



## Industrial Management & Data Systems

Social information landscapes: Automated mapping of large multimodal, longitudinal social networks

Eugene Ch'ng

### Article information:

To cite this document:

Eugene Ch'ng , (2015),"Social information landscapes", Industrial Management & Data Systems, Vol. 115 Iss 9 pp. 1724 - 1751

Permanent link to this document:

<http://dx.doi.org/10.1108/IMDS-02-2015-0055>

Downloaded on: 02 November 2016, At: 21:37 (PT)

References: this document contains references to 21 other documents.

To copy this document: [permissions@emeraldinsight.com](mailto:permissions@emeraldinsight.com)

The fulltext of this document has been downloaded 283 times since 2015\*

### Users who downloaded this article also downloaded:

(2015),"Big Data promises value: is hardware technology taken onboard?", Industrial Management & Data Systems, Vol. 115 Iss 9 pp. 1577-1595 <http://dx.doi.org/10.1108/IMDS-04-2015-0160>

(2015),"Understanding community citizenship behavior in social networking sites: An extension of the social identification theory", Industrial Management & Data Systems, Vol. 115 Iss 9 pp. 1752-1772 <http://dx.doi.org/10.1108/IMDS-05-2015-0211>

Access to this document was granted through an Emerald subscription provided by emerald-srm:563821 []

### For Authors

If you would like to write for this, or any other Emerald publication, then please use our Emerald for Authors service information about how to choose which publication to write for and submission guidelines are available for all. Please visit [www.emeraldinsight.com/authors](http://www.emeraldinsight.com/authors) for more information.

### About Emerald [www.emeraldinsight.com](http://www.emeraldinsight.com)

Emerald is a global publisher linking research and practice to the benefit of society. The company manages a portfolio of more than 290 journals and over 2,350 books and book series volumes, as well as providing an extensive range of online products and additional customer resources and services.

Emerald is both COUNTER 4 and TRANSFER compliant. The organization is a partner of the Committee on Publication Ethics (COPE) and also works with Portico and the LOCKSS initiative for digital archive preservation.

\*Related content and download information correct at time of download.

## REGULAR PAPER

# Social information landscapes Automated mapping of large multimodal, longitudinal social networks

Eugene Ch'ng

*School of Computer Science,  
University of Nottingham Ningbo China, Ningbo, China*

### Abstract

**Purpose** – The purpose of this paper is to present a Big Data solution as a methodological approach to the automated collection, cleaning, collation, and mapping of multimodal, longitudinal data sets from social media. The paper constructs social information landscapes (SIL).

**Design/methodology/approach** – The research presented here adopts a Big Data methodological approach for mapping user-generated contents in social media. The methodology and algorithms presented are generic, and can be applied to diverse types of social media or user-generated contents involving user interactions, such as within blogs, comments in product pages, and other forms of media, so long as a formal data structure proposed here can be constructed.

**Findings** – The limited presentation of the sequential nature of content listings within social media and Web 2.0 pages, as viewed on web browsers or on mobile devices, do not necessarily reveal nor make obvious an unknown nature of the medium; that every participant, from content producers, to consumers, to followers and subscribers, including the contents they produce or subscribed to, are intrinsically connected in a hidden but massive network. Such networks when mapped, could be quantitatively analysed using social network analysis (e.g. centralities), and the semantics and sentiments could equally reveal valuable information with appropriate analytics. Yet that which is difficult is the traditional approach of collecting, cleaning, collating, and mapping such data sets into a sufficiently large sample of data that could yield important insights into the community structure and the directional, and polarity of interaction on diverse topics. This research solves this particular strand of problem.

**Research limitations/implications** – The automated mapping of extremely large networks involving hundreds of thousands to millions of nodes, encapsulating high resolution and contextual information, over a long period of time could possibly assist in the proving or even disproving of theories. The goal of this paper is to demonstrate the feasibility of using automated approaches for acquiring massive, connected data sets for academic inquiry in the social sciences.

**Practical implications** – The methods presented in this paper, together with the Big Data architecture can assist individuals and institutions with a limited budget, with practical approaches in constructing SIL. The software-hardware integrated architecture uses open source software, furthermore, the SIL mapping algorithms are easy to implement.

**Originality/value** – The majority of research in the literature uses traditional approaches for collecting social networks data. Traditional approaches can be slow and tedious; they do not yield adequate sample size to be of significant value for research. Whilst traditional approaches collect only a small percentage of data, the original methods presented here are able to collect and collate entire data sets in social media due to the automated and scalable mapping techniques.

**Keywords** Online communities, Big Data, Longitudinal network, Multimodal network, Social network analysis, Social information landscapes

**Paper type** Research paper



## 1. Introduction

What is social information landscapes (SIL) and how will it be useful for studying social networks? Here, SIL can be defined as “the automated mapping of large topological networks from instantaneous contents, sentiments and users reconstructed from social media channels, events and user generated contents within blogs and websites, presented virtually as a graph that encompasses, within a timescale, contextual information, all connections between followers, active users, comments and conversations within a social rather than a physical space”. The key phrase “automated mapping” is essential here as the mapping of very large networks is essential as network behaviour may differ in massive networks in relation to their emergent behaviour and the way they tend to self-organise. Larger networks are also ideal as subjects for studying geodesic distance, centrality, and density. SIL was first mentioned in an article dealing with a Big Data funnelling approach (Ch’ng, 2014), and subsequently in corresponding articles on the study of online community formation and decline when research results were published (Ch’ng, 2015b, c). In comparison to traditional networks, SILs carry much more information. SILs encapsulate activities, which complements the traditional follower-followee network format, which contains only human nodes as opposed to activity nodes that are multimodal, e.g., contents and context. As the mapping is novel, so therefore the “landscapes” created will be new in the context that they contain a higher resolution of, and broader context of information.

The onset of the internet age has made our world smaller. As offline communities connect to virtual communities, and become virtual communities, space and time are, in a sense compressed to within the social medium that facilitates community interaction. Insights into the behaviour of virtual communities in the age of social media requires a Big Data approach in mapping the interactions as social networks, for the interaction of online communities for a single topic can occur over large geographical distances and may span many months involving thousands to millions of participants from highly diverse demographics. Such interactions may also be multimodal, involving not only the users, but also the content of the interactions linked between multiple parties. As such, the traditional approach of manually mapping such networks can be tedious and would not necessarily yield a large enough sample for social network analysis. This is due to the fact that the collection, collation, and pre-processing of data sets are frequently too large to manage with manual data processing. This paper presents a Big Data solution as a methodological approach to the automated collection, cleaning, collation, and mapping of multimodal, longitudinal data sets from social media. The methods and algorithms presented here are generic, and can be applied to diverse types of social media or user-generated contents involving user interactions, such as within blogs, comments in product pages and other forms of media, so long as a formal data structure can be constructed.

The mapping of extremely large networks involving hundreds of thousands to millions of nodes could possibly prove or disprove theories in the social sciences. The goal of this paper is to demonstrate the feasibility of using automated approaches for acquiring massive, connected data sets for academic inquiry in the social sciences.

The paper begins with a background of the research. It continues with the methodology for mapping social networks, covering the need for a Big Data architecture, and identifying suitable asynchronous and distributable open source technologies that are scalable in terms of volume of data and velocity of incoming data streams. The latter section of the methodology presents the focus of this paper – the mapping of SIL. The section describes the data structures and algorithms for mapping,

reconfiguration and storage of dynamic and static networks. Section 4 demonstrates the feasibility of the methods by presenting real-world data sets. The paper concludes with a discussion and future work.

## 2. Background

The limited presentation of the sequential nature of content listings within social media, as viewed on Web browsers or on mobile devices, do not necessarily reveal nor make obvious an unknown nature of the medium; that every participant, from content producers, to consumers, to followers and subscribers, including the contents they produce or subscribed to, and the context they are in, are intrinsically connected in a hidden but massive network. Such networks when mapped, could be quantitatively analysed using methods in social network analysis (e.g. centralities), and the semantics and sentiments could equally reveal valuable information with appropriate analytics. Yet that which is difficult is the tediousness of collecting, cleaning, collating, and mapping such data sets into a data structure that could yield important insights into the community structure and the directional, and polarity of interaction on diverse topics. It seems that the majority of research that investigate social networks have not looked beyond the traditional approach of manually collecting and mapping such networks, to a Big Data approach where mapping processes are automated. The research community that invariably focused their attention on social networks should genuinely consider the fact that any online activity is a sample of a population, and that manual data collection, due to its tediousness, will perhaps harvest a tiny sample and, as a consequence, yield a minor percentage of the full online content – a sample of a sample of a population. Unfortunately, such data sets will invariably be biased. A Big Data approach using scalable technologies is a viable approach in the era of Big Data, using machines to collect machine-collated contents.

