face problems in combating spam. *EventIdentityInfoGraph* (Section 5.7.1) is presented and explained in this dissertation that defines mutually reinforcing chains in Twitter for identifying the most event-specific informative references and filters out the spam tweets after ranking its nodes using *EventIdentityInfoRank* (Section 5.7.2). The algorithm can also identify users who are producing event-specific informative content and URLs (images, videos, news articles) that are extremely relevant to the event.

4.3 Multiple Data Sources with Variety of Content

The APIs (Application Programming Interfaces) of the different social media websites returns data in different formats (JSON, XML) using different web standards (REST, HTTPS). Moreover, the information obtained from a social media website is dependent upon the type of content it produces. A video sharing website might return an entirely different set of information from a blogging website. Thus, integrating the data obtained from various social media platforms for the purpose of tracking and extraction of event related information is one of the many challenges.

Although, type and format of the data returned by the social media services vary, yet most of them have certain common meta-data associated with messages obtained in response to the API requests. These are hashtags, text units extracted from the messages/descriptions, the messages itself, the users posting them and the URLs that lead to external sources of information. Some of the websites where hashtags are not present, the associated tags can be used instead. Table 4.2 shows the presence of these meta-data in the most popular social media platforms.

TABLE 4.2: Presence of different meta-data in popular social media websites.

| Social Media/ Website | Hashtags | Users | Text Units | Messages/ Description | URLs |
|--------------------------|----------|-------|------------|--------------------------|------|
| <i>Facebook</i> | Yes | Yes | Yes | Yes | Yes |
| <i>LinkedIn</i> | Yes | Yes | Yes | Yes | Yes |
| <i>Pinterest</i> | Yes | Yes | Yes | Yes | Yes |
| <i>Instagram</i> | Yes | Yes | Yes | Yes | Yes |
| <i>Twitter</i> | Yes | Yes | Yes | Yes | Yes |
| <i>Flickr</i> | Yes | Yes | Yes | Yes | Yes |
| <i>Google Plus</i> | Yes | Yes | Yes | Yes | Yes |
| <i>YouTube</i> | Yes | Yes | Yes | Yes | Yes |

The developed techniques that are explained in this dissertation rely only upon the above meta-data. This makes the techniques generic and useful for most of the social media channels.

4.4 Informal Text

Unlike sources of news media and edited documents on the web, the textual content of the social media references are highly colloquial

and pose great difficulties in extracting information. One of the most important sources of information about events, prevalent in the domain of social media are the micro-blogging platforms. Micro blogs pose additional challenges due to their brevity, noisiness, idiosyncratic language, unusual structure and ambiguous representation of discourse [116]. Variation in language, less grammatical structure of sentences, unconventional uses of capitalization, frequent use of emoticons, and abbreviations have to be dealt by any system processing social media content. Moreover, various signals of communications embedded in the text in the form of hash-tags (eg. #sochi), retweets (RT) and user mentions (@) should be understood by the system in order to extract the contextual information hidden in the text. Intentional misspellings sometimes demonstrate examples of intonation in written text [117]. For instance, expressions like, ‘this is so cooool’, emphasizes stress on the emotions and conveys more information that should be captured. 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 [101]. For example, named entity recognition methods typically show 85-90% accuracy on longer texts, but 30-50% on tweets [102]. Status messages in social networking websites, content in question answering websites, reviews, and discussions in blogs, and forums exhibit similar nature and present similar challenges to information

extraction and text mining procedures.

In order to solve some of these problems additional resources are compiled in the dissertation. Some of these resources are list of slang words, list of acronyms and list of stop words commonly used in short social media text (refer Chapter 5). These resources are used to aid the natural language processing techniques, resulting in their better performance. Experiments were also performed using the state-of-the-art entity extraction service, AlchemyAPI⁷ (Chapter 6). The final results obtained using the data processing pipeline developed in this dissertation gave better results than AlchemyAPI for short textual content in social media.

4.5 Searching for Information in Long Tail

Due to the power law distribution of the Internet [118], and the present search engine technology, the ‘Short Head’ is generally dominated by the mainstream media websites. As illustrated in Figure 4.1 the top 10 search results for “Egyptian Revolution”, “Libyan Revolution”, and “Tunisian Revolution” by Google, visualized using Touchgraph⁸, mostly retrieved mainstream media references. Consequently, the social media sites get buried in the Long Tail [119] as

⁷<http://alchemyapi.com>

⁸<http://touchgraph.com>

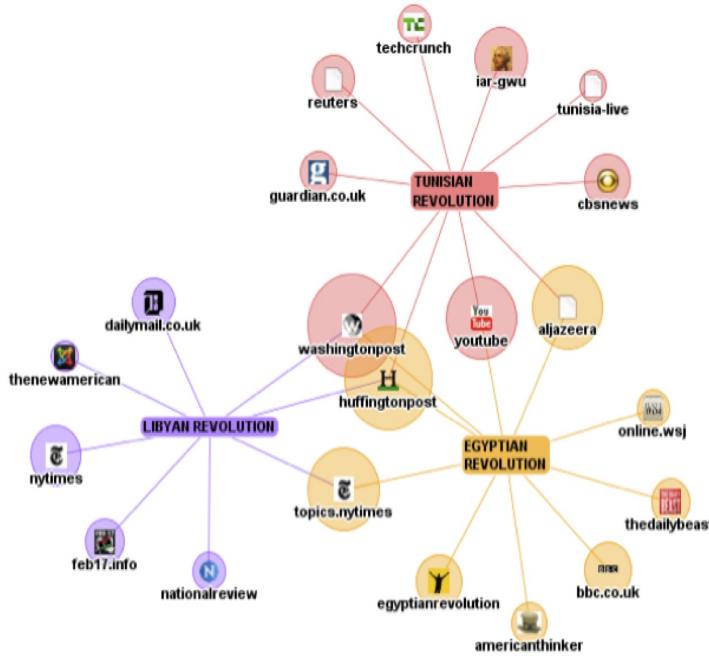


FIGURE 4.1: Top 10 Google search results for “Egyptian Revolution”, “Libyan Revolution”, and “Tunisian Revolution”, visualized using TouchGraph.

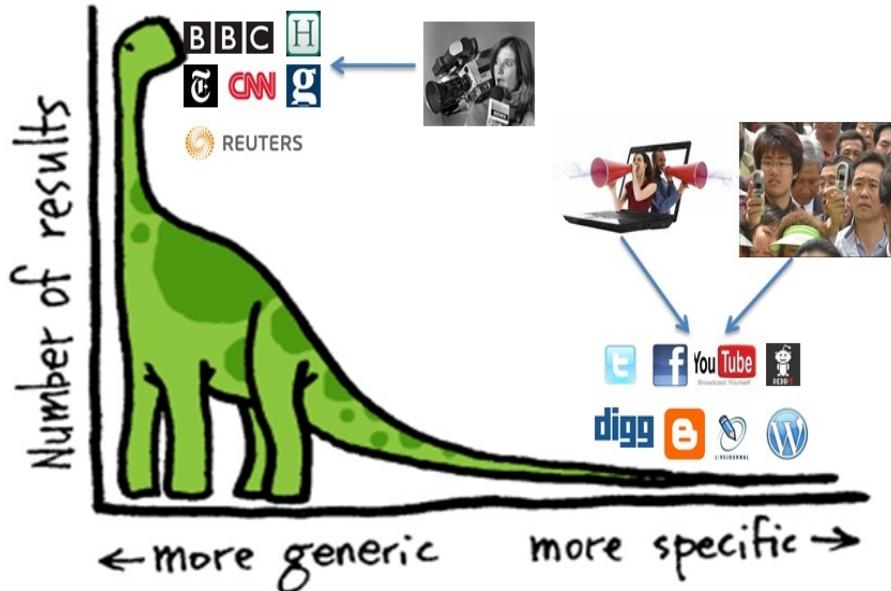


FIGURE 4.2: Short Head Vs Long Tail media sources.

shown in Figure 4.2. However, references from the social media channels, act as hubs of relevant information about real-life events [120]. On the other hand, the mainstream media sources often gloss over

the intricate details while covering a real-life event. They are often biased, regulated by the government, and may not portray the true picture of an event [19]. While, social media references often contain unbiased, uninhibited, and unedited opinions from people. Political blogs have been accepted as more credible sources of information over mainstream media references by the weblog users [121]. Thus the references, which are obtained from social media could potentially provide a rather ‘closer’ or an “on-ground” view of the events with novel information. The “on-ground” information gleaned from the social media affords opportunities to study various online social phenomenon from methodological and theoretical perspectives including, social movements, crowdsourcing, citizen journalism, collective behavior, collective action [6, 14, 122], and more.

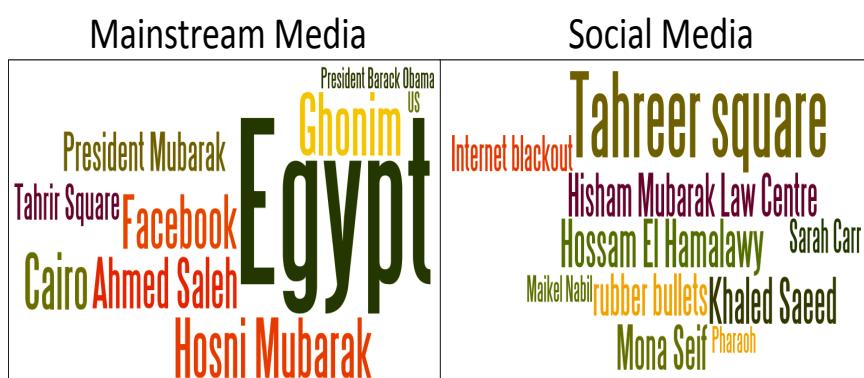


FIGURE 4.3: Top 10 entities from mainstream media and blogs.

An initial analysis of the top 10 entities obtained from the top 10 search results related to “Egyptian Revolution” from two mainstream media channels (BBC and CNN), and from blogs during the time of the revolution is shown in Figure 4.3. The top entities from

the mainstream media channels are generic and are quite obvious for the event. In contrast, the top entities from the blogs are very specific to the event. The activists like ‘Mona Seif’, ‘Sarah Carr’, ‘Maikel Nabil’ and ‘Hosam El Hamalwy’ were very closely involved, and were responsible for mobilizing the event. The entities like ‘Internet Black-out’ and ‘Khaleed Saeed’ were central to the event. Moreover, the presence of entities like ‘Facebook’ and ‘Ghonim’ (who was responsible for spreading the event in Facebook) among the top mainstream media entities also indicates the significance of social media in the event.

A person interested to know about an event in detail, may miss out the novel and specific information available in social media by relying on the top results from the popular search engines. It is a challenge for the current ranking schemes to retrieve the event-specific social media content from the long tail. Moreover, in the words of Chris Anderson [123], “*With an estimated 15 million bloggers out there, the odds that a few will have something important and insightful to say are good and getting better.*” This also motivated to look for techniques in this dissertation, that would help in identifying these otherwise buried sources providing highly specific information related to an event.

4.6 Sparse Link Structure

The casual nature of the users posting content in social media channels gives rise to the challenge of sparsity in link structures. Most often the users who posts content, do not provide links to the original source of information. Also, the behavior of linking to other similar content or building citation networks between information posted about the same topic is completely absent among the social media users. This creates an extremely sparse link structure between the user-generated posts. This further creates problems for the traditional and state-of-the-art searching techniques such as PageRank [124] that performs well in ranking web pages.

This dissertation, introduces and defines novel implicit relationships, intrinsic to content available in the social media channels for solving this issue, and ranks them showing better performances than some of the popular techniques.

4.7 Sampling Bias

Most commonly used method for obtaining data samples from social media websites is by using their application programming interfaces (APIs). Given the humungous amounts of data produced in real-time, the APIs cannot provide all the data to every single

API requests. The requests are often made through a query interface by passing certain query parameters to the APIs. The amount of data returned against the queries may vary. This depends upon the popularity of the content related to the query. For example, in Twitter, studies have estimated that by using Twitter's Streaming API users can expect to receive anywhere from 1% of the tweets to over 40% of tweets in near real-time⁹. The only way to get access to all the tweets is to buy the firehose service, which is seldom done for academic purposes. Other real-time social media publishing services mostly follow the same model. Therefore, this might lead to biasness in the samples collected for studying event related phenomenon and for tracking all the important event related information being produced in real-time.

4.8 Lack of Evaluation Datasets

There is a lack of ground truth evaluation data for most of the social media text mining tasks. In traditional data mining research, there is often two types of datasets. One of them is known as training dataset and the other is known as test dataset. The models are trained or developed using the training datasets and are evaluated on test datasets. Thus, the test datasets act as the ground truth. The

⁹<https://www.brightplanet.com/2013/06/twitter-firehose-vs-twitter-api-whats-the-difference-and-why-should-you-care/>

test dataset for various text mining tasks is mostly not available for social media. It is often the duty of the researchers to create new test datasets in order to solve a specific task in social media. Sometimes this data might not be a benchmark dataset due to various unwanted noise and human error or perception in annotating the data. This might lead to wrong assumptions and false results.

Popularly used validation techniques that are widely accepted by the research community for circumventing the evaluation challenges are used in this dissertation. Both automated as well as manual methods of evaluation, are considered according to the scenario. Independent annotators are used for all the annotation tasks. They are properly educated about the tasks and the events. Please refer (Chapter 6) for more details.

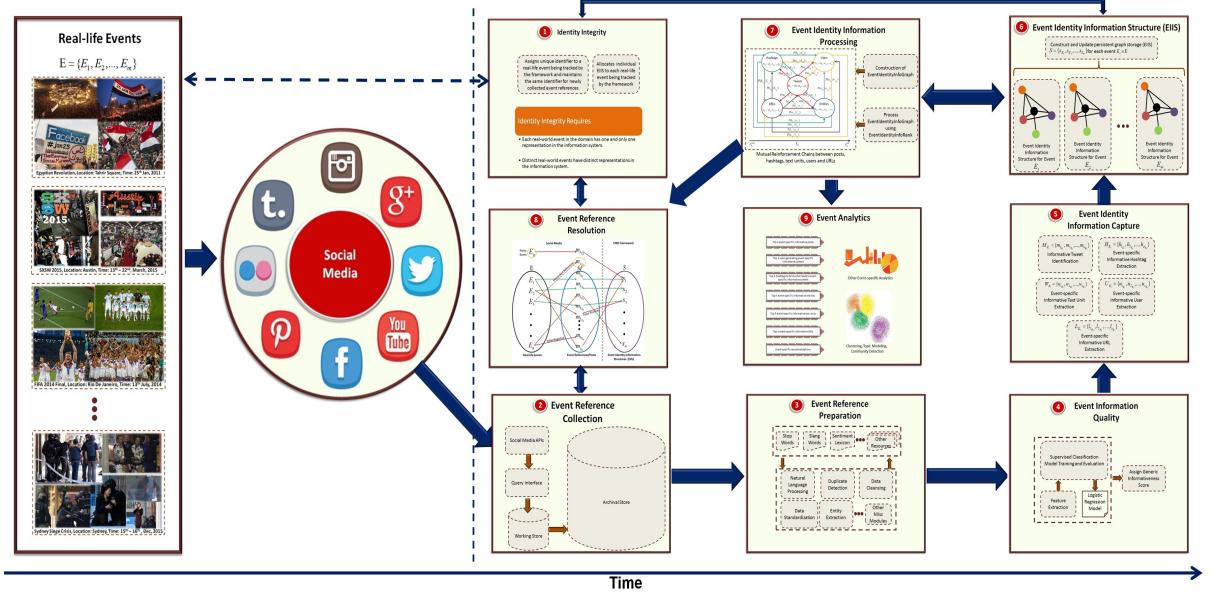
Chapter 5

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

This chapter discusses about the Event Identity Information Management framework in order to solve the problems posed in Chapter 2, Section 2.5. A high level view of the EIIM components and processes for social media is shown in Figure 5.1. These components provide a generic framework on which any EIIM system based on social media references could be built. The various components of the framework go through cycles of interactions with each other over time, which is known as EIIM life cycle. At the heart of the framework lies the Event Identity Information Structure (EIIS) and Event Identity Information Processing, which manages the identity information related to a particular event and also helps in resolving high

quality, informative references to the event from social media. Next, we give a detailed explanation of each component.