Massive quantities of information generated by people are being tapped by diverse groups within the domains of computer science, physics, mathematics, social sciences, business, and economics with the hope that it will answer questions in their disciplines. In the physics domain, for example, the Large Hadron Collider at CERN churns out vast quantities of information. Within computer science and the engineering disciplines, the development of connected devices defined as the Internet of Things generate massive information on human activities, systems and environments from sensors and video surveillance. Computational social science on the other hand offers “the capacity to collect and analyse data with an unprecedented breadth and depth and scale” (Lazer *et al.*, 2009). Regardless of which disciplines, the fact is that “these massive amount of information can be tracked and measured with unprecedented fidelity” (Anderson, 2008). Anderson continues, “This is a world where massive amounts of data and applied mathematics replace every other tool that might be brought to bear. Out with every theory of human behaviour, from linguistics to sociology. Forget taxonomy, ontology, and psychology. Who knows why people do what they do? The point is they do it, and we can track and measure it with unprecedented fidelity. With enough data, the numbers speak for themselves”. Whilst Anderson’s statement appeared to be controversial to many, there is truth in the fact that taking a Big Data approach could have significant contribution to a new instrument for the vehicle of knowledge – “change the instrument, and you will change the entire social theory that goes with them” (Latour, 2010, p. 9).

Data can be static, recorded in memory or on a physical medium and remain unchanged. Data can also be dynamic and continuous, with incoming streams from

various sources. Data can grow linearly, or exponentially, particularly with the viral nature of social media with volumes and velocity that will render a conventional computer useless. The nature of such data is in the realm of Big Data. Certain types of data can also “grow”, and evolve in a way that structural changes become different in short spans of time. These types of data are driven by biological or human activities and perhaps, machine-generated data, and are essentially complex adaptive systems. SILs belong to the type that grows exponentially and evolve over time.

Data can also be “Big” in a number of ways, and is observed by Manovich (2011) as data sets that are sufficiently large to require supercomputers. It is loosely defined by Jacobs as “data whose size forces us to look beyond the tried-and-true methods that are prevalent at that time” (Jacobs, 2009). However, Big Data is not only characterised by its size, but by its relationality with other data (Boyd and Crawford, 2011) – “Big Data is fundamentally networked”. Big Data can be characterised by three Vs – Volume, Variety, and Velocity. These Vs, together with four others are briefly cover here, and how they are related to the context of this paper – relationality and the networked nature of data.

Volume refers to the enormity of data that have been created from the past to the future in ever increasing rate. Velocity refers to the speed of data that is being created, stored and analysed. Variety refers to the number of types of raw data that is scattered, unstructured and complex. The “Big” in the term “Big Data” describes the three Vs, for all can potentially be “Big” to an extent that the data sets will render a standard computer insufficient for storing, processing and visualising incoming data in real-time (or near real-time), without which data becomes meaningless as the big picture will not be known. Four recent Vs have been proposed as part of the goal to define the properties of Big Data. The first is Veracity, which refers to the credibility, or accuracy of the data that a research may acquire or process, such as false information that are spread widely online via social media costing huge loss of resources (Bontcheva, 2014a, b). Variability refers to the meaning in which a piece of information may be associated with, and how that meaning may differ in different context or at different times. The third is Visualisation, which extracts patterns in data in a human-consumable format to extract meanings and inform data consumers. The previous Vs may eventually add Value to institutions and organisations within the digital economy when raw data is transformed via the Big Data process into information that facilitates decision making. The permutations of a set of research problems associated with Big Data’s Volume, Variety, Velocity, Veracity, Variability, Visualisation, and Value are manifold, and the mapping of SIL is one such area. Mapping SILs potentially faces challenging issues related to all these terms, especially during the process of analysis and interpretation. The present paper focuses on the two largest issues in Big Data – volume and velocity, this paves the way for resolving the other terms.

SIL of a similar nature may be obtained, but they are perhaps a limited subset of the definition. For example, Java *et al.* (2007) investigated a very small community network based on key terms within a gaming circle who also shared daily experiences with each other. Jansen *et al.* (2009) examine Starbucks’ customer networks with Starbucks as a central node based on the frequency of tweets in “word of mouth” communication. Huberman *et al.* (2008) reveals a hidden network of actual friends within a subset network via the removal of followers who were not part of an actual friendship. Sakaki *et al.* (2010) compared information diffusion in three Twitter network types: earthquake, typhoon, and new Nintendo game news. A larger study by Kwak *et al.* (2010) constructs retweets trees of “Air France Flights” and examined their temporal

and spatial characteristics for retweets pattern, depth, and speed. The mapping and analysis of large scale SIL may indeed provide deeper and broader context of information on a social network than conventional statistical analysis may allow. It may also reveal greater information as compared to previous studies related to traditional networks.

Social media research has attracted considerable attention in the academia, and has become highly important on the business sectors (Ngai *et al.*, 2015). The hidden social information within channels of social media is important for revealing insights on the spatial organisation of the social network. This requires the mapping of the “social space”. Bourdieu (1985) states that the social world, a “social topology”, can be represented as a space with several dimensions constructed on the basis of principles of differentiation or distribution constituted by the set of properties active within the social universe capable of conferring strength and power within the sphere of that universe. Agents and groups, according to Bourdieu, are defined by their relative positions within space as a set of objective power relations, i.e., field of forces as economic, cultural, social capital, or symbolic capital (prestige, reputation, etc.) that affect other agents who enter that space. Social networks are similar in nature. Whilst one may suggest that physical interactions are irreplaceable, the benefits that computer-based social networking sites have provided to society cannot be ignored. Social networks increase our range of human connectedness beyond the boundary of a participant’s geographical location. With regards to the time of interaction, communications sent now may be retrieved and responded to, much later by other participants. This invariably opens up a broad range of opportunities as space and time, in the eye of a user are “compressed” to within a digital display, allowing diverse communications from large demographics groups. In this context, the study of SIL is an important aspect of social science inquiries.

### 3. Methods for mapping SIL

Two basic types of social network analysis exist. The first is the ego network analysis, the second is complete network analysis. Ego network analysis is conducted via surveys where the “ego” responds to questions asked about the ego’s immediate network. Each ego is sampled from a large population in order to assess the person’s network size, diversity, etc. The second type of analysis views the complete networks, such as all existing relationships between a set of nodes. The research here is posited in the latter type, although data from large-scale surveys used for collecting ego networks can also be automatically mapped.

There are two types of social networks which can be mapped:

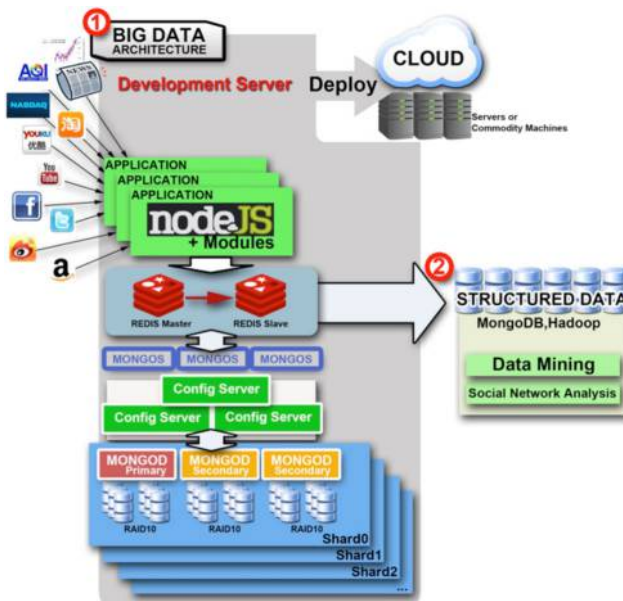
- (1) Follower – followee or “Content Provider” – “Content Consumer” network. This network consists mainly of inactive nodes with more or less permanent connections between them. Followers such as within Twitter, LinkedIn, or subscribers, and in YouTube and Instagram, for example, are posited within this type of network. The growth of these types of network is slower. Such networks have a larger temporal distance between followers and followees.
- (2) Activity networks consist of active participant interactions, which have a shorter temporal distance between each link. The growth of the network size can be rapid, particularly with viral contents. Examples of activity networks are active conversations within Twitter, Facebook, under the comments section of YouTube, Google+, and product pages or blogs.