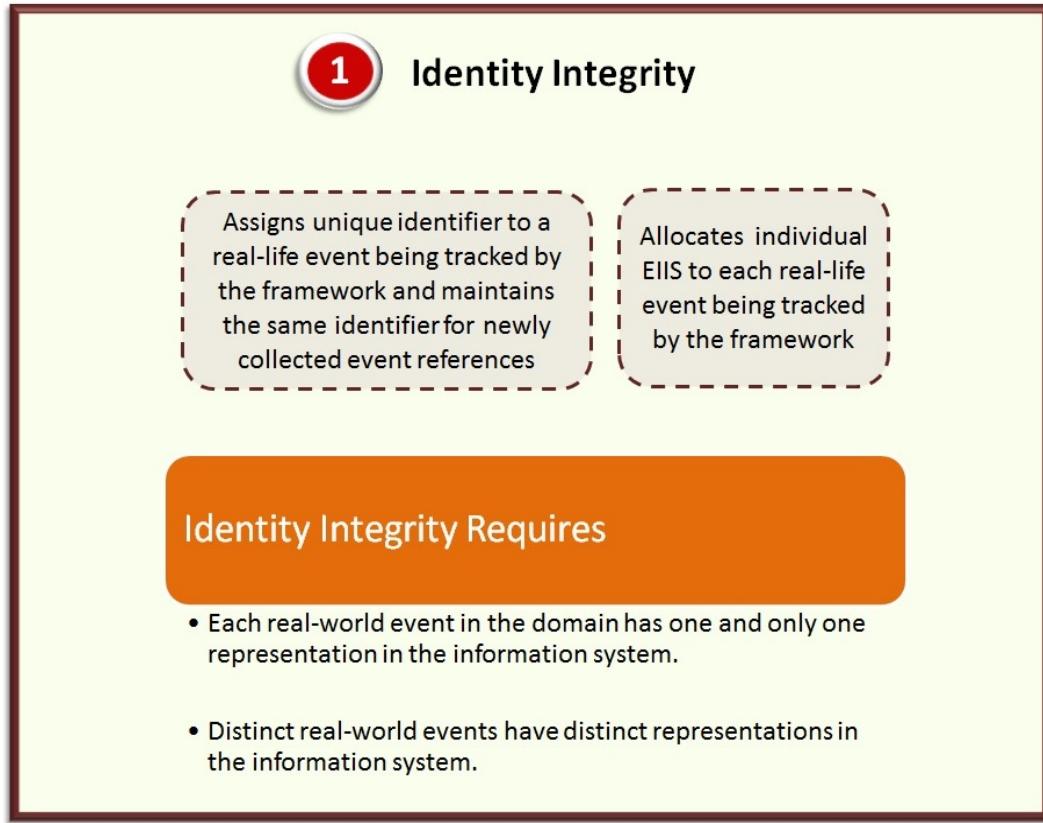
FIGURE 5.1: Event Identity Information Management Life Cycle for Textual Content in Social Media.



5.1 Identity Integrity

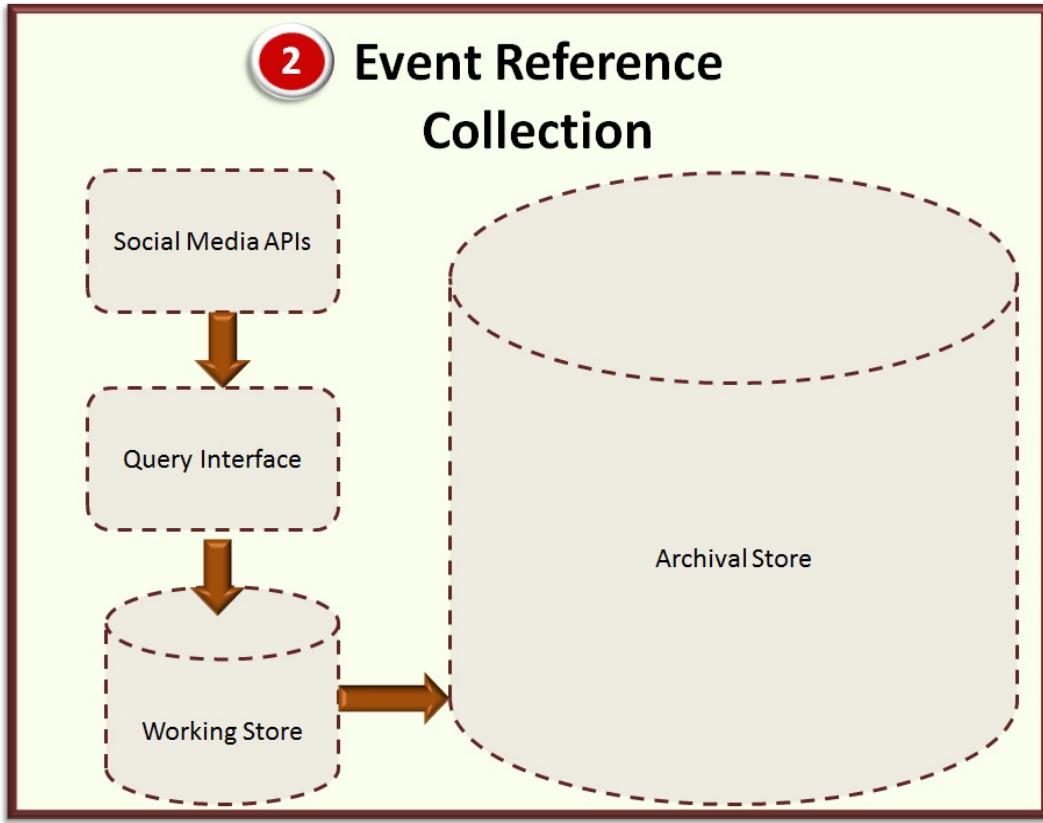
One of the fundamental goals of the proposed framework is to maintain a one-to-one correspondence between real-world events being monitored and the Event Identity Information Structure (EIIS) of the corresponding events for ensuring identity integrity. Therefore, a separate EIIS is maintained corresponding 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

FIGURE 5.2: Identity Integrity component of the EIIM life cycle.



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.

FIGURE 5.3: Event Reference Collection component of the EIIM life cycle.



5.2 Event Reference Collection

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. A query interface is implemented that allows an

user of the system to pass query parameters (for example event related hashtags and key words) that bootstraps the data collection and event tracking process. As already shown in Table 4.2 most of the popular social media channels allow hashtags, the data for the experiments were collected using a popular hashtag for the respective events.

Due to extreme popularity of Twitter, data from it was collected for representing the microblog genre and short textual social media references. Four million tweets (approx) related to five different events were collected. Details of the collected event references are provided in Table 5.1. The tweets were collected over the given period of time, by providing a popular hashtag to the Twitter streaming API as shown in Table 5.1 (for details about Twitter Data Collection please refer Appendix A). Only English language tweets are considered for the experiments as the available natural language toolkits performs well for English, and the annotators used for different tasks are only proficient in English language.

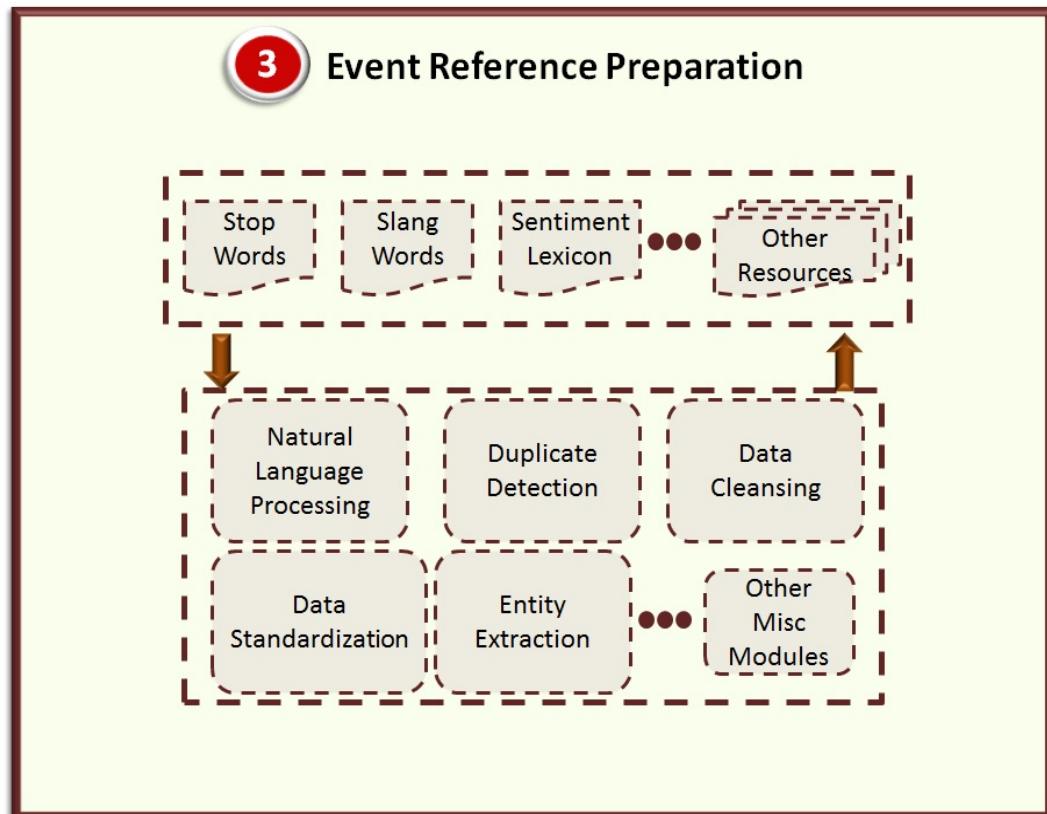
5.3 Event Reference Preparation

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

TABLE 5.1: Details of data collected for analyzing event related tweet content.

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

FIGURE 5.4: Event Reference Preparation component of the EIIM life cycle.



suitable for further processing by the other components of the EIIM life cycle. Several resources are compiled in order to tackle the challenges posed by short informal text prevalent in social media.

The tweets collected using the previous component goes through the following pre-processing steps:

5.3.1 Parts-of-speech tagging

The Natural Language Toolkit (<http://nltk.org>) POS tagger is used for tagging the raw text of the tweets. All the words of a tweet is assigned one of the following parts-of-speech:

- Noun
- Adjective
- Verb
- Adverb
- Preposition
- Interjection
- Pronoun
- Article

The tagged tweet is also separately stored and maintained. These tags are used later in different components down the pipeline. The Penn Treebank tags¹ are used for tagging and is parsed accordingly for identifying words with a certain parts-of-speech.

5.3.2 Special Character Detection

All the special characters that are not alphanumeric are detected and the total number of special characters in a tweet is stored.

5.3.3 Data Cleansing

The raw text of the tweet is extracted and cleaned. All the user mentions, hashtags, retweet symbols, URLs and special characters are removed and the entire tweet is converted into lowercase.

5.3.4 Duplicate Detection

The tweets after cleaning are assigned a md5 hash code, which helps in detecting duplicate content. Tweets having the same hash code are considered to be redundant copies of each other, and only a single copy of the tweet is finally stored in the database. This technique also helps in detecting the retweets of a tweet that contain the same content. Also there are certain tweets that are shared with the same

¹https://www.ling.upenn.edu/courses/Fall_2003/ling001/penn_treebank_pos.html

content but has different user mentions and with different combination of the words expressing the content. For example, the following tweets talks about the same video shared by mashable and does not present any new information.

- RT @mashable: Timelapse video reveals massive size of New York City protests <http://t.co/zhqHpkDLk1> #MillionsMarch-
NYC <http://t.co/WktxssAfDp>

- RT @dianebhartford: “@mashable: Timelapse video reveals massive size of New York City protests <http://t.co/CE0VIyHnLe> #MillionsMarchNYC ht...

After the data cleansing step, the duplicate detection scheme identifies both the tweets to be same and only maintain a single copy. This process occurs in real-time. Whenever, a new tweet is obtained from the straming API, the md5 hashcode is calculated after going through the previous data pre-processing steps. A hashtable is maintained in the memory that is constantly searched for the presence of the generated hashcode. If the hashcode is already present then the tweet is dropped and not stored.

5.3.5 Stop Word Detection and Elimination

A list of English stop words is compiled that is publicly shared in the following URL :

- <https://github.com/dxmahata/EIIMFramework/blob/master/CodeBase/EventIdentityInformationManagement/Resources/englishStopwords.txt>

This list is used for detecting the stop words in English language tweets. The stop words are eliminated and the number of stop words detected is recorded.

5.3.6 Slang Word Detection

Slang words commonly used on the Internet and twitter specific slang publicly shared by FBI² is combined together for compiling a list of English slang words. This list is used for detecting and extracting the slang words from the tweets. The number of slang words detected is recorded. This list is also used later to detect the slang hashtags and slang text units.

The compiled list of twitter specific slang words is publicly shared and can be obtained from the following URL :

²<https://www.documentcloud.org/documents/1199460-responsive-documents.html#document/p1>

- <https://github.com/dxmahata/EIIMFramework/blob/master/CodeBase/EventIdentityInformationManagement/Resources/slangWords.txt>

5.3.7 Feeling Word Detection

A list of words expressing feelings on the Internet, obtained from wefeelfine.org is used for detecting and extracting the feeling words from a tweet. The number of feeling words detected is recorded and the extracted feeling words are stored. This list is also publicly available and can be obtained from the following URL :

- <https://github.com/dxmahata/EIIMFramework/blob/master/CodeBase/EventIdentityInformationManagement/Resources/feelingWords.txt>

5.3.8 Tokenization

The tweet obtained after performing the data cleansing steps and elimination of stop words are tokenized into unigram and bigram tokens using the tokenizer module available in NLTK. The set of tokens thus obtained are stored separately.

5.3.9 Stemming

The unigram tokens obtained after tokenization are stemmed and a separate list of stemmed tokens are stored. A standard Porter stemmer available with NLTK library is used for the purpose.

5.3.10 Tweet Meta-data Extraction

Several meta-data that are associated with each tweet obtained from the JSON response of the streaming API are extracted. Some of these meta-data are :

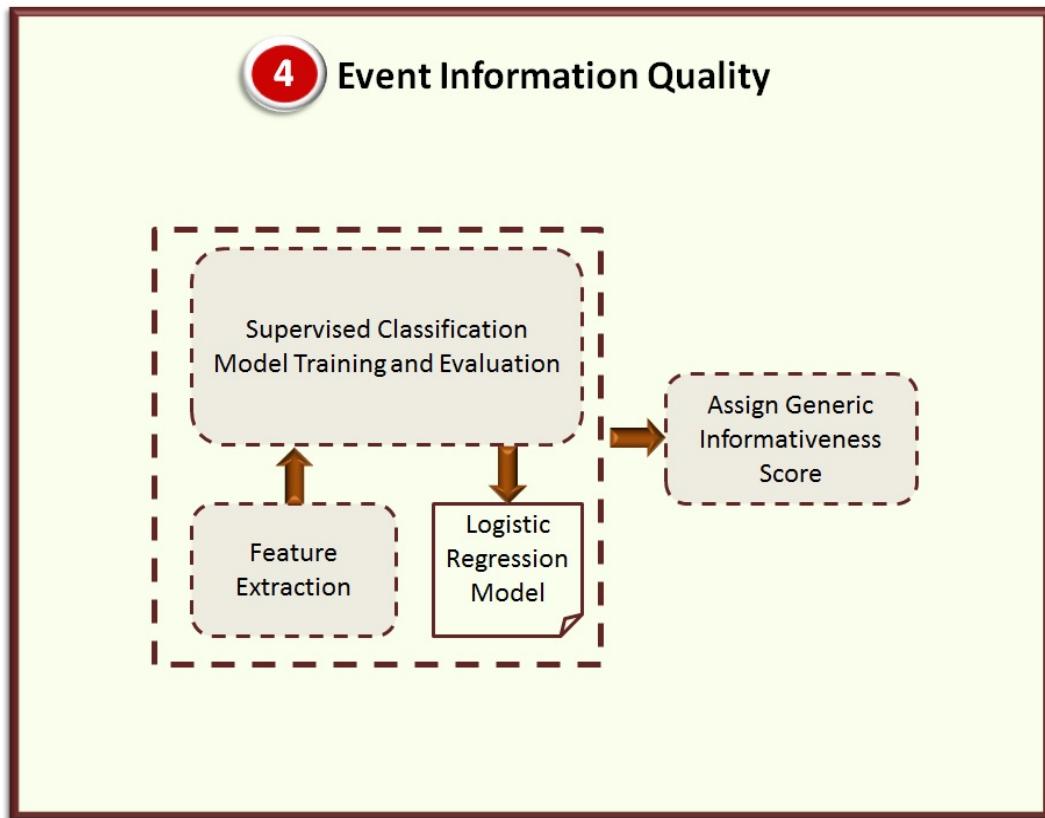
- Expanded URLs
- Hashtags
- Retweet Counts
- Favorite Counts
- User Mentions
- User Follower Counts
- Verification Information
- Time Information

5.3.11 Named Entity Extraction

Named entities such as name of persons, animals, places, cars and organizations are extracted from the raw tweets. For this purpose the entity extraction service of AlchemyAPI³ is used.

5.4 Event Information Quality

FIGURE 5.5: Event Information Quality component of the EIIM life cycle.



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

³<http://alchemyapi.com>

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 related informative references we implemented a logistic regression classifier. Once the classifier is trained, it is used for assigning generic informativeness score to the tweets in real-time as they are collected using the streaming API. The component aids in solving the problem of information overload. Just like an user searching for relevant informative content about an event faces the challenging situation of information overload as discussed in Chapter 4, Section 4.1, it is also a challenge for automated systems to process the huge amount of content coming at high velocity and extract useful information out of it. By filtering out the tweets that are less likely to contain useful and high quality information it solves the problem of information overload for the other components of the EIIM framework.

5.4.1 Annotated Dataset

We use a publicly available annotated dataset from the CrisisLex⁴ website shared by Olteanu et al. [125]. The collection includes tweets collected during 26 large crisis events in 2012 and 2013, with about 1,000 tweets labeled per crisis for informativeness (i.e. “informative”, or “not informative”), information type, and source. 28,000 tweets

⁴<http://crisislex.org/data-collections.html>

(about 1,000 in each collection) were labeled by crowdsourced workers according to informativeness (informative or not informative), information types (e.g. caution and advice, infrastructure damage), and information sources (e.g. governments, NGOs).

For example, for the Colorado wildfire⁵ event, the following tweets were assigned labels of “related and informative”, “related but not informative”, and “not related”, respectively.

- *Related and Informative* - #Media Large wildfire in N. Colorado prompts evacuations: Crews are battling a fast-moving wildfire
<http://t.co/ju1BGTKH> #Politics #News
- *Related but not Informative* - RT @LarimerSheriff: #HighPark-Fire update
<http://t.co/hBy5shen>
- *Not Related* - #Intern #US #TATTOO #Wisconsin #Ohio #NC #PA #Florida #Colorado #Iowa #Nevada #Virginia #NV #mlb Travel Destinations ;
<http://t.co/TIHBJKF2>

Not all the tweets were in English language. We selected English language tweets for training the logistic regression model. There were only 9729 tweets.

⁵http://en.wikipedia.org/wiki/2012_Colorado_wildfires

5.4.2 Feature Selection and Training

In order to train the model we assigned a score of 1 to the tweets that were labeled ‘related and informative’, and all the other tweets labeled as ‘related-but not informative’, and ‘not related’ were assigned a score of 0. The choice of features was governed by previous works related to identifying high quality information from Twitter [38, 41, 51, 126]. The list of features selected for the model are:

1. **Has URL** - has a value of 1 if the tweet contains URL or else has a value of 0.
2. **Number of Words** - total number of unigram tokens extracted from the raw tweet text.
3. **Number of Stop Words** - total number of English stop words detected in the raw tweet.
4. **Number of Feeling Words** - total number of feeling words detected in the raw tweet.
5. **Number of Slang Words** - total number of slang words detected in the raw tweet.
6. **Number of Hashtags** - total number of hashtags used in the raw tweet.

7. **Number of User Mentions** - total number of user mentions detected in the raw tweet.
8. **Tweet Length** - total number of characters used in the raw tweet.
9. **Unique Characters** - total number of unique characters used in the raw tweet.
10. **Special Characters** - total number of special characters detected in the tweet.
11. **Favorite Count** - total number of favorite count of the tweet at the time it was collected.
12. **Retweet Count** - total number of retweet count of the tweet at the time it was collected.
13. **Verified** - has a value of 1 if the user posting the tweet is a verified user by Twitter or else has a value of 0.
14. **Number of Nouns** - total number of nouns detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as nouns.
15. **Number of Adjectives** - total number of adjectives detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as adjectives.

16. **Number of Verbs** - total number of verbs detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as verbs.
17. **Number of Adverbs** - total number of adverbs detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as adverbs.
18. **Number of Pronouns** - total number of pronouns detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as pronouns.
19. **Number of Interjections** - total number of interjections detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as interjections.
20. **Number of Articles** - total number of articles detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as articles.
21. **Number of Prepositions** - total number of prepositions detected in the tweet, without considering the hashtags and the user mentions whenever they are tagged as prepositions.
22. **Formality** - which is defined as follows, $\text{Formality} = (\#\text{nouns} + \#\text{adjectives} + \#\text{prepositions} + \#\text{articles} - \#\text{pronouns} - \#\text{verbs} - \#\text{adverbs} - \#\text{interjections} + 100)/2$ and is proposed in [127].

#nouns, denotes the number of nouns detected in the tweet, and so on.

TABLE 5.2: Evaluation measures for logistic regression model.

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

5.4.3 Model Evaluation

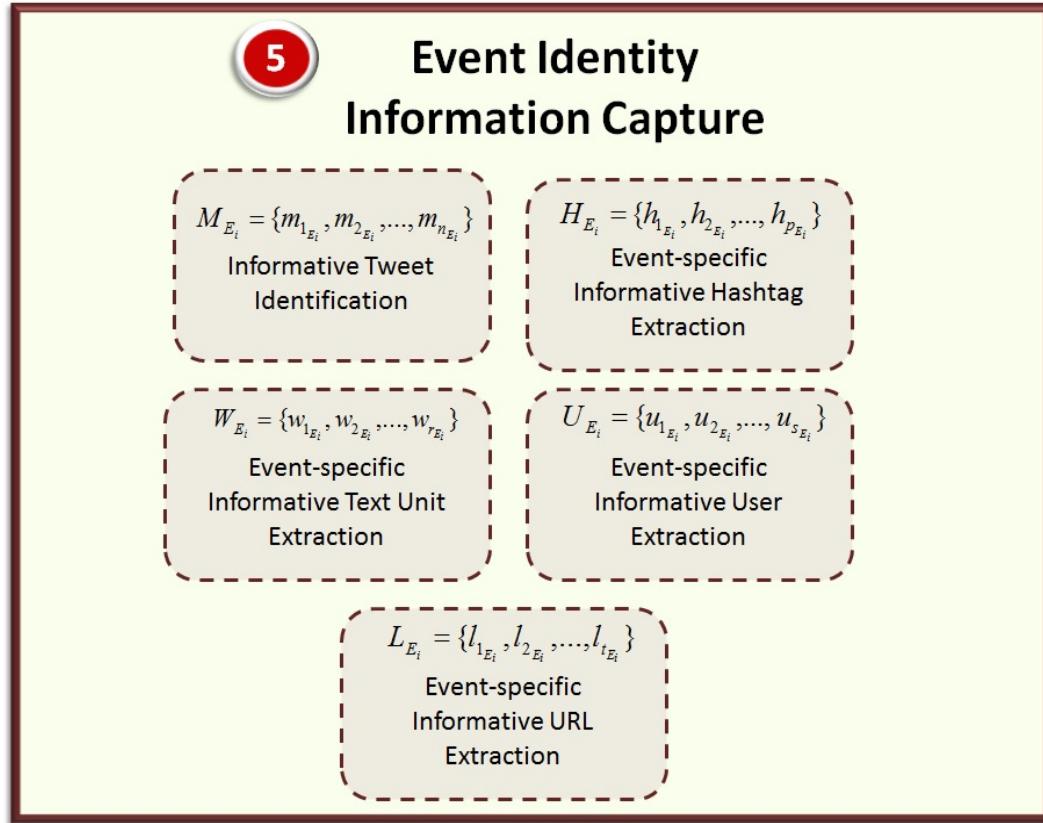
10-fold cross validation was performed, with ‘l1’ penalty, resulting in a model with an accuracy of 76.64%. Table 5.2 lists the evaluation measures obtained while training the classifier. The ROC AUC Score of the classifier is 0.6934.

5.4.4 Assignment of Generic Informativeness Score

The 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.

5.5 Event Identity Information Capture

FIGURE 5.6: Event Identity Information Capture component of the EIIM life cycle.



The main functions of this component are:

- This component 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 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 following three events.

- CPAC 2014.
- SXSW 2014.
- Sochi Winter Games 2014.

The analysis and the conclusions we made from it, is presented next.

5.5.1 Content Analysis of Event Related Tweets

Details of the data collected for the analysis are provided in Table 5.1. The data collection task was accomplished by ‘Event Reference Collection’ component and was then preprocessed by the ‘Event Reference Preparation’ component.

Twitter allows its users to post short messages with a limitation of 140 characters. Users not only post plain textual content in their messages but also share URLs, linking to other external websites, images and videos. The images and videos are labeled as media

elements by Twitter. Apart from curating new content, the users also share content produced by others. This activity is known as *retweeting*, and such tweets are preceded by the special characters *RT*. The messages are normally written by a single person and are read by many. The readers in the context of Twitter are known as *followers*, and the user whom the other users follow is considered as their *friend*. Any user with good intent either share messages that might be of interest to his followers, or for joining conversations on topics of his interest. The ‘@’ symbol followed by the username commonly known as *user mentions*, is used for mentioning other users in tweets for initiating conversation with them.

The concise and informal content of a tweet is often contextualized by the use of a crowdsourced annotation scheme called *hashtags*. Hashtags are a sequence of characters in any language prefixed by the symbol ‘#’ (for e.g. #icwsm2015). They are widely used by the users in order to add context to the tweets, categorizing the content based on a topic, join conversations related to a topic, and to make the tweets easily searchable by other interested users. They also act as strong identifiers of topics [128]. When tweeting about real-life events the users also tend to use hashtags in order to post event-specific content. For example ‘#Egypt’ and ‘#Jan25’, were among the most popular hashtags in Twitter used for spreading, organizing and analyzing information related to ‘Egyptian Revolution of 2011’

[129].

Given the mechanisms of user interactions and content production in Twitter, we started our analysis with the assumption that the content of a tweet is primarily composed of hashtags, words for expressing and conveying information, and URLs that lead to additional information about the content. We plotted the distribution of occurrences of all the hashtags, tokens and shared URLs for each event. Due to the short length of the tweets we only considered unigram tokens. We also plotted the distribution of the number of tweets posted by the users. We observed a power law distribution for all of them (refer Figures 5.7, 5.8, 5.9 and 5.10). This gave us an intuition that there are skewed sets of hashtags and tokens that are widely used for posting content related to an event. There is also a specific set of URLs that become popular in comparison to others and a set of users who are more active than others in posting event related content.

Our second step was to use the logistic regression model developed for the ‘Event Information Quality’ component and assign informativeness 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.

FIGURE 5.7: Distribution of hashtags in event related tweets.

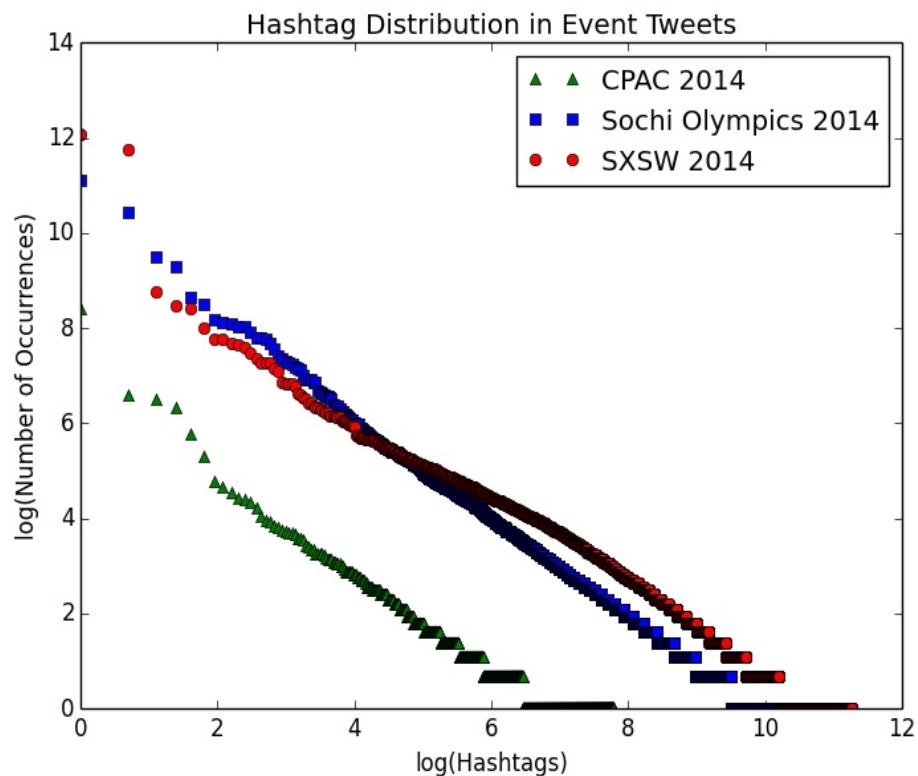


FIGURE 5.8: Distribution of tokens in event related tweets.

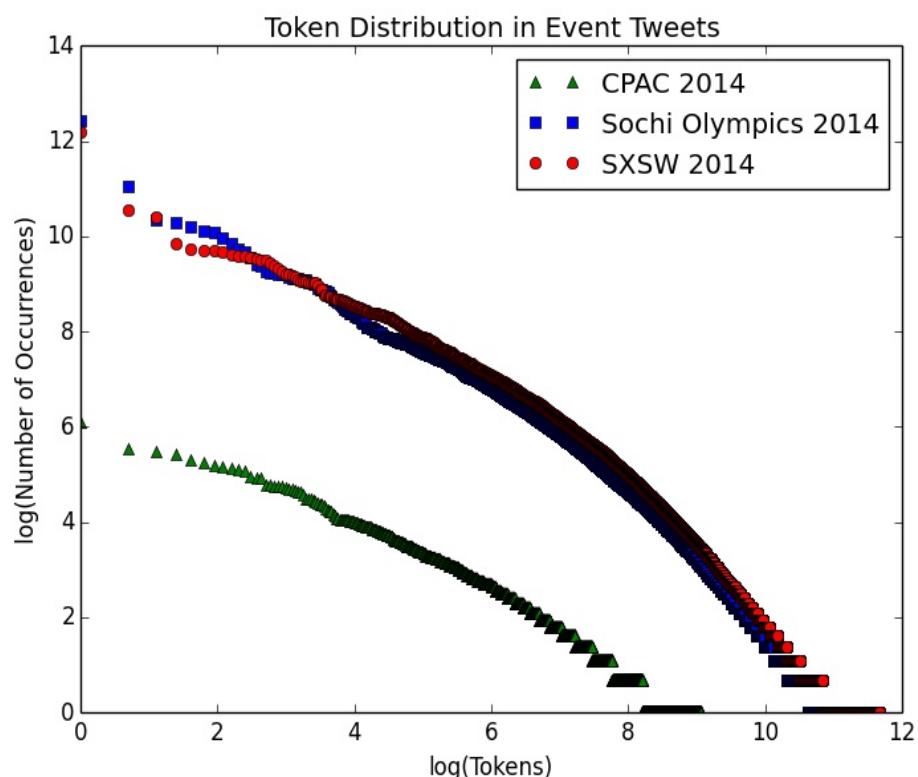


FIGURE 5.9: Distribution of URLs in event related tweets.