Both network types are possible with the method described here. Section 3.1 begins the methodology by suggesting appropriate technology for such work that makes it possible to deal with the volume and velocity of data in the Big Data context.

3.1 Accessible big data software-hardware integrated architecture

This section introduces and suggests accessible technologies for Big Data scaling issues but do not delve into the details of their implementation. Implementation details are available within the technical manuals of the technologies and a summary guide in the same context is available (Ch'ng, 2014). For reasons of practicality, it is necessary to set two prerequisites prior to mapping social networks. The first of which is that such a system must be accessible to a wide range of individuals and institutions with low-budget constraints. The second requirement must ensure that all technologies used are scalable in terms of volume and velocity. When volume and velocity can be managed, data integrity, and completeness can be assured.

Fortunately, both requirements are achievable for individuals and institutions. Open source scalable technologies, packaged as application programming interfaces (API) and software libraries are available for the development and deployment on inexpensive hardware, i.e., commodity machines, simple server setup, etc. Alternately they may be deployed on Cloud services. Figure 1 is a summary diagram of the software architecture deployed on the author's Linux platform (and works with Windows and Mac OSes). In first, unstructured data enters through the server "edge nodes" via Node.js. Here, various algorithms work together to clean and structure data into schemas of key-value pairs.



**Notes:** (1) A Bid Data architecture integrating various open source scalable technologies; (2) structured data used for research. In the context of this paper, social network analysis, and data (graph) mining

**Figure 1.** A summary diagram illustrating the Big Data software-hardware integrated architecture

Data either enters through Redis (if real-time in-memory processing is needed), or directly into MongoDB, distributed across MongoDB Shards. The asynchronous nature of the “edge nodes” using Node.js ensures that data streams that scale in volume or velocity (viral contents) are all received. The software integration is similar across developmental servers, Cloud services, or within in-house servers or commodity machines. Commodity machines (low-budget clusters of PCs) are the single most inexpensive way for deploying such a system – each Node.js application uses only a single CPU core, thus the suitability of using recycled hardware. In second, data sets are queried through either Redis, or directly from MongoDB, via conventional algorithms that do not need high-efficiency algorithms as the velocity of data streams can now be controlled. Mapping algorithms can be programmed in either Node.js, or using any other languages that support connections with MongoDB or Redis.

Secondary data types are important to SILs. The availability of secondary data types is dependent on the social media API and the target web page. These statistical data attributes are usually available (friends count, number of subscribers, number of followers, number of posts, etc.) and could be used to compare with the growth and evolution of a network over time. Geo-location, time zones, language, gender, and textual contents provide other means of informing the status of each node in relation to the SILs. The research questions asked and the hypotheses tested within each research are different, and therefore, the collection of secondary data types is dependent on the project at hand.

*3.1.1 Accessing social media data and web crawling.* The volume and velocity of incoming data associated with mapping SILs may involve accessing social media API such as Twitter, Facebook, Weibo, QQ, YouTube or Youku, and online stores such as Amazon, JD.com or Taobao.com, for example. It is expected that some APIs may limit the availability of their repositories due to internal financial motive of the company at stake. Data after all is the new oil and has great value to business profitability. Where APIs limit data access, web-crawling and scraping algorithms can be used to extract data from JavaScript’s Document Object Model from web sites, which any graduate programmer from a reputable institution would have learned. This paper will not cover the details of accessing social media data via API, nor will it cover Web-scraping technologies as the topics are beyond the scope of this paper.

It is imperative that we use social media data for testing the methods presented here. The choice of social media should generate sufficiently large volumes of data with potentially high velocity. One of these types of social media is Twitter. Twitter generates around 500 million tweets per day. In any topical viral hash tags, the tweets in Twitter could easily reach tens of thousands of posts in a very short span of time. This meets the first two categorical Vs of Big Data – Volume and Velocity. As such, Twitter data will be used for testing the methods covered in this paper, using the streaming API, which captures more complete data sets than the limited REST API for polling endpoints.

*3.1.2 Asynchronous I/O.* The volume of data is expected to be in the orders of magnitude, the velocity will have a higher rate, considering that hundreds of millions of contents are generated every second. Such volume of data multiplied by the velocity of incoming streams from the Web will definitely crash conventional server-side applications. A scalable network layer that is event driven, with a non-blocking I/O model that is capable of data-intensive application with the ability to push and pull data from multiple platforms is needed. Node.js built on Google’s V8 JavaScript



runtime, an open source JavaScript engine developed for Google Chrome's Web browser is used. V8 compiles to native code prior to execution, optimised and re-optimised dynamically at runtime and is built for browser-based and stand-alone high-performance applications. Node.js supports clustering but a PC with a single core is generally sufficient for supporting large input streams. Node.js is purely a server-side application and therefore, integration algorithms will need to be written for various API and Web Services. A suite of basic libraries and modules accompany Node.js (HTTP, Sockets.io, etc., amongst which are modules that connects to social media). The asynchronous nature of Node.js algorithms makes it capable of pulling large volumes of data in increasingly high velocity on a server side application (Ch'ng, 2014).

*3.1.3 Scalable data storage.* The big truth about Big Data is that it is easier to get the data in than out (Jacobs, 2009, p. 4) from relational database management systems. The storage of Big Data in the Terabytes require a format that can be stored, but accessed quickly and processed on the fly in real-time or near real-time. Raw data from web sites and social media will be inconsistent and unstructured. As such, indexing may be difficult. As noted by Jacobs, "To achieve acceptable performance for highly order-dependent queries on truly large data, one must be willing to consider abandoning the purely relational database model".

A NoSQL (key-value pairs) with non-relational features that scales massively and horizontally with a number of solutions that makes it easy to shard across distributed systems is the key to storing structured social media data. There are benefits in using NoSQL databases of scientific data storage. One of the most important aspects of statistical correlation is the storage of as much data as a project can possible obtain. Storing all data for filtering later is a better approach than not having the full amount of data, for data lost means opportunity lost forever. This is important as web sites conduct routine archive and removal of old data. Social media and social network API limit access to past data (e.g. Twitter restricts access to stored data and only releases 10 per cent of their API search). Streaming real-time data from API however is unrestricted and therefore is mandatory, which means that all incoming data streams have to be stored.

MongoDB, a cross-platform open source freeform document database that stores document in dynamic schemas as JavaScript Object Notation (JSON) (MongoDB terms it BSON). JSON is a lightweight format that can be used for structuring transmitting data between a server and Web application and is convenient as a data format for Node.js' JavaScript. In this project, the replication of MongoDB is used for increasing the distributed storage and throughput of data. All incoming data are parsed, categorised, and stored in MongoDB shards via Node.js.

Redis is another open source, networked, in-memory, advanced key-value store data structure server for fast access of structured data. This is possible as Redis holds the whole data set in memory with configuration to allow semi-persistence of data via snapshotting. Redis can be positioned as an intermediary between MongoDB and other necessary services (e.g. Hadoop) within the Big Data architecture, such as the processing of data sets, data mining, data and social networks, mapping, and visualisation. Both MongoDB and Redis are scalable and supports clustering.

### *3.2 Mapping social networks*

Social media contains hidden networks and can be constructed as SILs. The components of an SIL consist of users, contents, and links that made up the landscape. The subsequent sub-section covers the methods and algorithms for mapping, reconfiguring, and storing SIL.