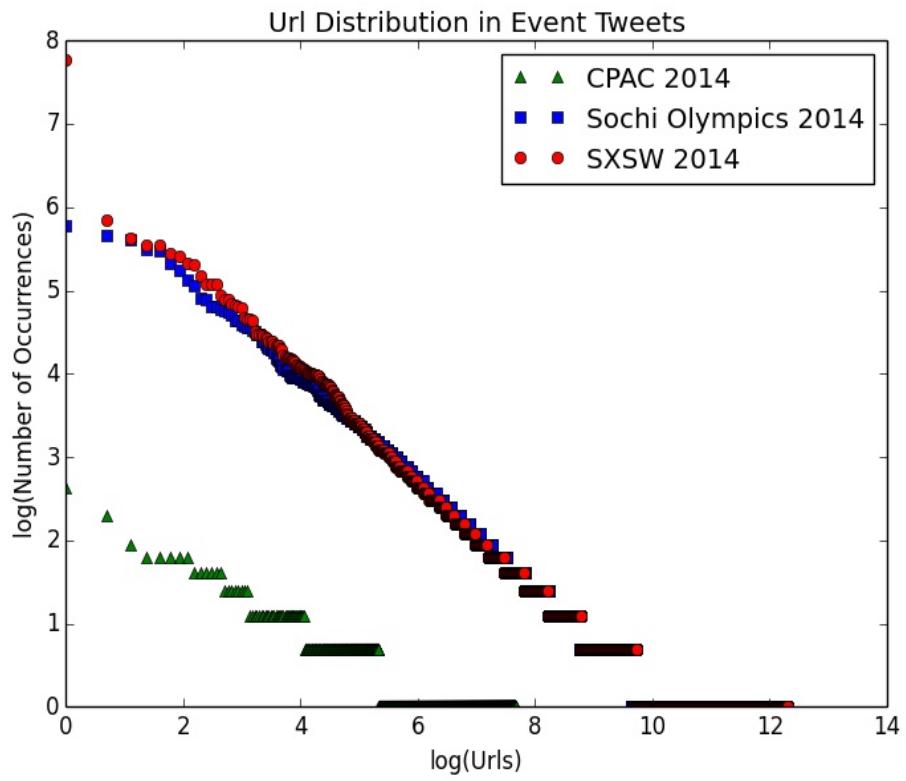
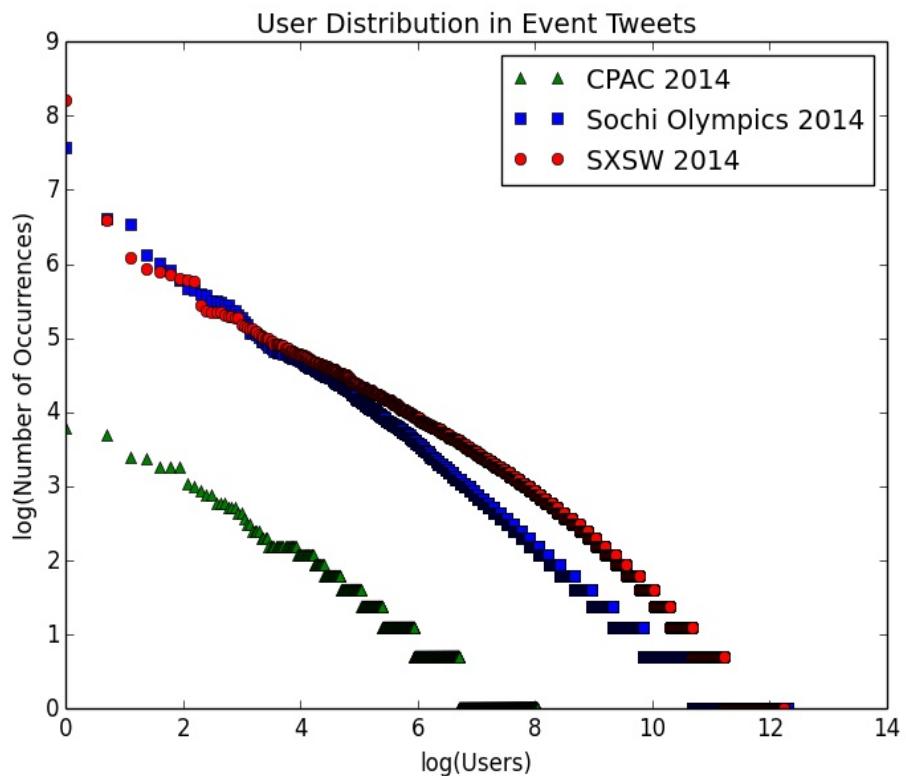


FIGURE 5.10: Distribution of users in event related tweets.



We calculated the average values of different content characteristics of the tweets. 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 5.11.

FIGURE 5.11: 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% |

As presented in the table, some of the observations for all the three events are,

- 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.
- 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 pointed out in Chapter 4, Section 4.2.
- 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. We made the following conclusions based on the observations:

- 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.

5.5.2 Event Identity Information Units

After the observations in the previous section we conclude that the informative tweets in general are characterized by wordiness, occurrences of URLs and are posted by users with high follower count. These characteristics are also the primary features that distinguish informative from non-informative content. Although, presence of hashtags is not a good indicator of informativeness, yet it is a strong

identifier of a topic as already pointed by [128]. Popular hashtags for an event might be used maliciously. On the other hand, the presence of a popular hashtag in a wordy tweet consisting of words popular for the event, along with a popular URL, posted by an influential user is highly likely to contain event-specific content. Therefore, it is intuitive that given a stream of tweets for an event an optimal combination of event related popular text units (words, unigrams, bigrams etc), hashtags, and URLs, posted by an influential user in a tweet, is one of the key indicators for identifying event-specific informative content. It would be highly unlikely for a tweet to contain all of these and yet not convey useful event-specific information. Based on the above analysis we decided to build the EIIS for an event E_i , composed of the following event identity information units:

1. A set of tweets $M_{E_i} = \{m_{1_{E_i}}, m_{2_{E_i}}, \dots, m_{n_{E_i}}\}$, related to the event E_i , having high chances of containing informative content.
2. A set of hashtags $H_{E_i} = \{h_{1_{E_i}}, h_{2_{E_i}}, \dots, h_{p_{E_i}}\}$, used for annotating the tweets ($\in M_{E_i}$) related to event E_i .
3. A set of text units $W_{E_i} = \{w_{1_{E_i}}, w_{2_{E_i}}, \dots, w_{r_{E_i}}\}$, used for expressing textual content in tweets ($\in M_{E_i}$), related to event E_i .
4. A set of URLs $L_{E_i} = \{l_{1_{E_i}}, l_{2_{E_i}}, \dots, l_{t_{E_i}}\}$, shared in the tweets ($\in M_{E_i}$) related to event E_i .

5. A set of users $U_{E_i} = \{u_{1_{E_i}}, u_{2_{E_i}}, \dots, u_{s_{E_i}}\}$, tweeting the tweets ($\in M_{E_i}$), about the event E_i .

5.5.3 Extracting Event Identity Information Units

The event identity information units for an event E_i , as defined above are extracted from the event dataset. Following steps are taken:

- A threshold for the informativeness score assigned in the previous component is set between 0.0-1.0, for extracting the tweets ($\in M_{E_i}$). We set a threshold of 0.7. Therefore, the tweets having an informativeness score greater than or equal to 0.7 are filtered out and comprises the set M_{E_i} .
- The hashtags that were extracted in the ‘Event Reference Preparation’ step from the tweets $\in M_{E_i}$, are used for populating the set H_{E_i} . However, the hashtags that matches the slang words and stop words are not considered. This is done in order to ensure good quality of information contextualized by the hashtags.
- The nouns that were extracted in the ‘Event Reference Preparation’ step from the tweets $\in M_{E_i}$, are considered as the text units ($\in W_{E_i}$). The nouns that matches the slang words are not considered. This is done in order to ensure good quality of textual content represented by the nouns. In another experiment,

we considered the extracted named entities as the text units.

We report our results and compare the results obtained in both the cases in the next Chapter.

- The expanded URLs from the meta-data of the tweets $\in M_{E_i}$ populates the set L_{E_i} .
- The meta-information of the users posting the tweets $\in M_{E_i}$, represented by their user ids, is extracted for populating the set U_{E_i}

These event identity information units forms the Event Identity Information Structure (EIIS), as explained next.

5.6 Event Identity Information Structure (EIIS)

This is the component that maintains a persistent EIIS as introduced in Chapter 2, Section 2.5, for each individual event tracked by the framework and updates the metadata of the EIIS throughout the EIIM life cycle. Due to the unstructured nature of the social media references and evolving nature of the events, we store the event identity information units extracted by the previous component along with their associated meta-data, in a persistent graph data structure stored in the database. We update the meta-data

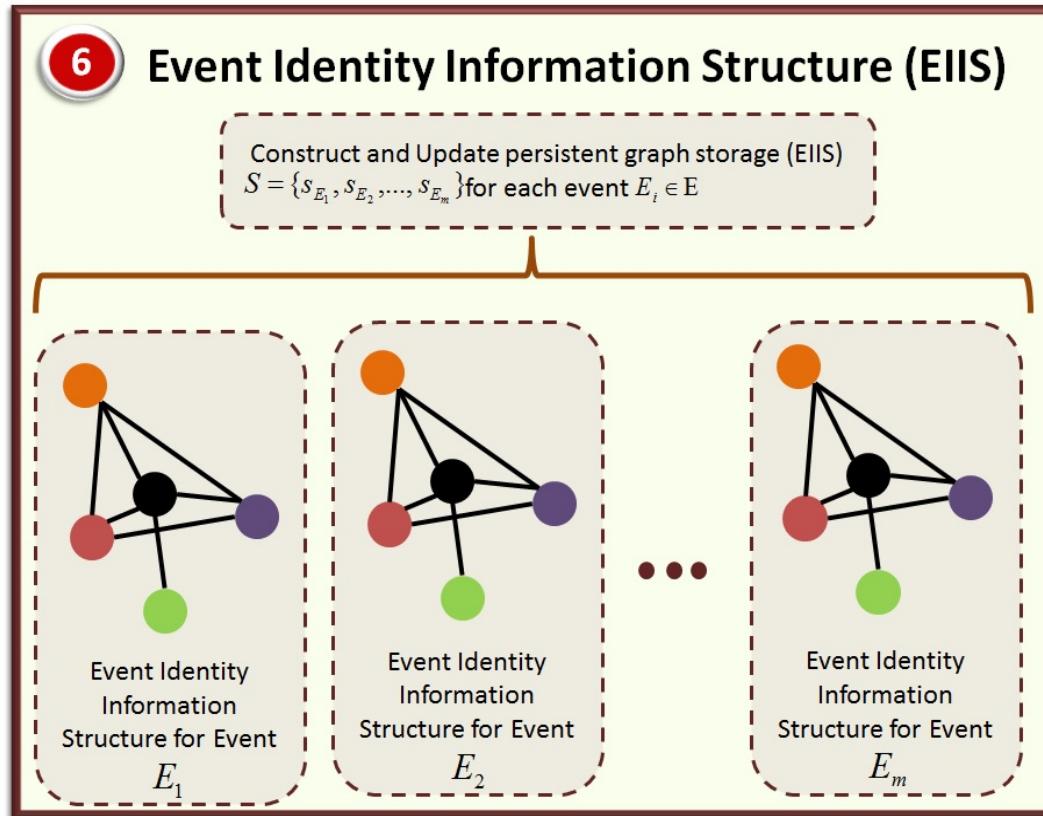
related to each node of the graph using the normal database updation queries. Adjacency lists are used for representing the graph. The choice of storing the event identity information units in a graph structure is also motivated by the wide array of graph processing algorithms used for natural language processing and text mining operations. We show the efficacy and the advantages of a graph in the next section.

- Therefore the EIIS is a graph $\mathbf{G}_{\mathbf{E}_i} = (\mathbf{V}_{\mathbf{E}_i}, \mathbf{D}_{\mathbf{E}_i})$, where $\mathbf{V}_{\mathbf{E}_i} = \mathbf{M}_{\mathbf{E}_i} \cup \mathbf{H}_{\mathbf{E}_i} \cup \mathbf{V}_{\mathbf{E}_i}$, $\mathbf{L}_{\mathbf{E}_i}$, is the set of vertices and $\mathbf{D}_{\mathbf{E}_i}$ is the set of directed edges between different vertices. Whenever two vertices are associated, there are two edges between them that are oppositely directed. For example, if a tweets consist of hashtags, text units, URLs and is posted by an user, then there are bi-directed edges between each one of them.

5.7 Event Identity Information Processing

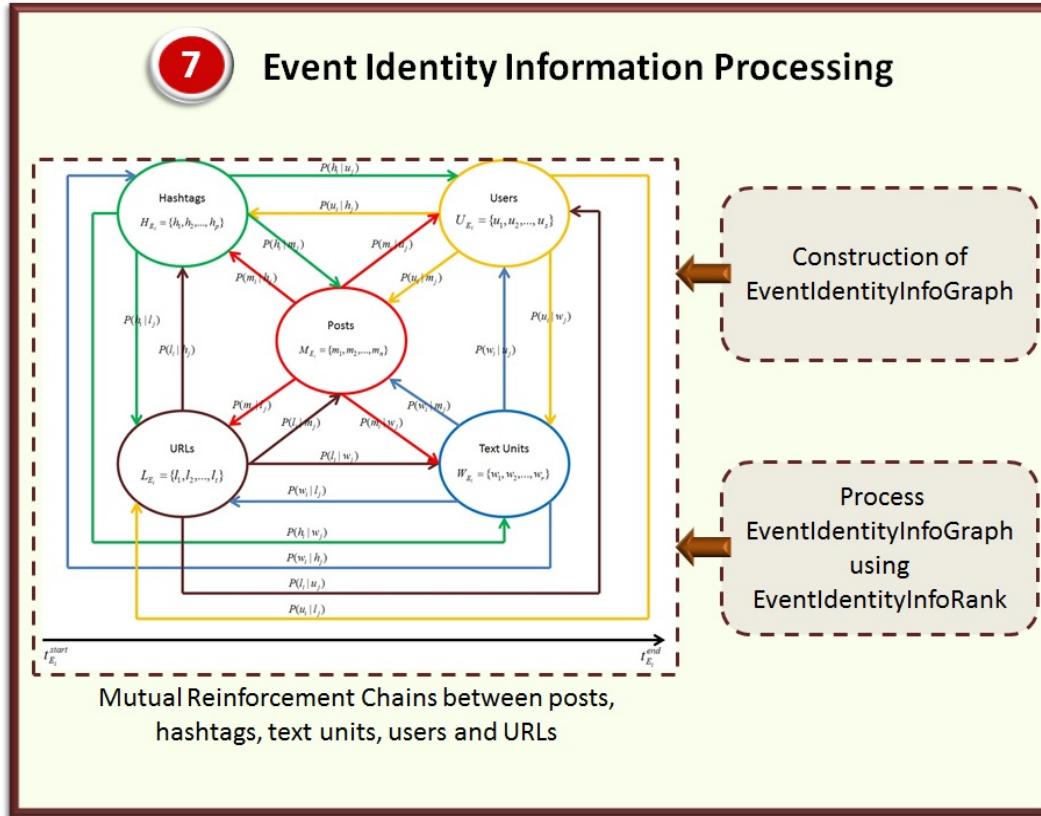
This is the most important processing component of the EIIM life cycle and is at the heart of the entire framework. We make our most novel contributions in this component. The component is mainly divided into two sub-components:

FIGURE 5.12: Event Identity Information Structure component of the EIIM life cycle.



- **EventIdentityInfoGraph** - that represents and defines novel relationships between the vertices of the graph \mathbf{G} representing the EIIS.
- **EventIdentityInfoProcess** - processes the *EventIdentityInfo-Graph* in order to rank its nodes and identify the top most informative event identity information units that acts as inputs to the next two components of the EIIM life cycle.

FIGURE 5.13: Event Identity Information Processing component of the EIIM life cycle.



5.7.1 EventIdentityInfoGraph

We implement a novel graph structure - *EventIdentityInfoGraph*, which is dynamically generated from the graph \mathbf{G} (EIIS), after a configurable interval of time, using following assumptions.

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.*
-
- 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,*

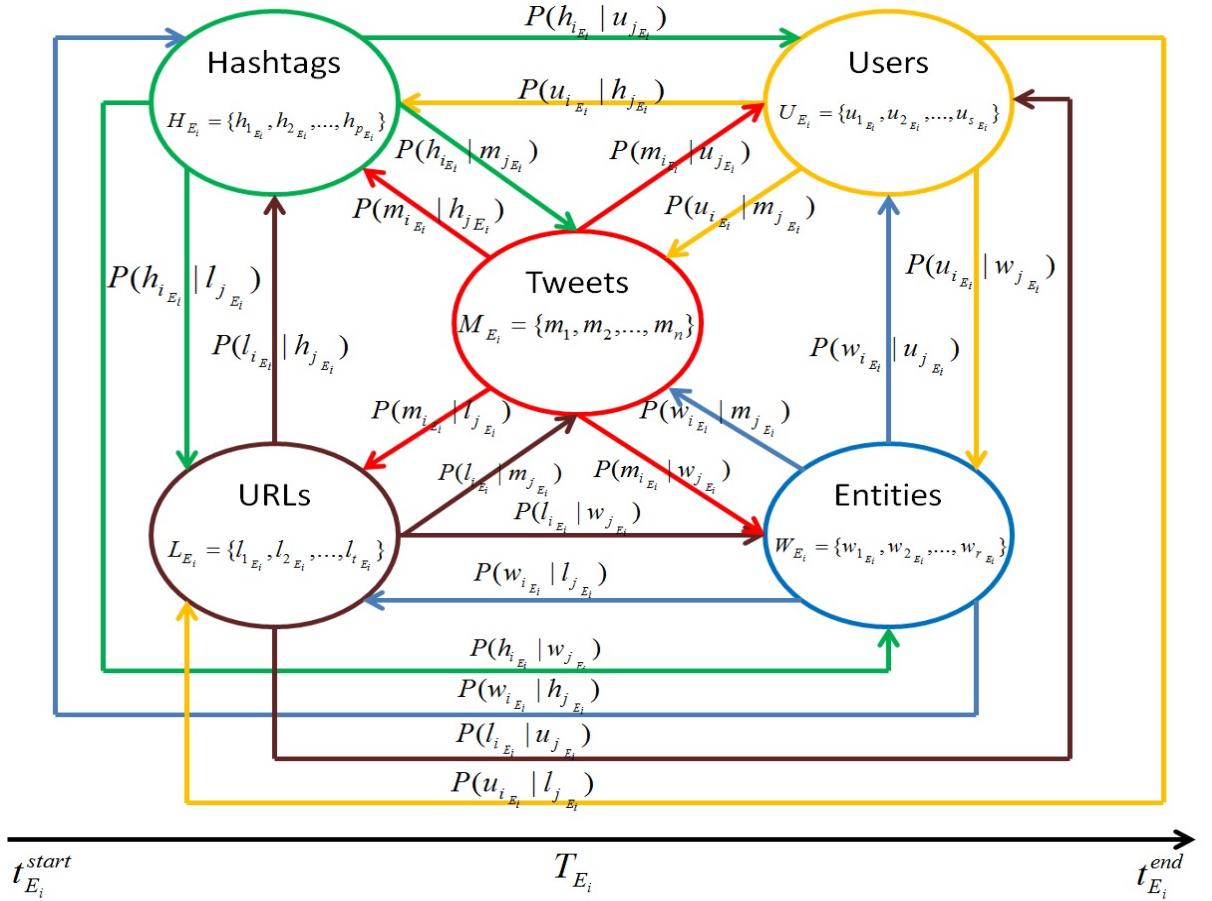
(d) *event-specific informative URLs.*

- a *URL is an event-specific informative URL* if it is strongly associated with:
 - (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* [130] for the event E_i as shown in Figure 5.14. We represent this relationship in a graph $\mathbf{G}_{\mathbf{E}_i(t)} = (\mathbf{V}_{\mathbf{E}_i(t)}, \mathbf{D}_{\mathbf{E}_i(t)})$, which we call as *EventIdentityInfoGraph*, where $V_{E_i(t)} = M_{E_i} \cup H_{E_i} \cup W_{E_i} \cup U_{E_i} \cup L_{E_i}$, is the set of vertices and $\mathbf{D}_{\mathbf{E}_i(t)}$ is the set of directed edges between different vertices. The graph $G_{E_i(t)}$ is basically, a snapshot of the EIIS structure of the event E_i at time t.

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 as given by equations 5.1-5.18.

FIGURE 5.14: Mutual Reinforcement Chains in Twitter for an event.



$$P(h_{i_{E_i}} | w_{j_{E_i}}) = \frac{\text{No. of tweets } h_{i_{E_i}} \text{ and } w_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } w_{j_{E_i}} \text{ occurs}} \quad (5.1)$$

$$P(w_{i_{E_i}} | h_{j_{E_i}}) = \frac{\text{No. of tweets } w_{i_{E_i}} \text{ and } h_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } h_{j_{E_i}} \text{ occurs}} \quad (5.2)$$

$$P(h_{i_{E_i}} | l_{j_{E_i}}) = \frac{\text{No. of tweets } h_{i_{E_i}} \text{ and } l_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } l_{j_{E_i}} \text{ occurs}} \quad (5.3)$$

$$P(l_{i_{E_i}} | h_{j_{E_i}}) = \frac{\text{No. of tweets } l_{i_{E_i}} \text{ and } h_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } h_{j_{E_i}} \text{ occurs}} \quad (5.4)$$

$$P(h_{i_{E_i}} | u_{j_{E_i}}) = \frac{\text{No. of tweets } h_{i_{E_i}} \text{ and } u_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } u_{j_{E_i}} \text{ occurs}} \quad (5.5)$$

$$P(u_{i_{E_i}} | h_{j_{E_i}}) = \frac{\text{No. of tweets } u_{i_{E_i}} \text{ and } h_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } h_{j_{E_i}} \text{ occurs}} \quad (5.6)$$

$$P(w_{i_{E_i}} | l_{j_{E_i}}) = \frac{\text{No. of tweets } w_{i_{E_i}} \text{ and } l_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } l_{j_{E_i}} \text{ occurs}} \quad (5.7)$$

$$(l_{i_{E_i}} | w_{j_{E_i}}) = \frac{\text{No. of tweets } l_{i_{E_i}} \text{ and } w_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } w_{j_{E_i}} \text{ occurs}} \quad (5.8)$$

$$P(u_{i_{E_i}} \mid l_{j_{E_i}}) = \frac{\text{No. of tweets } u_{i_{E_i}} \text{ and } l_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } l_{j_{E_i}} \text{ occurs}} \quad (5.9)$$

$$P(l_{i_{E_i}} \mid u_{j_{E_i}}) = \frac{\text{No. of tweets } l_{i_{E_i}} \text{ and } u_{j_{E_i}} \text{ occur together}}{\text{No. of tweets } u_{j_{E_i}} \text{ occurs}} \quad (5.10)$$

$$P(h_{i_{E_i}} \mid m_{j_{E_i}}) = 1.0 \quad (5.11)$$

$$P(m_{i_{E_i}} \mid h_{j_{E_i}}) = 1.0 \quad (5.12)$$

$$P(w_{i_{E_i}} \mid m_{j_{E_i}}) = 1.0 \quad (5.13)$$

$$P(m_{i_{E_i}} \mid w_{j_{E_i}}) = 1.0 \quad (5.14)$$

$$P(u_{i_{E_i}} \mid m_{j_{E_i}}) = 1.0 \quad (5.15)$$

$$P(m_{i_{E_i}} \mid u_{j_{E_i}}) = 1.0 \quad (5.16)$$

$$P(l_{i_{E_i}} \mid m_{j_{E_i}}) = 1.0 \quad (5.17)$$

$$P(m_{i_{E_i}} \mid l_{j_{E_i}}) = 1.0 \quad (5.18)$$

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 as explained in Chapter 4, Section 4.2.

Next, we explain *EventIdentityInfoRank*.

5.7.2 EventIdentityInfoRank

EventIdentityInfoRank is an iterative algorithm that takes into account the mutually reinforcing relationships between the vertices of *EventIdentityInfoGraph* as explained in the previous section and propagates event-specific scores of each vertex to connected vertices across the graph for ranking its vertices ($\in V_{E_i(t)}$) 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 5.19-5.22. The tweets ($\in M_{E_i}$) are assigned an initial informativeness score as obtained from the logistic regression model in ‘Event Information Quality’ component. 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 *EventIdentityInfoGraph*. We aim to refine the initial scores and assign a final score for ranking the vertices by leveraging the mutually reinforcing relationships between them.

$$Score(h_{i_{E_i}}) = \frac{freq(h_i)}{\max\{freq(h_1), freq(h_2), \dots, freq(h_p)\}} \quad (5.19)$$

$$Score(w_{i_{E_i}}) = \frac{freq(w_i)}{\max\{freq(w_1), freq(w_2), \dots, freq(w_r)\}} \quad (5.20)$$

$$Score(u_{i_{E_i}}) = \frac{followers(u_i)}{\max\{followers(u_1), \dots, followers(u_r)\}} \quad (5.21)$$

$$Score(l_{i_{E_i}}) = \frac{freq(l_i)}{\max\{freq(l_1), freq(l_2), \dots, freq(l_r)\}} \quad (5.22)$$

The relationships between two different subsets of vertices in graph $\mathbf{G}_{E_i(t)}$ is denoted by an affinity matrix. For e.g., $\mathbf{A}_{E_i}^{MH}$ denotes the $\mathbf{M}_{E_i} - \mathbf{H}_{E_i}$ affinity matrix for event E_i , where $(i,j)^{th}$ 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 equations 5.1-5.18. Similarly, $\mathbf{A}_{E_i}^{WH}$ denotes the $\mathbf{W}_{E_i} - \mathbf{H}_{E_i}$ 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 ($\in H_{E_i}$), set of text units ($\in W_{E_i}$), set of users ($\in U_{E_i}$), and set of URLs ($\in L_{E_i}$), respectively. Therefore, the Mutual Reinforcement Chain ranking for the k^{th} iteration can be formulated as follows:

$$R_{E_i}^{M(k+1)} = A_{E_i}^{MM(k)} R_{E_i}^{M(k)} + A_{E_i}^{MH(k)} R_{E_i}^{H(k)} + A_{E_i}^{MW(k)} R_{E_i}^{W(k)} + A_{E_i}^{MU(k)} R_{E_i}^{U(k)} + A_{E_i}^{ML(k)} R_{E_i}^{L(k)} \quad (5.23)$$

$$R_{E_i}^{H(k+1)} = A_{E_i}^{HM(k)} R_{E_i}^{M(k)} + A_{E_i}^{HH(k)} R_{E_i}^{H(k)} + A_{E_i}^{HW(k)} R_{E_i}^{W(k)} + A_{E_i}^{HU(k)} R_{E_i}^{U(k)} + A_{E_i}^{HL(k)} R_{E_i}^{L(k)} \quad (5.24)$$

$$R_{E_i}^{W(k+1)} = A_{E_i}^{WM(k)} R_{E_i}^{M(k)} + A_{E_i}^{WH(k)} R_{E_i}^{H(k)} + A_{E_i}^{WW(k)} R_{E_i}^{W(k)} + A_{E_i}^{UU(k)} R_{E_i}^{U(k)} + A_{E_i}^{WL(k)} R_{E_i}^{L(k)}$$
(5.25)

$$R_{E_i}^{U(k+1)} = A_{E_i}^{UM(k)} R_{E_i}^{M(k)} + A_{E_i}^{UH(k)} R_{E_i}^{H(k)} + A_{E_i}^{UW(k)} R_{E_i}^{W(k)} + A_{E_i}^{UU(k)} R_{E_i}^{U(k)} + A_{E_i}^{UL(k)} R_{E_i}^{L(k)}$$
(5.26)

$$R_{E_i}^{L(k+1)} = A_{E_i}^{LM(k)} R_{E_i}^{M(k)} + A_{E_i}^{LH(k)} R_{E_i}^{H(k)} + A_{E_i}^{LW(k)} R_{E_i}^{W(k)} + A_{E_i}^{LU(k)} R_{E_i}^{U(k)} + A_{E_i}^{LL(k)} R_{E_i}^{L(k)}$$
(5.27)

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

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

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 Δ_{E_i} .

$$\Delta_{E_i} \cdot R_{E_i} = \lambda \cdot R_{E_i} \quad (5.28)$$

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

To make Δ_{E_i} stochastic we divide the value of each element in a column of Δ_{E_i} by the sum of the values of all the elements in that column. This finally makes Δ_{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

$$\bar{\Delta}_{E_i} = \alpha \hat{\Delta}_{E_i} + (1 - \alpha) E \quad (5.29)$$

$$E = p \times [1]_{1 \times k} \quad (5.30)$$

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, $\overline{\Delta}_{E_i}$ is stochastic and irreducible and it can be shown that it is also primitive by checking $\overline{\Delta}_{E_i}^2$ 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

$$(R_{E_i}^M, R_{E_i}^H, R_{E_i}^W, R_{E_i}^U, R_{E_i}^L)$$

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. The entire procedure is presented step-by-step in an Algorithm 1.

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

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 Δ_{E_i} ;

Make matrix Δ_{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 $\hat{M}_{E_i}, \hat{H}_{E_i}, \hat{W}_{E_i}, \hat{U}_{E_i}, \hat{L}_{E_i}$;

Algorithm 1: EventIdentityInfoRank algorithm

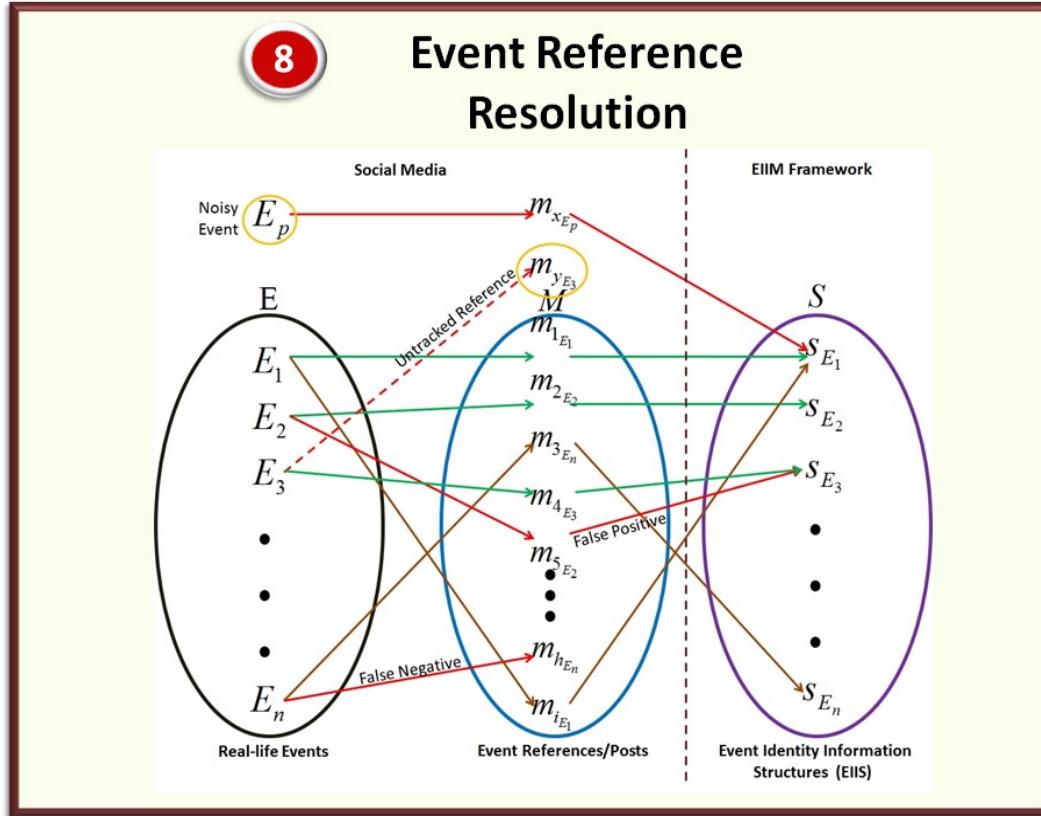
Since the algorithm uses power iteration method for ranking the vertices of a graph, it could be easily made scalable using mapreduce paradigm [131]. We plan to work on it in the future and implement our framework using hadoop and mapreduce environment. Also, the EIIM framework takes a hybrid approach 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.