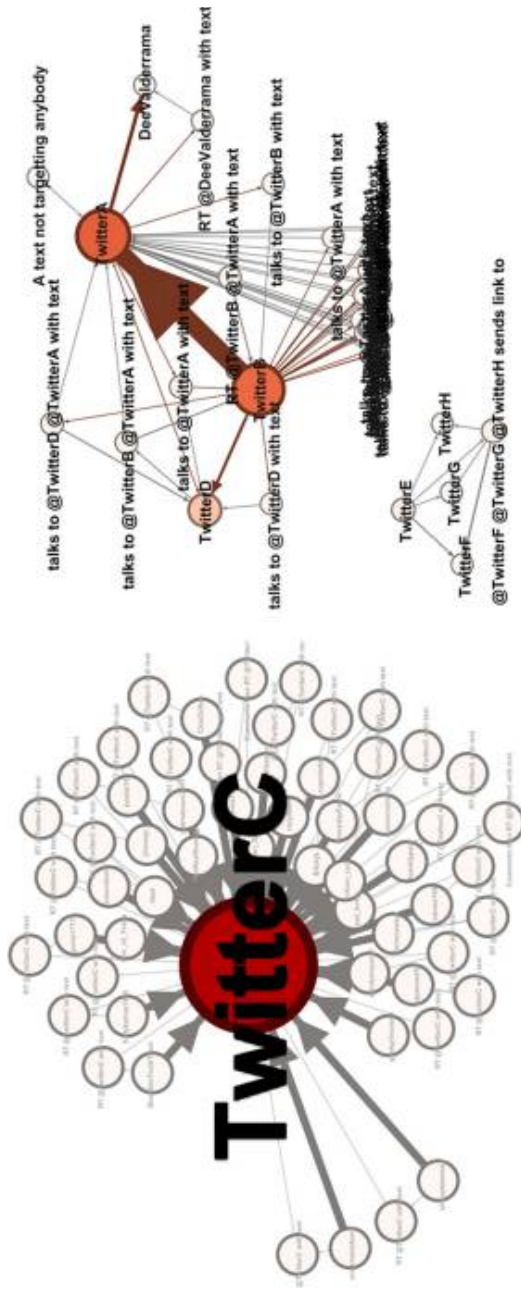
*3.2.1 Mapping concepts.* Figure 2 presents three models of extremely small multimodal networks, which can be used to compare the differences between various node communication scenarios. The Twitter social network is used as an example here for reasons stated earlier, and for a clearer understanding of the data to be presented in Section 4, but any multimodal network can be mapped in a similar fashion. Visualisation uses the Gephi graph visualisation package, but any network visualisation application can be used. Gephi's layout feature (ForceAtlas and ForceAtlas 2) and ranking for setting up the colour scheme. The average degree and average path length is used for calculating the statistics of the network, these are degree and betweenness centralities.

In Figure 2, on the left is a larger node with many retweets in-degrees (directional edges) pointing towards node "TwitterC". On the left of TwitterC are four protruding nodes, which were deliberately "pulled out" of the cluster for visual clarity. They are Twitter users who retweeted TwitterC's tweet. There is a thick arrow coming out of each of the Twitter user node, with another arrow pointing towards the retweet node (the post) from it, which points back towards TwitterC. All retweets have a similar pattern. This records both the Twitter user and the retweets as a multimodal connection.

At the bottom right of the figure is a small cluster consisting of five interactions, with TwitterE tweeting to TwitterF, G, and H. TwitterG tweets a post to TwitterF, TwitterH, and TwitterG.

The cluster of nodes with two larger nodes TwitterA and TwitterB at the top right of the figure is a conversational network, typical of an online activity network. Such a configuration may reveal a community of some sort (Ch'ng, 2015c). At the bottom of the community cluster is a group of nodes with directed edges from TwitterB showing conversations "talks to @TwitterA with text" directed towards TwitterA (directed edges from the nodes), these group of nodes are responsible for the thick arrow pointing towards TwitterA. At the top right of the cluster, TwitterA tweeted a node with a double directional edge with the dummy text "A text not targeting anybody". This is a tweet (or a comment) without a target. The amount of the interaction defines the weight of the directed edge. The edge between TwitterB and TwitterD is not as thick but is thicker than the other edges, showing more interactions between the two nodes. Other edges and nodes show the direction of conversation, which gives it a distinct visual pattern as compared to the retweets. The unique arrangements of the nodes give rise to distinct features within a large network, as we shall see in Section 4. Activity nodes such as the clusters show here may allow us to identify and isolate communities from small and large-scale networks.

*3.2.2 Data structures for dynamic and static graphs.* A standard data format Graph Exchange XML Format (GEXF) can be used for representing social networks. GEXF is a language for describing complex network structure and their associated data and dynamics. Whilst data can be output to any of the network-based formats (GraphML, XGMML, SVG, GDF), GEXF preceded the other formats as it is relatively mature, robust, flexible, and includes time as one of the dynamics feature. The GEXF format provides a way to visualise time-based networks that other formats do not. The algorithms within the mapping system should consist of nodes representing users and posts (tweets, comments, posts, etc.), edges representing links between users, and between users and their posts. In multimodal networks, posts are included as nodes so that the expression of the landscape has distinctive features when visualised using a combination of pre-set algorithms in some visualisation software (e.g. Gephi). All informational entities should be recorded with a timestamp so that the evolution of occurrence and growth of the network can be reconstructed in a dynamical way.



**Notes:** The unique configuration of nodes and edges make it possible to isolate activities in social media. The visualisation uses the Gephi graph visualisation software

The data structure of a typical GEXF file is in the listings below, with a graphical representation in Figure 5(f). In the figure,  $N$  refers to a node (a person) with its corresponding number whilst  $P$  refers to a post (a tweet, a comment, etc.).

Listing 1 shows the head and tail of the data structure of a typical GEXF file, with a graphical representation of the network. Between the `<nodes></nodes>` and `<edges></edges>` tags are all the nodes and edges of the graph, shown in Figures 3 and 4.

Listing 1: the top and tail of the data structure of a typical GEXF file, with a graphical representation of the network. Between the `<nodes></nodes>` and `<edges></edges>` tags are all the nodes and edges of the graph, shown in Figures 3 and 4:

```
<?xml version = '1.0' encoding = 'UTF-8'? >
< gexf  xmlns = 'http://www.gexf.net/1.2draft' xmlns:viz = 'http://www.gexf.net/1.2draft/viz' version = '1.2' >
  < meta lastmodifieddate = '113-6-4' >
    < creator > Eugene Chng </creator >
    < description > Twitter Map </description >
  </meta >
  < graph defaultedgetype = 'directed' idtype = 'string' mode = 'dynamic' timeformat = 'dateTime' start = '2013-05-18T05:00:58.975+0100' end = '2013-05-18T06:58:25.584+0100' >
    < attributes class = 'node' mode = 'dynamic' >
      < attribute id = 'description' title = 'Description' type = 'string'/ >
    </attributes >
    < attributes class = 'edge' mode = 'dynamic' >
      < attribute id = 'weight' title = 'Weight' type = 'float'/ >
    </attributes >
    < nodes count = '12' >
      ... // all nodes here
    </nodes >
    < edges count = '16' >
    </edges >
      ... // all edges here
    </graph >
  </gexf >
```

Particular attention should be paid to a few lines of code and the attributes of the tags in Listing 1 and Figures 3 and 4. Specifically, due to the nature of the dynamic graph, the element of time (start and end attributes) is important. In the `<graph>` tag (bold), the start and end attributes shows the entire timespan of the network. The other tags are the `<nodes>` and `<edges>` tags, which shows the size of the nodes (Figure 3) and edges (Figure 4) in the network.

```
<node id='N0' label='N0' start='2013-05-18T05:01:10.901+0100' end='2013-05-18T05:07:10.901+0100'>
  <attvalues>
    <attvalue for='description' value='Sat May 18 05:01:15 BST 2013' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='N1' label='N1'>
  <attvalues>
    <attvalue for='description' value='N1' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='P0' label='P0' start='2013-05-18T05:01:10.901+0100' end='2013-05-18T05:07:10.901+0100'>
  <attvalues>
    <attvalue for='description' value='null' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='N2' label='N2' start='2013-05-18T05:01:25.684+0100' end='2013-05-18T05:07:25.684+0100'>
  <attvalues>
    <attvalue for='description' value='Sat May 18 05:01:30 BST 2013' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='P1' label='P1' start='2013-05-18T06:01:25.684+0100' end='2013-05-18T06:07:25.684+0100'>
  <attvalues>
    <attvalue for='description' value='null' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='N3' label='N3' start='2013-05-18T06:41:25.684+0100' end='2013-05-18T06:47:25.684+0100'>
  <attvalues>
    <attvalue for='description' value='N3' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='P2' label='P2' start='2013-05-18T06:51:25.684+0100' end='2013-05-18T06:57:25.684+0100'>
  <attvalues>
    <attvalue for='description' value='null' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='N4' label='N4' start='2013-05-18T06:51:25.684+0100' end='2013-05-18T06:57:25.684+0100'>
  <attvalues>
    <attvalue for='description' value='N3' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='P4' label='P4' start='2013-05-18T06:51:25.684+0100' end='2013-05-18T06:57:25.684+0100'>
  <attvalues>
    <attvalue for='description' value='null' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='N5' label='N5' start='2013-05-18T06:52:25.584+0100' end='2013-05-18T06:58:25.584+0100'>
  <attvalues>
    <attvalue for='description' value='NN' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='P3' label='P3' start='2013-05-18T06:52:25.584+0100' end='2013-05-18T06:58:25.584+0100'>
  <attvalues>
    <attvalue for='description' value='null' />
  </attvalues>
  <viz:size value='1.0' />
</node>
<node id='P5' label='P5' start='2013-05-18T06:58:15.584+0100' end='2013-05-18T07:15:15.584+0100'>
  <attvalues>
    <attvalue for='description' value='null' />
  </attvalues>
  <viz:size value='1.0' />
</node>
```