5.8 Event Reference Resolution

One of the main aims of the EIIM framework is to assign the reference of an event E_i to its corresponding Entity Identity Information Structure (s_{E_i}). This aim is executed by this component. The main functions are as follows:

- The output of the previous component assigns a final ranked event-specific informativeness scores to the tweets. Using this score the component allows to choose top k tweets ranked in

FIGURE 5.15: Event Reference Resolution component of the EIIM life cycle.



terms of their event-specific informative content. This enables the identification and resolution of high quality tweets providing useful information about the event E_i , solving the problems of *noisy events* and *noisy references* as already discussed in Chapter 2, Section 2.5. The tweets after the top K are thrown away from the EIIS. This results in extremely high quality, event-specific informative tweets in the EIIS. We performed our experiment on the two events:

- Sydney Siege Crisis
- Millions March NYC

Excerpts of the top 5 tweets identified for both the events are given below.

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.c...>
3. RT @PatDollard: #SydneySiege: Hostages Held By Jihadis In Australian Cafe - WATCH LIVE VIDEO COVERAGE <http://t.co/uGxmd7zLpc> #tcot #pjnet sydney-siege-scene/index.html
4. RT @FoxNews: MORE: Police confirm 3 hostages escape Sydney cafe, unknown number remain inside [#SydneySiege](http://t.co/pcAt91LIdS)
5. Watch #sydneySiege police conference live as hostages are still being held inside a central Sydney cafe [#c4news](http://t.co/OjulBqM7w2)

- The other task that can be achieved using this component is the extraction of extremely informative features from the ranked results of the previous component and use them to form an evolutionary classifier or a feature vector for constantly tracking the event tweets w.r.t time. But the feature may get updated as the event progresses and a new set of ranked features needs to be obtained from the previous component after an interval of time that should be configurable. This, functionality of the component is currently not implemented and tested. We consider it as one of our future works.

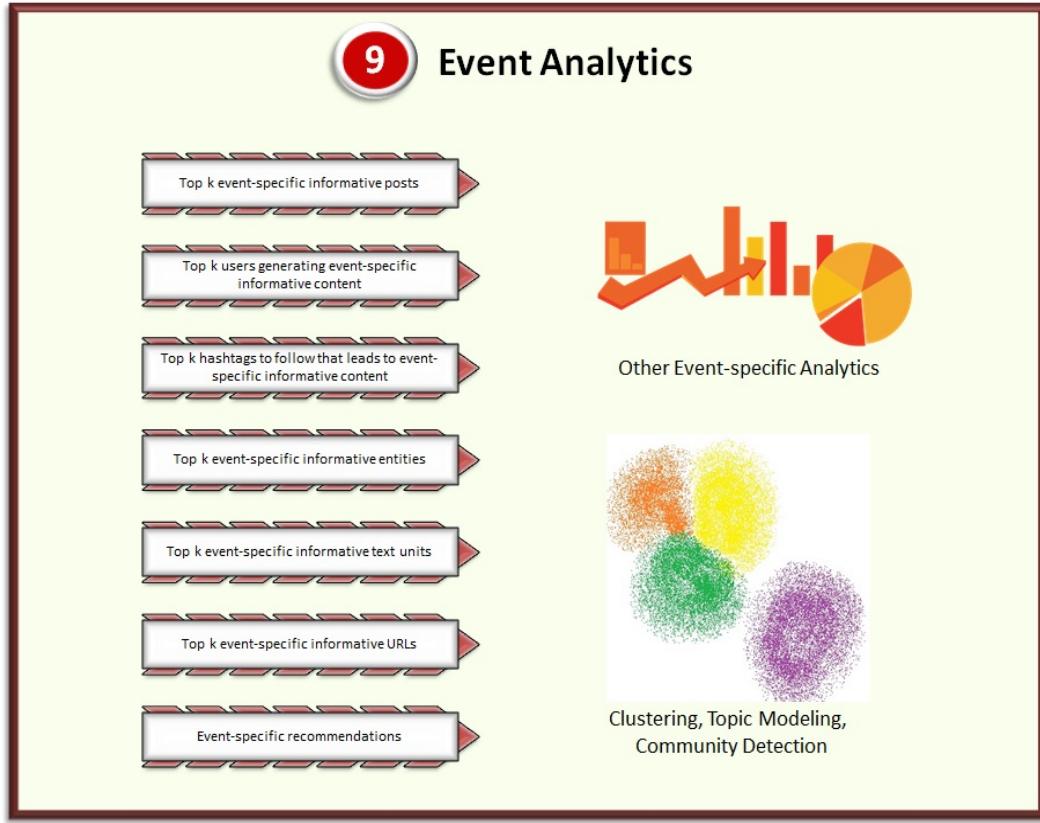
5.9 Event Analytics

The outputs of the Event Identity Information Processing component after processing the EventIdentityInfoGraph of a particular event is used for generating different analytics related to the event for the chosen time period, by this component. Some of the immediately available analytics that help in extracting deeper insights from the event related content are shown below for our sample events.

Top Five Event-specific Informative Hashtags for Sydney Siege Event

1. #sydneysiege

FIGURE 5.16: Event Analytics component of the EIIM life cycle.



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>

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/a2vgrn30X> #sydneysiege

- (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 #sydneySiege

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> #sydneySiege
- (c) RT @hughwhitfeld: British PM David Cameron informed of #sydneySiege ..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
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
5. <https://www.youtube.com/watch?v=Iz7hkfNmftY&feature=youtu.be>

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/WkttxssAfDp>

- (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> [#MillionsMarchNYC](http://t.co/X4QZrfOISh)
- (b) RT @NYjusticeleague: Hundreds killed. Ten Demands. One Continued Fight. Sign our petition at: [#MillionsMarchNYC](http://t.co/KETNo6bS0V) [htt...](http://t.co/ht...)
- (c) RT @cobismith: Union Square now with NYPD in foreground, #MillionsMarchNYC protesters at right and; US national debt ticker on the left [http:/...](http://t/...)

3. User 3

- (a) RT @mashable: Timelapse video reveals massive size of New York City protests [#MillionsMarchNYC](http://t.co/zhqHpkDLk1) <http://t.co/WktxssAfDp>

- (b) RT @KeeganNYC: LOTS of NYPD waiting for protesters on the BK side of the Brooklyn Bridge #MillionsMarch-NYC #ShutItDown #ICantBreathe <http://...>
- (c) RT @Zegota42: . @KeeganNYC Protesters on Brooklyn Bridge leaving Manhattan Skyline behind. #MillionsMarch-NYC #ICantBreathe <http://t.co/UPvN...>

Some of the other event analytics that can be readily done using the *EventIdentityInfoGraph*, and we would like to explore in the near future are:

- Topic modeling.
- Event-specific recommendations.
- Clustering and Cluster Analysis.
- Community detection.
- Trend Analysis.

The event related data stored in the database by the ‘Event Reference Collection’ would allow to do all types of analysis that are popular in social media.

Chapter 6

Evaluations

6.1 Evaluation Baselines

In order to evaluate the performance of *EventIdentityInfoRank* we selected six different techniques that acted as our baselines. The six techniques along with their brief explanation and the reason behind their choice is discussed below.

1. **LexRank** - LexRank is a popular graph based algorithm commonly used in summarization of textual documents [56]. It uses a stochastic graph-based method for computing the relative importance of textual units for Natural Language Processing. The task of extractive text summarization is based on the concept of identifying the most important sentences in a document or a set of documents. Importance is defined in terms of the presence of particular important words or in terms of similarity to

a centroid pseudo-sentence. LexRank, computes sentence importance based on the concept of eigenvector centrality in a graph representation of sentences. In this model, a connectivity matrix based on intra-sentence cosine similarity is used as the adjacency matrix of the graph representation of sentences. We implement the LexRank algorithm using an open-source python module named sumy¹, and rank the event related tweets considering them as individual sentences.

The objective of ranking natural language sentences in terms of their importance, makes it very similar to the *EventIdentityInfoRank* algorithm as proposed in this dissertation for the purpose of ranking tweets instead of textual documents. *EventIdentityInfoRank*, additionally ranks other units of information such as hashtags, text units, users and URLs, simultaneously, and does not takes into account the similarity of the tweets with a centroid pseudo-tweet.

2. **TextRank** - TextRank is another popularly used technique used for summarization of textual documents [132]. It is also a graph-based ranking model for text processing, that can be successfully used in natural language applications. The mechanism of its working is very similar to PageRank. However, instead of ranking web pages based on their linking structure,

¹<https://pypi.python.org/pypi/sumy/0.1.0>

it ranks text units based on their linking structure. It can be used for identifying salient sentences as well as key words of a document. Its objective of identifying key words and important sentences is also similar to our objective of finding important tweets.

In our implementation, we modified the algorithm in order to make it suitable for our context. Apart from creating heterogeneous relationships in *EventIdentityInfoGraph* we also created homogeneous relationships between the *event identity information units*. Cosine similarity (≥ 0.10) was used as the measure of relatedness between tweets, and the association scores of the hashtags, text units, users and URLs were based on their co-occurrence normalized between 0 and 1. The users were associated whenever they mentioned each other in the tweets, and the association score was measured by the number of mentions normalized between 0 and 1.

3. **Centroid** - Centroid is one of the techniques that was previously used in the literature for solving a part of our problem that ranks tweets. The technique is used for identifying high quality informative and useful tweets related to an event [52]. In order to implement it as a baseline we considered the tweets for the event in the given time period as one cluster. After pre-processing the tweets, we calculated the centroid of the cluster

and ordered the tweets in the decreasing order of their similarities with the centroid.

4. **SeenRank** - SeenRank is a proprietary algorithm commercially used by Seen.co for generating event summaries and highlights from Twitter. We considered *SeenRank* as the state-of-the-art technique. Although the true working of the algorithm is unknown, yet the task that the algorithm achieves is similar to the task of *EventIdentityInfoRank*. In order to use this technique as one of our baselines. we also collected the tweets about the events tracked by our framework from Seen.co. The collection task was achieved using their API found at <http://developer.seen.co/>, and a freely available python wrapper pySeen², for collecting data from Seen.co. Each tweet collected from the website has a score assigned to it using SeenRank. We use this score for arranging the tweets in descending order. The ordering of the tweets were confirmed from the company's co-founder in order to be sure that greater score reflects higher ranking.

5. **RTRank** - Number of retweets is a good measure of popularity of a tweet and is also used by Twitter for ranking its search results. It is also commonly used by other platforms for ranking tweets as already pointed out in Chapter 3. Therefore, we

²<https://github.com/dxmahata/pySeen>

also considered tweets ordered in decreasing order of number of retweets as one of our baselines. We name this scheme as RTRank.

6. Logistic Regression Model - This technique is the logistic regression model that we implemented for initializing the informativeness score of the tweets in the ‘Event Information Quality’ component. The generic informativeness score assigned by the logistic regression model is different from the final event-specific informativeness score assigned by *EventIdentityInfoRank*. Also the logistic regression model acts as a good representative of supervised approaches. We explain this with an example.

On manually analyzing the informative tweets we tried to assess if it is good enough to train a classifier for detecting informative tweets for an event in order to identify valuable event-specific information. Although the tweets on which we trained our logistic regression model were related to events yet we came across tweets like, *RT @BFDealz: http://t.co/TSJAigrVJI WHEELS SUPER TREASURE HUNT SUPERIZED HARLEY DAVIDSON FAT BOY LONG CARD 2014 #cpac2014 #sxsw*, which were classified as informative, even when it did not contain any event-specific information.

This was probably because of the choice of features for the model, which were generic and not event-specific. The model

did not take into account the presence of features that were popular and specific to the events, like popular hashtags, text units, etc. Popularity alone might not work as it is often mis-used by the spammers. It is also challenging to come up with a list of such event-specific features. Moreover, if one can compile such a list then it would be difficult to set thresholds on each such feature in order to qualify it as event-specific. Also, a supervised classification model does not have the ability to simultaneously rank tweets, hashtags, text units, URLs and users in terms of event-specific informativeness. After going through the existing literature we assume that the challenges discussed above would be a shortcoming of any supervised model and there is a need for an alternative feasible approach. It is also difficult to predict the event-specific informativeness in the URLs shared along with the tweets, as it might be necessary to analyze the content pointed to by the URLs. Also, not all the URLs contain text. They might be images or videos providing valuable information about an event. This motivated us to devise a novel framework that solves all the above problems.

Therefore, we considered the model as one of our baselines in order to make sure that our *EventIdentityInfoRank* improves upon the initial generic informativeness score already assigned

to the tweets at the start of the iteration and assigns event-specific informativeness scores on convergence. In other words the tweets having high score after the final ranking are more useful and informative than the initial ranking obtained using the logistic regression model.

Due to unavailability of proper baseline techniques for ranking hashtags, text units, URLs and users in terms of event-specific informativeness we do not compare the results obtained for them with any other approach. However, we report their average scores and sample results. Please refer the previous chapter for the sample results.

6.2 Evaluation Setup and Objectives

We evaluated the rankings obtained using *EventIdentityInfoRank* on the datasets (refer Chapter 5, ‘Event Reference Collection’ component), collected for events: “Millions March NYC” and “Sydney Siege Crisis”, by comparing its performance with the selected baselines. A subset of tweets for each event for a given time period (one hour) was selected. The choice of the time period was made on the basis of the intersection of the time period of the tweets collected by us and that provided by Seen for the same event. There were 21641 tweets for Millions March NYC and 37429 tweets for Sydney Siege, respectively. We obtained the ranked tweets for all the seven

approaches. For all the approaches except *SeenRank* the tweets were sorted in decreasing order on the basis of the ranking scores as the primary key and time of posting as the secondary key. This was done in order to get the most informative yet recent tweets at the top of the order. For *SeenRank* we sorted the tweets in terms of the scores assigned to them by Seen, as showing recent informative tweets for an event is one of the features of their platform.

We then followed a standard user evaluation approach to judge the event-specific informativeness of ranked tweets and also the hashtags, text units, URLs, and users. A team of three independent annotators comprising of graduate students, having taken the course of Information Retrieval, were assigned the task of annotation. Necessary background of the events were given to the annotators along with suitable resources for learning more about the events. Next, we present the annotation schemes.

TABLE 6.1: Avg IIC scores and total avg scores of annotations for Millions March NYC event.

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

TABLE 6.2: 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 |

6.2.1 Tweet Annotation

The ranked tweets were annotated on an event-specific informativeness-scale of 1 to 3 by the three independent annotators. We provide sample tweets for each of them taking the Sydney Siege event as our example.

- The value of 1 was assigned to tweets that does not contain any event related information (for e.g. *SteveSmith becomes Australias 45th Test captain http://t.co/nYh9DqRXxh #sydneyseige #MartinPlace Lindt #MYEFO #sieg Ray Hadley Mus-lims ISIS*).
- Value of 2 was assigned to tweets that were related to the event yet they did not provide useful event-specific information (for

e.g. *RT @TheDavidStevens: It wasn't just the policeman grabbing that girl in his arms, it was every Australian watching on too #sydneyseige*).

- A value of 3 was assigned to tweets that not only provided useful event-specific informative content but also led the user to more detailed information following the URLs mentioned in the tweet (for e.g. *RT @FoxNews: MORE: Police confirm 3 hostages escape Sydney cafe, unknown number remain inside <http://t.co/pcAt91LIdS> #Sydneyseige*).

The annotators assigned scores to top 100 tweets ranked according to each of the seven strategies. Thereafter, we computed *Inter Indexer Consistency* (IIC) values [133] for the annotations of the two datasets. The average IIC scores obtained for the two events are shown in Table 6.1 and Table 6.2, respectively. The IIC values for both the events fall in the acceptable range of accuracy of annotations. A tweet might be assigned three different scores by the annotators. In that scenario we find the average of the three scores and round it off to the smallest positive integer and assign a single score to each tweet. We also report the total average scores for top 100 tweets for both the events in the tables.

6.2.2 Hashtags, Text Units and URL Annotations

A similar annotation strategy was taken for annotating the top 50 hashtags, text units and URLs obtained using EventIdentityInfoRank. For hashtags and text units the annotators were asked to look at the tweets that consisted them. Following strategy was followed for scoring.

- If the tweets containing them primarily led to event-specific informative content then a score of 3 was assigned.
- If the tweets containing them led to related but not so informative content about the event then they were assigned a score of 2.
- Hashtags and text units that were irrelevant and did not lead to any event related content, were assigned a score of 1.

Similarly, the annotators visited the links for each URL, and based on the content they assigned them a score between 1-3. If the URLs were videos and images, then they further visited the tweet containing them in order to understand the context and scored them accordingly. Table 6.1 and Table 6.2 shows their average IIC scores and total average scores for top 50 ranks.

6.2.3 User Annotations

For annotating users we selected 5 random tweets for each of the top 50 users ranked according to *EventIdentityInfoRank*. An user was assigned a score of 3 if more than three of his tweets out of five got a score of 3 in the event-specific informativeness scale as already explained earlier. If three of his tweets get a score of 3 then the user gets a score of 2. Otherwise, a score of 1 is assigned to the user. Table 6.1 and Table 6.2 shows average IIC scores and total average scores for top 50 users.

6.2.4 NDCG@n and Precision@n

After being assured about consistency and accuracy of annotations, we moved to compute the *Normalized Discounted Cumulative Gain* (NDCG) [134] and Precision [135] values at each of the hundred recall levels. The NDCG values consider both the position and event-specific informativeness scores of the tweets. The NDCG value up-to position p in the ranking is given by equation 6.2, where DCG_p denotes the *discounted cumulative gain up-to position p* and is calculated using equation 6.1, and $IDCG_p$ denotes the *ideal discounted cumulative gain* value till position p in the ranking, or in other words the maximum possible DCG_p value till position p . rel_i denotes the graded relevance of the result at position i . In the context of our

evaluation rel_i represents the average rounded score in the scale of (1-3) that has been assigned by the annotators to the tweet at position i in the ranked list of top 100 tweets.

$$DCG_p = \sum_{i=1}^p \frac{2^{rel_i} - 1}{\log(i + 1)} \quad (6.1)$$

$$nDCG_p = \frac{DCG_p}{IDCG_p} \quad (6.2)$$

Precision@n is measured using equation 6.3. A tweet was considered to be relevant if it has a score of either 3 or 2 and was considered irrelevant if it has a score of 1.

$$\frac{\text{No. of relevant tweets at position } n}{n} \quad (6.3)$$

NDCG@n and Precision@n values were calculated for all the seven approaches for each of the datasets. Figures 6.1 and 6.2 shows the NDCG curves for all the seven approaches on the Millions March NYC and the Sydney Siege events, respectively, for up-to 20 recall levels. Tables 6.4 and 6.6 presents the NDCG@n values and Precision@n values for different recall levels upto 100 for the Sydney Siege Crisis event. Similarly, Tables 6.3 and 6.5 presents the NDCG@n values and Precision@n values for different recall levels upto 100

for the Millions March Nyc event. It is quite evident from the figures and the tables that EventIdentityInfoRank approach outperforms all the baselines including the state-of-the-art approach of *SeenRank* in gaining event-specific information.

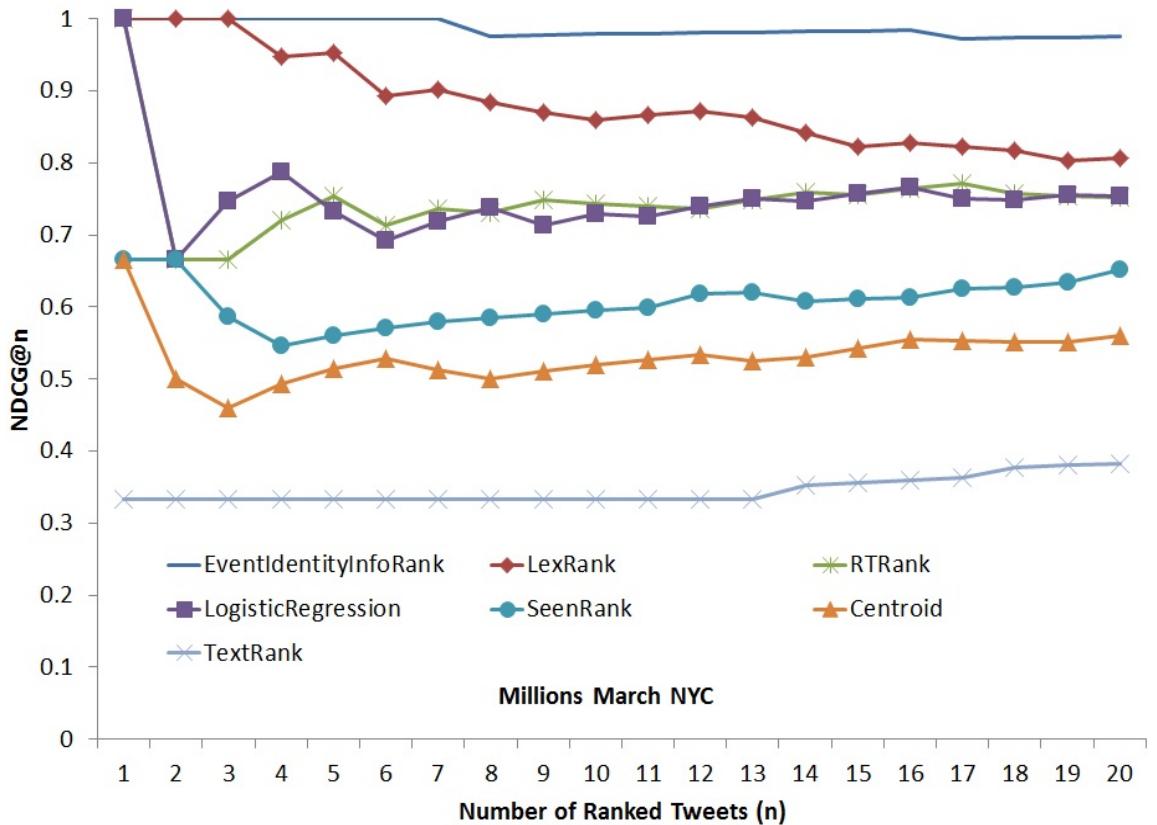


FIGURE 6.1: Performance comparison of ranking techniques using NDCG scores.

On considering only the top 10 tweets we observed a substantial information gain of our algorithm over the state-of-the-art (*SeenRank*) and the baseline that performed second best for both the events. On comparing the values of NDCG@10 for the two events we found that our algorithm performs 13.96% (Millions March NYC)

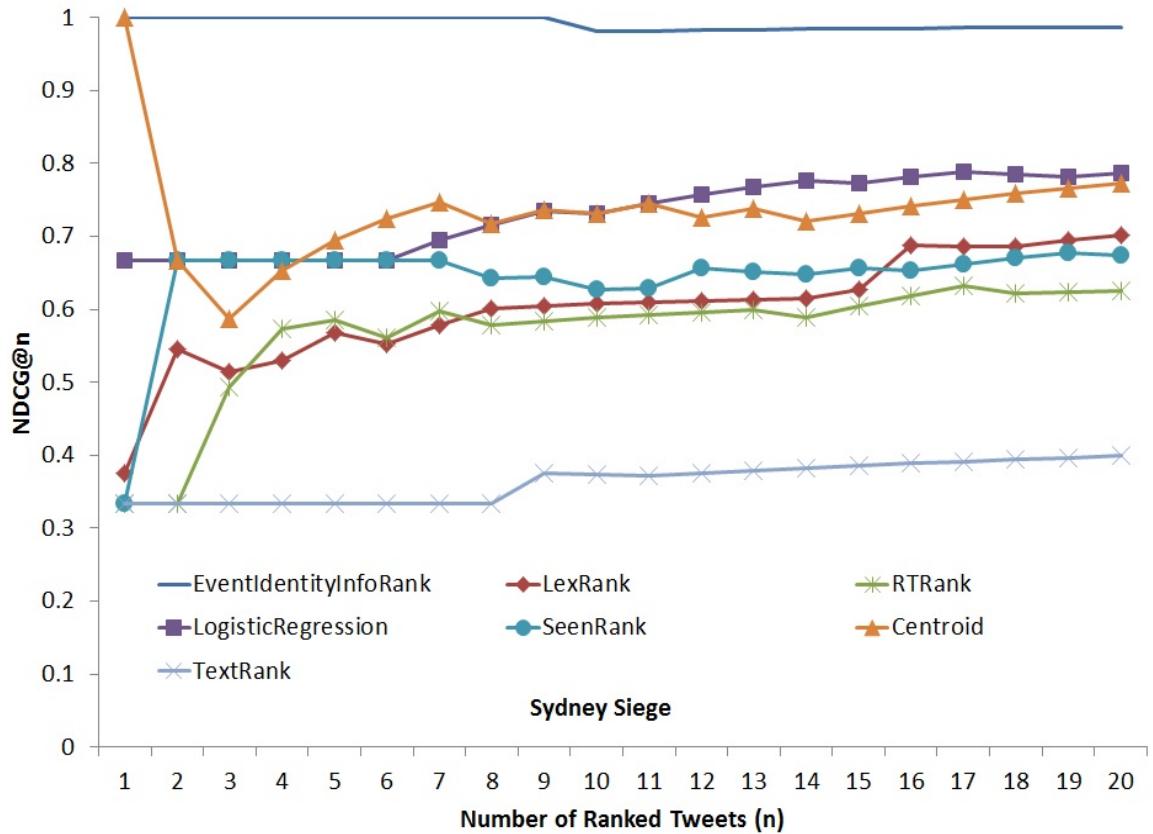


FIGURE 6.2: Performance comparison of ranking techniques using NDCG scores.

| Technique | @ 10 | @ 20 | @ 30 | @ 40 | @ 50 | @ 60 | @ 70 | @ 80 | @ 90 | @ 100 |
|-----------------------|-------|-------|-------|-------|-------|-------|-------|-------|-------|-------|
| EventIdentityInfoRank | 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 6.3: Performance comparison of ranking techniques using NDCG scores.

| Technique | @ 10 | @ 20 | @ 30 | @ 40 | @ 50 | @ 60 | @ 70 | @ 80 | @ 90 | @ 100 |
|-----------------------|---------|---------|---------|---------|---------|---------|---------|---------|---------|----------|
| EventIdentityInfoRank | 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 6.4: Performance comparison of ranking techniques using NDCG scores.

| Technique | @ 10 | @ 20 | @ 30 | @ 40 | @ 50 | @ 60 | @ 70 | @ 80 | @ 90 | @ 100 |
|-----------------------|---------|---------|---------|---------|---------|---------|---------|---------|---------|----------|
| EventIdentityInfoRank | 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 6.5: Performance comparison of ranking techniques using precision scores.

| Technique | @ 10 | @ 20 | @ 30 | @ 40 | @ 50 | @ 60 | @ 70 | @ 80 | @ 90 | @ 100 |
|------------------------------|---------|---------|---------|---------|---------|---------|---------|---------|---------|----------|
| EventIdentityInfoRank | 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 6.6: Performance comparison of ranking techniques using precision scores.

and 34.07% (Sydney Siege) better than the second best baseline technique, in identifying event-specific informative tweets. When compared with *SeenRank*, our algorithm was 64.53% (Millions March NYC) and 56.59% (Sydney Siege) better.

We also reasoned about the poor performance of *TextRank* in both the events. Since *TextRank* allowed random walks between homogeneous nodes, the strong association of non-informative nodes with the informative ones might have lowered the final scores of the informative nodes. The strong association of non-informative nodes with informative ones can be attributed to the spamming activity as already explained earlier Chapter 3, Section 4.2 This also proves that our framework is robust against spams and is very effective in identifying the most informative content related to events from the noisy stream of tweets in Twitter.

Chapter 7

Potential Applications of the EIIM Framework

As users share more and more event related information in social media, there is a growing need of identifying, organizing and analyzing these content streams. One of the most important challenges for both automated and user facing systems is to glean event-specific high quality information out of the humongous volumes of social media data. The design and implementation of the EIIM framework as presented in this dissertation can be integrated as the core engine for different applications. Various applications of the EIIM framework is discussed next.

7.1 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’ [136] to international sporting events like the ‘Winter Olympics’ [137] to socio-political [5] and socio-economical [138] events that shook the world such as presidential elections [21], ‘Egyptian Revolution’ [139], 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 Geofeedia¹, TwitterStand², Twitris³, Truthy⁴, and Tweet-Tracker⁵ 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 [140]. Social media is being used for tracking terrorism activities [141], collective actions [13], and countering cyber-attack threats⁶. 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

¹<http://geofeedia.com/>

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

³<http://twitris.knoesis.org/>

⁴<http://Truthy.indiana.edu/>

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

⁶<https://www.recordedfuture.com/>

identifying, tracking and analyzing events and its related references in an organized manner over time.

7.2 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 Seen⁷. 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 benifitted 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 dissertation can act as the core component of information retrieval systems retrieving and organizing information related to real-life events from social media.

⁷<http://seen.co>

7.3 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 references 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 as discussed in the next chapter.

7.4 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 beneficial⁸. 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* along with *EventIdentityInfoRank* as proposed in this dissertation can aid in identifying event-specific informative content and users producing 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.

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

7.5 Social Media Data Integration

Organizations have increasingly started integrating the data available in social media with the enterprise data⁹. 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 [142] 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.

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

Chapter 8

Conclusion and Future Work

8.1 Conclusion

8.2 Future Work

8.2.1 Summarizing Event Related Content

Given the huge amount of content produced in social media related to real-life events, summarization of the content can be very useful in such a scenario. It can help the users to overcome the problem of information overload. One of the most important characteristics of the summarization techniques is to identify the most salient units of information from the textual posts. The *EventIdentityInfoGraph* and the *EventIdentityInfoRank* algorithm as proposed and implemented in this dissertation can be tuned for such a purpose. One of the steps that needs to be taken is to find how the salient event identity

information units obtained as an output of *EventIdentityInfoRank* can be used for constructing event summaries from short textual social media messages. We have started looking into this problem and look forward to solve it using the framework as proposed in this dissertation. One of the main advantages and novelty in solving this problem using the EIIM framework would be the capability of generating event summaries as the event evolves.

8.2.2 Identifying Insightful Opinionated Content Related to Events

Users often share insightful opinionated content about different topics, people, organizations in social media. Such content is also generated in the context of an event. For example, in a sporting event the fans may post a lot of opinionated content about the players. Not all of them will be insightful. Similarly, in a product launch event, the prospective customers, or the reviewers may post very insightful and opinionated reviews about the new product. This type of content is extremely useful for the prospective customers, targeted marketing and for automated systems in order to identify the positive and negative aspects of the product that is creating buzz in social media. Identification of such insightful opinionated tweets can lead to the discovery of very useful and strategic information. 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 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 #Heat-Nation #HeatvsCavs 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. #HeatvsCavs #NBAXmas

The above tweets point to the reactions of the viewers on the game as well as the players participating in the event. We plan to work on this and take steps to tune our framework in order to make it better than the state-of-the-art techniques, for identifying insightful opinionated content from social media. This useful content once identified and ranked can also be used for generating opinion summaries.

8.2.3 Event-specific Recommendations

The graph based data structure used for storing the EIIS can be used for generating event related informative recommendations in near real-time. The graph structure aids in exploring relationships

between tweets, text units, hashtags, users and URLs. Moreover, the ‘Event Identity Information Process’ component processes the EIIS and assigns event-specific informativeness scores to its vertices. These scores combined with the relationships between the vertices can be leveraged for recommending users to other users who are producing event-specific informative content. Similarly, event-specific informative tweets, URLs and hashtags can be recommended. A naive approach has been implemented. For example following is a refined tweet recommendation for an event obtained from a snapshot of the *EventIdentityInfoGraph* created for the event: “BlackLives-Matter”: Protest movement against the killing of Eric Garner.

Original Tweet:

- #BREAKING #NEWS — New York City Mayor Says, #Black-LivesMatter
<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.