**Figure 3.**  
Nodes of the  
graph in Figure 5,  
inserted within  
the <nodes>  
</nodes> tags  
in Listing 1

```
<edge id='N0-N1' source='N0' target='N1' type='directed' label='2'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
<edge id='N0-P0' source='N0' target='P0' type='directed' label='0'>
  <attvalues>
    <attvalue for='weight' value='1' start='2013-05-18T05:01:10.901+0100' />
  </attvalues>
</edge>
<edge id='P0-N1' source='P0' target='N1' type='directed' label='1'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
<edge id='N4-N1' source='N4' target='N1' type='directed' label='2'>
  <attvalues>
    <attvalue for='weight' value='2' />
    <attvalue for='weight' value='2' />
  </attvalues>
</edge>
<edge id='N4-P4' source='N4' target='P4' type='directed' label='0'>
  <attvalues>
    <attvalue for='weight' value='1' start='2013-05-18T05:01:10.901+0100' />
  </attvalues>
</edge>
<edge id='P4-N1' source='P4' target='N1' type='directed' label='1'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
<edge id='N2-N3' source='N2' target='N3' type='directed' label='2'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
<edge id='N2-P1' source='N2' target='P1' type='directed' label='0'>
  <attvalues>
    <attvalue for='weight' value='1' start='2013-05-18T05:01:25.684+0100' />
  </attvalues>
</edge>
<edge id='P1-N3' source='P1' target='N3' type='directed' label='1'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
<edge id='N2-N4' source='N2' target='N4' type='directed' label='2'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
<edge id='N2-P2' source='N2' target='P2' type='directed' label='0'>
  <attvalues>
    <attvalue for='weight' value='1' start='2013-05-18T06:51:25.684+0100' />
  </attvalues>
</edge>
<edge id='P2-N4' source='P1' target='N4' type='directed' label='1'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
<edge id='N5-P3' source='N5' target='P3' type='directed' label='0'>
  <attvalues>
    <attvalue for='weight' value='1' start='2013-05-18T06:52:25.584+0100' />
  </attvalues>
</edge>
<edge id='N4-N1' source='N4' target='N1' type='directed' label='2'>
  <attvalues>
    <attvalue for='weight' value='2' />
    <attvalue for='weight' value='2' />
  </attvalues>
</edge>
<edge id='N4-P5' source='N4' target='P5' type='directed' label='0'>
  <attvalues>
    <attvalue for='weight' value='1' start='2013-05-18T06:58:15.584+0100' />
  </attvalues>
</edge>
<edge id='P5-N1' source='P5' target='N1' type='directed' label='1'>
  <attvalues>
    <attvalue for='weight' value='1' />
  </attvalues>
</edge>
</edges>
```

**Figure 4.**  
Edge nodes of the  
graph in Figure 5,  
inserted within  
the <edges>  
</edges> tags  
in Listing 1



Figure 3 shows a listing of all the nodes of the network. The most important attributes in the `<node>` tag are the id and the label. The id sets the uniqueness of the id for graph processing software and the label is displayed visually as a label on the nodes (see Figure 5).

Figure 4 shows a listing of all the edges of the network. Three most important attributes are the id, the source and the target, which describes the connection between the source and the target node in Figure 3. The id “N4-N1” of an edge has two `<attvalue>` tags showing a value of 2:

`<attvalue for = “weight” value = “2” />`

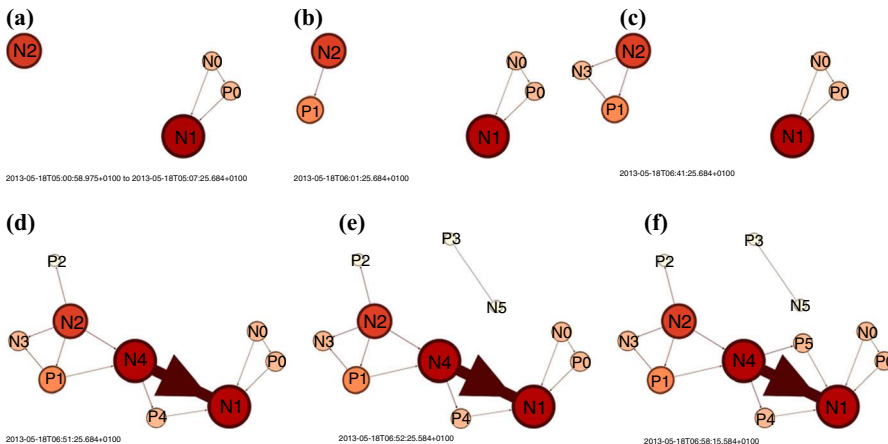
This is read by graph visualisation software (i.e. Gephi) as the weight of the connection. N4’s connection to N1 in Figure 5 has a thicker line because there were two interactions:

N4→P4→N1  
N4→P5→N1

The only disadvantage of the GEXF file format is perhaps the size attributed to the XML tags used for enclosing and describing the contents. In extremely large data set, an alternative can be used, by removing the time element if they are not important to the research inquiry. For static maps which are much smaller in file size, a GDF file can be generated (see Section 3.2.6 for algorithm) as an output. Listing 2 is a GDF version of the GEXF graph file. The file is self-explanatory. The first section beginning with `node def >` followed by a list of attributes describing each line of node properties are the nodes. The second section beginning with `edge def >` followed by the connection and the weight are edges.

Listing 2: the listing shows the static graph version of the GEXF file format. The file format is GDF:

```
node def > name VARCHAR,timeStamp VARCHAR,label VARCHAR,
N0,2013-05-18T05:01:10.901+0100,N0,
```



**Notes:** The graph is described in the GEXF data structure in Figure 3. The numbering shows the sequence. The dates are in the order of when new nodes and edges appear and connect to the network

**Figure 5.** A sequence diagram showing time-based dynamic graph, with nodes and edges appearing according to the timestamps

IMDS  
115,9**1738**

N1,2013-05-18T05:01:10.901+0100,N1,  
P0,2013-05-18T05:01:10.901+0100,P0,  
N2,2013-05-18T05:01:25.684+0100,N2,  
P1,2013-05-18T06:01:25.684+0100,P1,  
N3,2013-05-18T06:41:25.684+0100,N3,  
P2,2013-05-18T06:51:25.684+0100,P2,  
N4,2013-05-18T06:51:25.684+0100,N4,  
P4,2013-05-18T06:51:25.684+0100,P4,  
N5,2013-05-18T06:52:25.584+0100,N5,  
P3,2013-05-18T06:52:25.584+0100,P3,  
P5,2013-05-18T06:58:15.584+0100,P5,  
edgedef > node1 VARCHAR,node2 VARCHAR,weight INT,  
N0,N1,2,  
N0,P0,0,  
P0,N1,1,  
N4,N1,2,  
N4,P4,0,  
P4,N1,1,  
N2,N3,2,  
N2,P1,0,  
P1,N3,1,  
N2,N4,2,  
N2,P2,0,  
P1,N4,1,  
N5,P3,0,  
N4,N1,2,  
N4,P5,0,  
P5,N1,1,

*3.2.3 Mapping algorithm.* The mapping algorithm maps social media activities into an SIL. There are three components in such a network – the user, the post, and a collection of users mentioned in the post. A user  $v$  participates in social media on a topic (collected with a hash tag, or keywords, within a product page or a media channel). The total collection of users within the topic is  $V$ . Each user may comment, tweet, or post a statement  $p$ , which may contain mentions of a user  $m$ . Within a post  $p$ , there may be a collection of mentions of other users  $M$ .  $e$  is an edge between two nodes and  $E$  is the collection of all the edges. The summary of the sets is:

$v$ : a user object, also containing all the attributes associated with the user (timestamp, post, etc.)

$V$ : a collection of all the users and their attributes

$n$ : a node representing a user, containing all the information in  $v$



N: all nodes in the network  
 p: a post by a user  
 m: a user mentioned within the post  
 M: all users mentioned within the post  
 e: a single edge between two nodes  
 E: a collection of edges.

At the start of the algorithm, each user with the associated information is queried from the database and stored  $V \leftarrow \{\{v\} : v \in V\}$ . For each user, first assign the post  $p \leftarrow$  the post of  $v$ . Get all the mentions from the post  $M \leftarrow \{\{m\} : m \in M\}$ . Check if the user exists in the node collection  $N$ , if not, add the user to the node collection  $\text{Add}(v)$  to  $N$ . If no users  $m$  are mentioned in the post  $p$ , add the post as a node  $\text{AddNode}(p)$ . If users are mentioned in the post, and if any user is not in the collection  $N$ , add the user to the collect  $\text{AddNode}(m)$ . The procedure  $\text{CheckMentions}(M)$  is as follows: For all mentions, add an edge between the user and the post  $\text{AddEdge}(n,p)$ , add an edge from the post to the user mentioned  $\text{AddEdge}(p,m)$ , finally, increase the weight of the edge by 1 between  $n$  and  $m$   $\text{Increment}(w, 1)$  of  $n$  and  $m$ . If the user exists in the collection of nodes call the procedure  $\text{CheckMentions}(M)$ . Finally, the data is output as a GEXF file. The mapping algorithm:

```

V ← { {v} : v ∈ V }
FOR all v ∈ V DO
  p ← the post of v
  M ← { {m} : m ∈ M } in mentions m of the post p of v
  IF v ∉ N THEN Add(v) to N, DO
AddNode(p)
  IF count(M) < 1 DO
    AddEdge(n,p)
  ELSE
    proc CheckMentions(n);
  FOR all m ∈ M
    IF m ∉ N THEN AddNode(m) to N
      AddEdge(n,p), AddEdge(p,m), AddEdge(n,m)
    Increment(w, 1) of n and m
  END-IF
END-FOR
end-proc
END-IF
ELSE
  CheckMentions(n)
END-IF
END-FOR
Output(GEXF)

```

3.2.4 *Reconfiguration algorithm.* The SIL mapped with the algorithm in Section 3.2.3 has this configuration:

$$n \rightarrow p \rightarrow m, n \rightarrow m$$

The network can be reconfigured to  $n \rightarrow p \rightarrow m$  with the algorithm below, this simply “unwraps” the complete connection by severing  $n$  and  $m$ ’s edge. The reconfiguration algorithm: “unwrapping” the network:

```
N ← { {n} : n ∈ N } from GEXF file
```

```
FOR all n ∈ N DO
```

```
    RemoveEdge(n,m)
```

```
END-FOR
```

```
    Output(GEXF)
```

The network can also be reconfigured to remove all post, containing users only  $n \rightarrow m$ , with the algorithm below. The reconfiguration algorithm: retaining users and removing posts:

```
N ← { {n} : n ∈ N } from GEXF file
```

```
FOR all n ∈ N DO
```

```
    RemoveEdge(n, p), RemoveEdge(p, m), RemoveNodes(p)
```

```
END-FOR
```

```
    Output(GEXF)
```

3.2.5 *Dynamic graph output.* The GEXF dynamic graph output (Section 3.2.2) has a simple procedure Output (GEXF). The header is written first by streaming a text output and inserting information where needed into the header, such as the start and end timestamps of the first appearance of, and the last appearance of the final node, and subsequently the opening and closing tags with the node and edge counts that encloses the individual node and edge tags. For each node, write the node tag WriteNode( $n$ ) with associated information such as the start and end timestamps. For each edge, write the edge tag WriteEdge( $e$ ) with associated source and target attributes, and add the weight attributes for visualisation purposes. GEXF file output:

```
proc Output(GEXF);
```

```
    WriteHeader(N)
```

```
    WriteNodesTag(Count(N))
```

```
FOR all n ∈ N DO
```

```
    WriteNode(n)
```

```
END-FOR
```

```
CloseNodesTag()
```

```
WriteEdgesTag(Count(E))
```

```
    FOR all e ∈ E in N, DO
```

```
        WriteEdge(e)
```

```
    END-FOR
```

```
CloseNodesTag()
```

```
WriteFooter()
```

```
end-proc
```

---

*3.2.6 Static graph output.* The GDF format (Section 3.2.2, Listing 2) is a static format that has the same expression as the dynamic version, apart from the removal of the time element. The output and algorithm is straightforward. For the nodes and edges, respectively, the labels are being written first before the corresponding node and edge lists are appended. GDF file output algorithm:

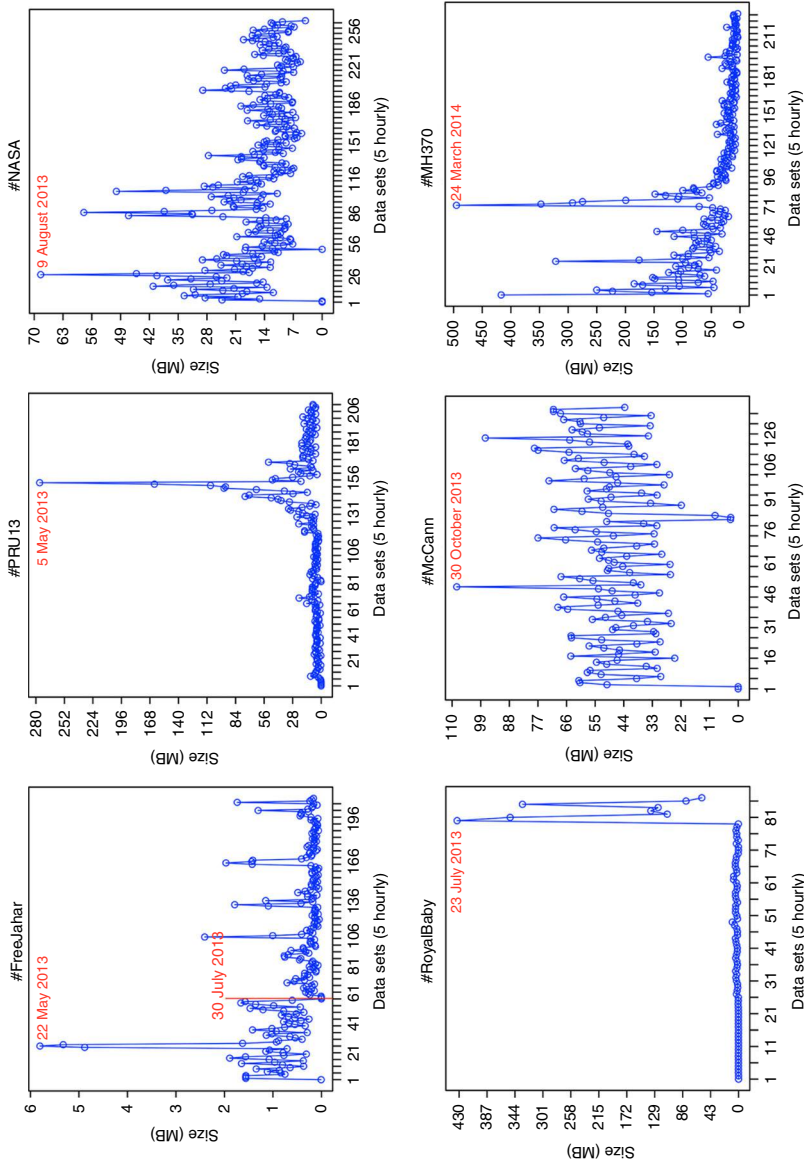
```
proc Output(GDF);
  writeNodeLabels(N);
  FOR all n ∈ N DO
    WriteNode(n)
  END-FOR
  writeEdgeLabels(N);
  FOR all e ∈ E in N, DO
    WriteEdge(e)
  END-FOR
end-proc
```

#### 4. Visualizing mapped SIL

This section demonstrates the feasibility of the Big Data architecture and the algorithms used for mapping SIL. The first subsection looks at the capability of the methods in managing the volume and velocity of topical data sets, followed by samples expressions of SIL retrieved from the topical data sets. The final subsections inspect the details of some of the graphs.

##### 4.1 Volume and velocity

Figure 6 are graphs of viral outbursts of trending data captured from Twitter with the Big Data architecture. Each point within the graphs is a five hourly recording of Twitter activities as a SIL with the associated hash tags (caption at the top), mapped as an SIL using the GEXF file format. The #MH370 data set has over 21 data points, for example, whilst the #RoyalBaby data set has 91 data points. The volume and velocity of each trending topic can be observed, particularly at the spikes of the data point at particular events released by news channels. #MH370 has one of the largest data point at around 500 Megabytes on the 24 March 2014, all the data points in total amounted to over 14 Gigabytes. The graphs show different signatures. These signatures are by-products of the mapping algorithms and are potentially a useful source of information. For example, the spikes in the #FreeJahar data set are news releases in associated Twitter accounts of news channels, out of which the majority of the data are retweets. However, in the case of the #FreeJahar, communities were identified within those data sets (see Ch'ng, 2015c), and community activities were heightened in the spiked data points. The #PRU13 is a Malaysian political election data set, it is easy to see that 5th May of 2013 has the highest activity as it was the day of the election. There is another spiked data point a little later, which is related to news of swapped ballot boxes during a supposedly deliberate electrical black out. The hills and valleys of data points for other data sets may just be the heightened activities during the day and the easing off of tweets during the night (e.g. the #McCann data set).