8.2.4 Distributed Processing of EventIdentityInfoGraph

The *EventIdentityInfoRank* algorithm processes the nodes and edges of the *EventIdentityInfoGraph* iteratively to come up with a simultaneous ranking of its heterogeneous vertices. The processing of the heterogeneous nodes can be distributed and then an aggregate score can be assigned to each vertex after an individual iteration. This is perfectly suitable for implementing the algorithm in a mapreduce paradigm. Similar steps are taken by the PageRank algorithm for ranking billions of web pages at scale. We plan to use the Giraph¹

¹<http://giraph.apache.org/>

distributed graph processing library on top of HDFS for implementing the process of ranking the vertices of *EventIdentityInfoGraph*.

8.2.5 Event Ontology for Social Media

Another research direction than can be explored in the future is to develop an ontology for representing the extracted event identity information units. This will enable a systematic categorization of the event identity information units into different concepts that can aid in formal reasoning. Reasoning on the relationships and characteristics of individual event identity information units can lead to extraction of deeper insights. For example, questions like “who are the people involved?”, “what are the places mentioned in the event related content?”, “how are the different people and places related to one another?”, and so on. Attempts have already been made on formulating ontological representation of the multimedia content produced during events by Troncy et al. [143], as well as, representing social media communities [144]. Ontologies for integrating information from different types of documents have also been proposed [145], that can be used for representing the social media references and the relationships between the content extracted from them. Another ontology, which is of great interest to us is the Basic Formal Ontology (BFO) [146]. This is because of the fact that BFO is an upper level ontology and has constructs for all types of entities

including events. It will give us the freedom of exploring the way events can be defined in social media and the representation of its related textual content. Also, BFO considers events as separate from the other types of named entities like, person or a place. This enables reasoning about the relationships between several events with a person, and vice versa. At the same time relationships between the events can be explored.

Appendix A

Appendix Title Here

Write your Appendix content here.

Bibliography

- [1] Yinle Zhou and John Talburt. Entity identity information management (eiim). In *International Conference on Information Quality (ICIQ-11), Adelaide, Australia*, pages 327–341, 2011.
- [2] 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.
- [3] 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.
- [4] 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.

- [5] 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.
- [6] 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.
- [7] 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.
- [8] 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.
- [9] Hemant Purohit and Amit P Sheth. Twitris v3: From citizen sensing to analysis, coordination and action. In *ICWSM*, 2013.
- [10] Nitin Agarwal and Debanjan Mahata. Grouping the similar among the disconnected bloggers. *Social Media Mining and*

- Social Network Analysis: Emerging Research: Emerging Research*, page 54, 2013.
- [11] 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.
- [12] Fatih Sen, Rolf T Wigand, Nitin Agarwal, Debanjan Mahata, and Halil Bisgin. Identifying focal patterns in social networks. In *CASoN*, pages 105–108, 2012.
- [13] Nitin Agarwal, Merlyna Lim, and Rolf T Wigand. *Online collective action: Dynamics of the crowd in social media*. Springer, 2014.
- [14] Nitin Agarwal, Merlyna Lim, and Rolf T Wigand. Raising and rising voices in social media. *Business & Information Systems Engineering*, 4(3):113–126, 2012.
- [15] 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.
- [16] 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.
- [17] Andreas M Kaplan and Michael Haenlein. Users of the world, unite! the challenges and opportunities of social media. *Business horizons*, 53(1):59–68, 2010.
- [18] Nitin Agarwal and Yusuf Yiliyasi. Information quality challenges in social media. In *International Conference on Information Quality (ICIQ 2010), Little Rock, Arkansas*, 2010.
- [19] N. Hamdy et al. Framing the egyptian uprising in arabic language newspapers and social media. *Journal of Communication*, 2012.
- [20] Eytan Bakshy, Itamar Rosenn, Cameron Marlow, and Lada Adamic. The role of social networks in information diffusion. In *Proceedings of the 21st international conference on World Wide Web*, pages 519–528. ACM, 2012.
- [21] Emily Metzgar and Albert Maruggi. Social media and the 2008 us presidential election. *Journal of New Communications Research*, 4(1):141–165, 2009.
- [22] Social media for events, 2014. URL <http://socialmediaforeventsebook.com/>.

- [23] Nina Eyrich, Monica L Padman, and Kaye D Sweetser. Pr practitioners' use of social media tools and communication technology. *Public relations review*, 34(4):412–414, 2008.
- [24] W Glynn Mangold and David J Faulds. Social media: The new hybrid element of the promotion mix. *Business horizons*, 52(4):357–365, 2009.
- [25] Sasa Petrovic, Miles Osborne, Richard McCreadie, Craig Macdonald, and Iadh Ounis. Can twitter replace newswire for breaking news? 2013.
- [26] How bin laden news spread on twitter, 2011. URL <http://www.cnn.com/2011/TECH/social.media/05/02/osama.bin.laden.twitter/>.
- [27] Theater shooting unfolds in real time on social media, 2012. URL <http://www.cnn.com/2012/07/20/tech/social-media/colorado-shooting-social-media/>.
- [28] Google and twitter launch service enabling egyptians to tweet by phone, 2011. URL <http://www.theguardian.com/technology/2011/feb/01/google-twitter-egypt>.
- [29] Edward N Zalta and Samson Abramsky. Stanford encyclopedia of philosophy, 2003.

- [30] Jeffrey M Zacks and Barbara Tversky. Event structure in perception and conception. *Psychological bulletin*, 127(1):3, 2001.
- [31] Hila Becker, Mor Naaman, and Luis Gravano. Beyond trending topics: Real-world event identification on twitter. *ICWSM*, 11:438–441, 2011.
- [32] John R Talburt. *Entity resolution and information quality*. Elsevier, 2011.
- [33] 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.
- [34] 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.
- [35] 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.

- [36] Sarita Yardi, Daniel Romero, Grant Schoenebeck, et al. Detecting spam in a twitter network. *First Monday*, 15(1), 2009.
- [37] 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.
- [38] 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.
- [39] 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.
- [40] Manish Gupta, Peixiang Zhao, and Jiawei Han. Evaluating event credibility on twitter. In *SDM*, pages 153–164. SIAM, 2012.
- [41] 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.

- [42] 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.
- [43] 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.
- [44] 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.
- [45] Lawrence Page, Sergey Brin, Rajeev Motwani, and Terry Winograd. The pagerank citation ranking: Bringing order to the web. 1999.
- [46] 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.

- [47] D Tunkelang. A twitter analog to pagerank. *The Noisy Channel*, 2009.
- [48] V Hallberg, A Hjalmarsson, J Puigcerver, C Rydberg, and J Stjernberg. An adaptation of the pagerank algorithm to twitter world. 2012.
- [49] 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.
- [50] 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.
- [51] 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.
- [52] Hila Becker, Mor Naaman, and Luis Gravano. Selecting quality twitter content for events. *ICWSM*, 11, 2011.

- [53] 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.
- [54] 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.
- [55] 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.
- [56] 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.
- [57] Rada Mihalcea and Paul Tarau. Textrank: Bringing order into texts. Association for Computational Linguistics, 2004.
- [58] Ivan P Fellegi and Alan B Sunter. A theory for record linkage. *Journal of the American Statistical Association*, 64(328):1183–1210, 1969.

- [59] Howard B Newcombe, James M Kennedy, SJ Axford, and Alison 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.
- [60] Thomas N Herzog, Fritz J Scheuren, and William E Winkler. *Data quality and record linkage techniques*. Springer Science & Business Media, 2007.
- [61] 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.
- [62] 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.
- [63] 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.
- [64] 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.

- [65] 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.
- [66] 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.
- [67] 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.
- [68] 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.
- [69] 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.
- [70] 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.
- [71] 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.
- [72] Jenny Rose Finkel. Named entity recognition and the stanford ner software, 2007.
- [73] Jason Baldridge. The opennlp project. *URL: <http://opennlp.apache.org/index.html>, (accessed 2 February 2012)*, 2005.
- [74] Hamish Cunningham. Gate, a general architecture for text engineering. *Computers and the Humanities*, 36(2):223–254, 2002.
- [75] Breck Baldwin and Bob Carpenter. Lingpipe. *Available from World Wide Web: <http://alias-i.com/lingpipe>*, 2003.
- [76] 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.

- [77] Jakub Piskorski and Roman Yangarber. Information extraction: Past, present and future. In *Multi-source, multilingual information extraction and summarization*, pages 23–49. Springer, 2013.
- [78] Sunita Sarawagi. Information extraction. *Foundations and trends in databases*, 1(3):261–377, 2008.
- [79] 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.
- [80] David Nadeau and Satoshi Sekine. A survey of named entity recognition and classification. *Lingvisticae Investigationes*, 30(1):3–26, 2007.
- [81] 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.
- [82] Sören Auer, Christian Bizer, Georgi Kobilarov, Jens Lehmann, Richard Cyganiak, and Zachary Ives. *Dbpedia: A nucleus for a web of open data*. Springer, 2007.

- [83] 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.
- [84] 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.
- [85] 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.
- [86] Lars Kolb, Andreas Thor, and Erhard Rahm. Dedoop: efficient deduplication with hadoop. *Proceedings of the VLDB Endowment*, 5(12):1878–1881, 2012.
- [87] John R Talburt and Yinle Zhou. *Entity Information Life Cycle for Big Data: Master Data Management and Information Integration*. Morgan Kaufmann, 2015.

- [88] 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.
- [89] Stefano Bortoli, Heiko Stoermer, Paolo Bouquet, and Holger Wache. Foaf-o-matic-solving the identity problem in the foaf network. In *SWAP*, 2007.
- [90] 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.
- [91] 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.
- [92] James Allan. *Topic detection and tracking: event-based information organization*, volume 12. Springer Science & Business Media, 2002.
- [93] Yiming Yang, Tom Pierce, and Jaime Carbonell. A study of retrospective and on-line event detection. In *Proceedings of the 21st annual international ACM SIGIR conference on Research and development in information retrieval*, pages 28–36. ACM, 1998.

- [94] 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.
- [95] Vasileios Hatzivassiloglou and Elena Filatova. Domain-independent detection, extraction, and labeling of atomic events. 2003.
- [96] Vasileios Hatzivassiloglou, Luis Gravano, and Ankineedu Manguanti. 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.
- [97] Juha Makkonen, Helena Ahonen-Myka, and Marko Salmenkivi. Simple semantics in topic detection and tracking. *Information Retrieval*, 7(3-4):347–368, 2004.
- [98] 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.

- [99] Heng Ji and Ralph Grishman. Refining event extraction through cross-document inference. In *ACL*, pages 254–262, 2008.
- [100] 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.
- [101] 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.
- [102] 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.
- [103] Jianshu Weng and Bu-Sung Lee. Event detection in twitter. *ICWSM*, 11:401–408, 2011.
- [104] 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.
- [105] 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.
- [106] 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.
- [107] 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.
- [108] 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.

- [109] 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.
- [110] Sarita Yardi and Danah Boyd. Tweeting from the town square: Measuring geographic local networks. In *ICWSM*, 2010.
- [111] 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.
- [112] 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.
- [113] Hila Becker, Feiyang Chen, Dan Iter, Mor Naaman, and Luis Gravano. Automatic identification and presentation of twitter content for planned events. In *ICWSM*, 2011.
- [114] Paul Hemp. Death by information overload. *Harvard business review*, 87(9):82–89, 2009.

- [115] Daniel Gayo-Avello. Nepotistic relationships in twitter and their impact on rank prestige algorithms. *Information Processing & Management*, 49(6):1250–1280, 2013.
- [116] Kalina Bontcheva, Leon Derczynski, Adam Funk, Mark A Greenwood, Diana Maynard, and Niraj Aswani. Twitie: An open-source information extraction pipeline for microblog text. In *RANLP*, pages 83–90, 2013.
- [117] Scott Prevost. An information structural approach to spoken language generation. In *Proceedings of the 34th annual meeting on Association for Computational Linguistics*, pages 294–301. Association for Computational Linguistics, 1996.
- [118] L.A. Adamic et al. Power-law distribution of the world wide web. *Science*, 287(5461):2115–2115, 2000.
- [119] LOmariba. Is new media posing a serious challenge to traditional media?”. Technical report, University of Westminster, 2009.
- [120] Z. Harb. Arab revolutions and the social media effect. *M/C Journal*, 14(2), 2011.
- [121] T.J. Johnson et al. Wag the blog: How reliance on traditional media and the internet influence credibility perceptions of weblogs among blog users. *Journalism & Mass Communication*

- Quarterly*, 81(3):622–642, 2004.
- [122] Nitin Agarwal, Merlyna Lim, and Rolf T Wigand. Finding her master’s voice: the power of collective action among female muslim bloggers. In *ECIS*, 2011.
- [123] C. Anderson. *Long Tail, The, Revised and Updated Edition: Why the Future of Business is Selling Less of More*. Hyperion, 2008.
- [124] S. Brin et al. The anatomy of a large-scale hypertextual web search engine. *Computer networks and ISDN systems*, 30(1):107–117, 1998.
- [125] Alexandra Olteanu, Carlos Castillo, Fernando Diaz, and Sarah Vieweg. Crisislex: A lexicon for collecting and filtering microblogged communications in crises. In *In Proceedings of the 8th International AAAI Conference on Weblogs and Social Media (ICWSM” 14)*, number EPFL-CONF-203561, 2014.
- [126] Minlie Huang, Yi Yang, and Xiaoyan Zhu. Quality-biased ranking of short texts in microblogging services. In *IJCNLP*, pages 373–382, 2011.
- [127] Mosquera Alejandro and Moreda Paloma. The use of metrics for measuring informality levels in web 2.0 texts. 2011.

- [128] David Laniado and Peter Mika. Making sense of twitter. In *The Semantic Web–ISWC 2010*, pages 470–485. Springer, 2010.
- [129] Genevieve Barrons. ’suleiman: Mubarak decided to step down# egypt# jan25 oh my god’: Examining the use of social media in the 2011 egyptian revolution. *Contemporary Arab Affairs*, 5(1):54–67, 2012.
- [130] 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.
- [131] 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.
- [132] R. Mihalcea et al. Textrank: Bringing order into texts. In *Proceedings of EMNLP*, volume 4, pages 404–411. Barcelona, Spain, 2004.
- [133] L Rolling. Indexing consistency, quality and efficiency. *Information Processing & Management*, 17(2):69–76, 1981.

- [134] Kalervo Järvelin and Jaana Kekäläinen. Cumulated gain-based evaluation of ir techniques. *ACM Transactions on Information Systems (TOIS)*, 20(4):422–446, 2002.
- [135] Ricardo Baeza-Yates, Berthier Ribeiro-Neto, et al. *Modern information retrieval*, volume 463. ACM press New York, 1999.
- [136] Huiji Gao, Geoffrey Barbier, Rebecca Goolsby, and Daniel Zeng. Harnessing the crowdsourcing power of social media for disaster relief. Technical report, DTIC Document, 2011.
- [137] Shaun Walker. Russia to monitor 'all communications' at winter olympics in sochi. *The Guardian, October*, 6, 2013.
- [138] Johan Bollen, Alberto Pepe, and Huina Mao. Modeling public mood and emotion: Twitter sentiment and socio-economic phenomena. *arXiv preprint arXiv:0911.1583*, 2009.
- [139] 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.
- [140] Jeffrey Ghannam. Social media in the arab world: Leading up to the uprisings of 2011. *Center for International Media Assistance*, 3, 2011.
- [141] 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.
- [142] Oren Tsur and Ari Rappoport. What’s in a hashtag?: content based prediction of the spread of ideas in microblogging communities. In *Proceedings of the fifth ACM international conference on Web search and data mining*, pages 643–652. ACM, 2012.
- [143] Raphaël Troncy et al. Linking events with media. In *Proceedings of the 6th International Conference on Semantic Systems*, page 42. ACM, 2010.
- [144] JG Breslin and U Bojars. Semantically-interlinked online communities.
- [145] Martin Doerr, Jane Hunter, and Carl Lagoze. Towards a core ontology for information integration. *Journal of Digital information*, 4(1), 2006.
- [146] Barry Smith and P Grenon. Basic formal ontology. *Draft. Downloadable at <http://ontology.buffalo.edu/bfo>*, 2002.