**Figure 6.** Graphs of viral outbursts of trending data captured from Twitter with the Big Data architecture

**Notes:** Each point in the graphs is a five hourly recording of a Twitter SIL with the associated hash tags (caption at the top). The #MH370 datasets have over 21 points for example, whilst the #RoyalBaby dataset has 91 data points

#### 4.2 Expressions of SIL

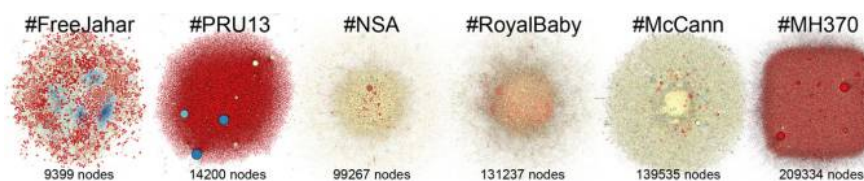
In Figure 7 are six SILs sampled from the topical data sets in Figure 6. Here, the value of the methods for mapping SIL is demonstrated. Each SIL (a data point) is a recording of five hourly activities in Twitter. Each SIL shows a different landscape expression. The degree centrality (number of in-degrees and out-degrees) is used for measuring each node, resulting in the more important nodes (larger in size) with a higher degree centrality within the multimodal activity networks. The number of nodes in each SIL is shown at the bottom of each graph. The graphs described with the GEXF format is visualised with the Gephi graph visualisation package.

The largest graph in Figure 7 has over 200,000 nodes (#MH370), the number of nodes and edges in the graph is a challenging problem as it may be difficult for certain graph visualisation packages to process. However, many of the nodes may not be important to social network analysis. The unimportant nodes are usually unconnected, or the nodes may be individuals who posted only a single tweet. Unless semantic or sentiment analysis is conducted on the full set, these nodes can be safely removed from the graph. Figure 8 is a set of graphs with nodes which have a limited number of edges removed from the graphs in Figure 7, the caption on the top right corner of each inset shows the reduced node and edge size, the parenthesis (e.g. < 3) is the filter used for reducing nodes with edges less than a number, in the case of #PRU is nodes with edges < 3. The smaller graphs are useful for efficient graph processing using graph visualisation software whilst maintaining the quality and meaning of the SIL.

The different signatures in each of the SIL are obvious here. The signature of each SIL is due to the nature of the activities, the background of the topic, and the intentions of the social media users. It can be observed that users who interacted more are closer together in that they form natural clusters. The #FreeJahar data set which has very interesting phenomenon has been analysed in detail (Ch'ng, 2015c).

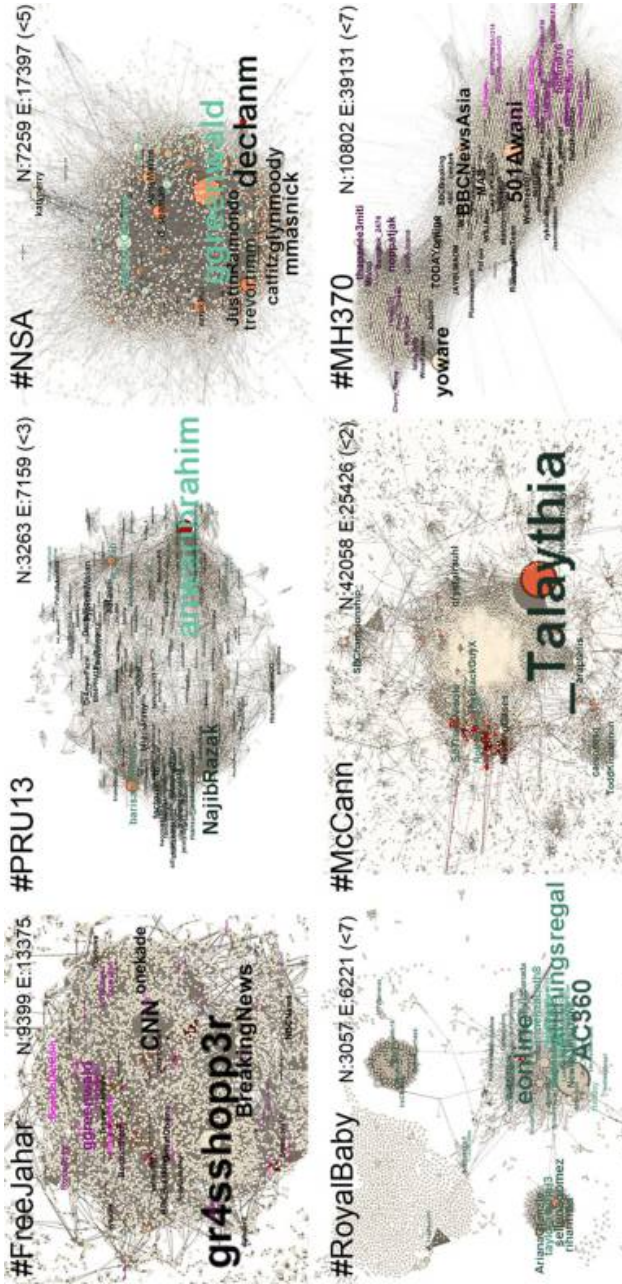
The other data sets are also of considerable interests to academics in the social sciences. Here, a brief description of the social network expression is sufficient until further results from the analysis are published.

In the #PRU graph, two large political parties contended during the Malaysian general election bridged by activists, the opposition leader @AnwarIbrahim is flanked by media channels and have large degree and betweenness centralities, an indication of heightened activities which greatly strengthened the opposition party's stand. In the controversial #NSA news, within the graph, a Guardian correspondent is highly active, with clusters of activities directly below in the SIL. Note that #KatyPerry is above in a separate but linked community. The glamorous birth of the #RoyalBaby were mentioned by celebrities (e.g. @selenagomez, etc.) with pockets of discussions. The reawakening of the Madeleine #McCann kidnapping case due to new evidence led to



**Notes:** Each SIL (a data point) is a recording of five hours of activities in Twitter. They show different expressions. The important nodes (larger in size) are nodes with a higher degree centrality within the multimodal activity networks

**Figure 7.** Each SIL here is a data point sampled from the topical graphs in Figure 6



**Figure 8.**  
The SIL datasets here are decimated graphs from Figure 7

**Notes:** Most singular nodes have been removed for more efficient graph processing using graph visualisation software whilst maintaining the quality and meaning of the SIL



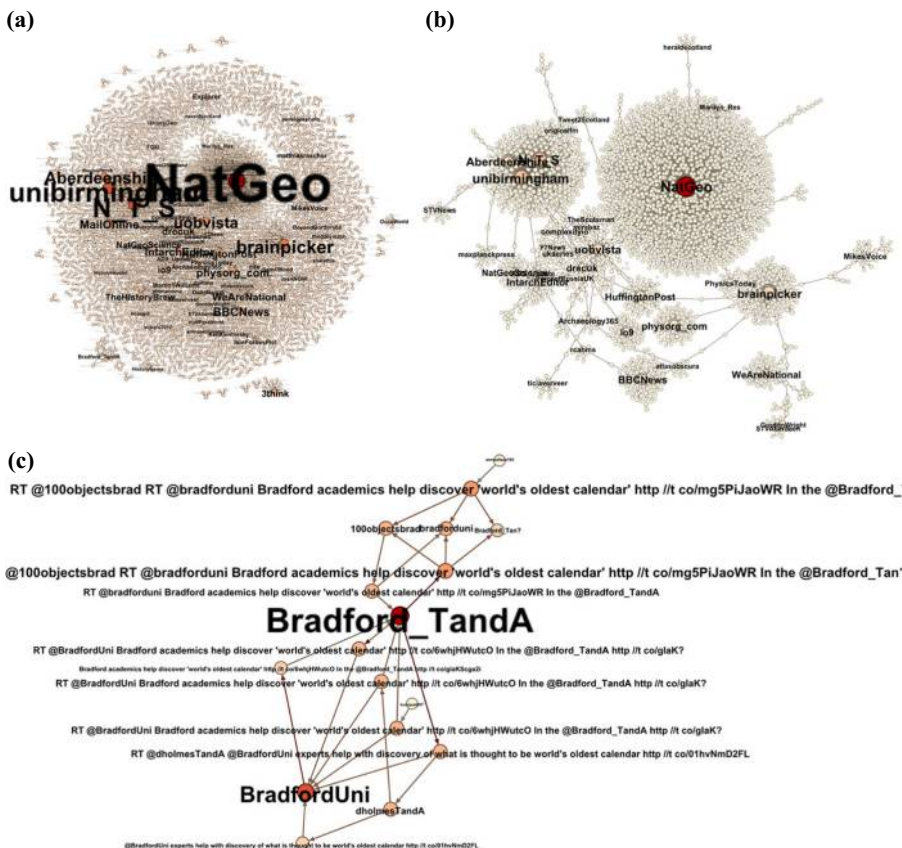
separate clusters of varying degree of intense conversations. A single MH370 SIL datapoint contains two separate clusters of important posts by activists.

The brief analysis and visualisation of the above data sets raise some interesting questions. For example, why are there many more distinct clusters of activities within the #McCann SIL as compared to the other data sets? Is this an indication that there were many kidnapping-related experiences in the separate clusters of social media users? Whilst the answers to these questions are interesting, they are beyond the scope of this paper.

#### 4.3 Network structures and centralities

In this section, three data sets captured from social media are shown. The data sets demonstrate the details of the algorithms in mapping SIL.

Figure 9 is a capture and mapping of a SIL 15 July 2013 for 9 days on the press release related to the publication of the discovery of “World’s Oldest Calendar”

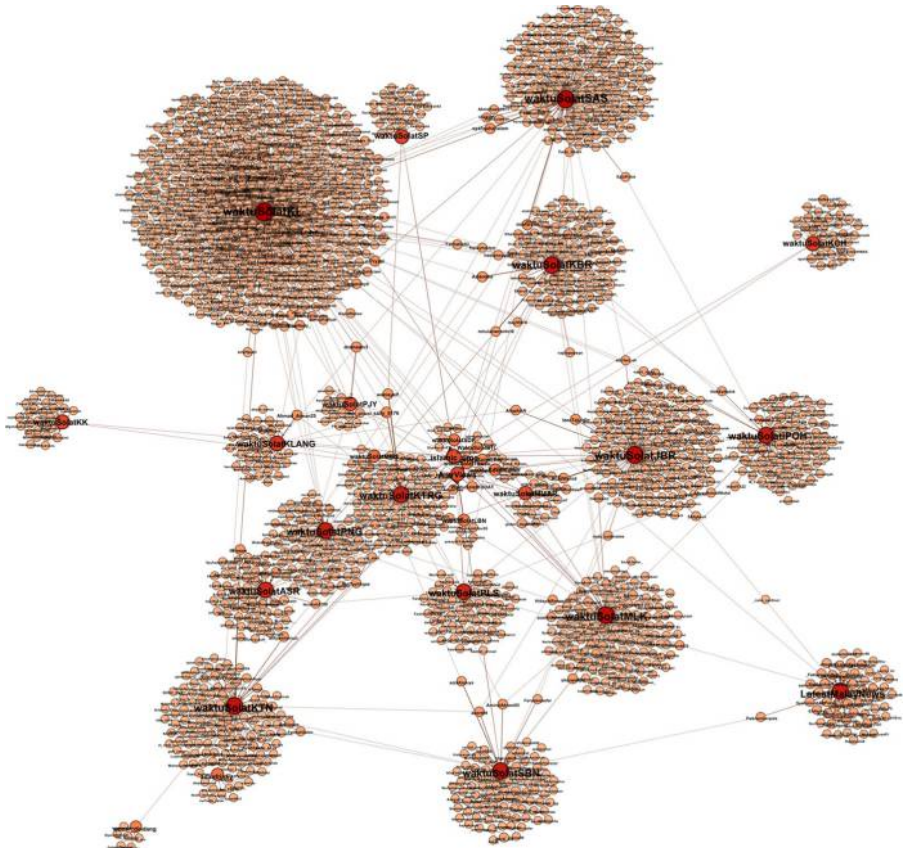


**Notes:** (a) The complete dataset composed of 9 days of activities since 15 July 2013 using the mapping algorithm (7187 nodes, 10322 edges); (b) an unwrapped version of the SIL using the reconfiguration algorithm in 3.2.4. Isolated nodes have been removed, leaving only connected nodes in the network (2060 nodes, 2203 edges); (c) a closer view of an isolated cluster within the SIL, showing the detailed connections between different types of nodes (users and posts)

**Figure 9.**  
Trending tweets of  
press releases related  
to the publication of  
discovery of the  
“World’s Oldest  
Calendar”

(Gaffney *et al.*, 2013). There were 44 data sets (5 hourly each data point), which have been combined into a single map in Figure 9. There were 5,973 nodes (users and their tweets) and 6,100 edges (connecting users and tweets). The relative importance of the nodes within the SIL in diffusing information uses the in- and out-degree centrality (shown as larger, redder nodes).

Figure 9(a) shows the complete data set composed of nine days of activities since 15 July 2013 using the mapping algorithm (7,187 nodes, 10,322 edges). Figure 9(b) shows an unwrapped version of the SIL using the reconfiguration algorithm in 3.2.4. The isolated nodes have been removed, leaving only connected nodes in the network (2,060 nodes, 2,203 edges). Figure 9(c) shows a closer view of an isolated cluster within the SIL, illustrating the detailed connections between different types of nodes (users and posts). Figure 10 is a visualisation of an SIL sampled from one of the MH370 data points. Posts have been removed, leaving only user nodes within the network of activities, which forms a small world network (Ch'ng, 2015b). The important nodes have higher degree centralities (larger and redder nodes). The nodes are news correspondence of major news channels. The main activities here are retweets from



**Figure 10.**  
An SIL sampled  
from one of the  
MH370 data points

**Notes:** Posts have been removed, leaving only user nodes within the network. The main activities here are retweets of news events

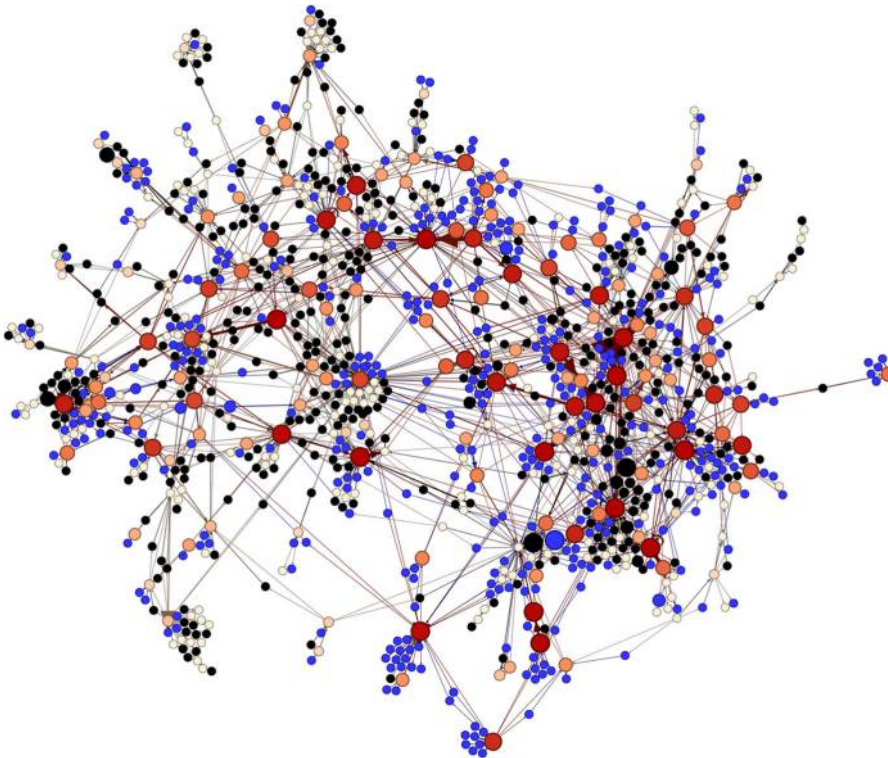


users of news events. Note that unlike the #FreeJahar SIL in Figure 11, there are no conversational clusters.

Figure 11 is a visualisation of the largest five-hourly #FreeJahar SIL. The expression of the SIL is typical of a network with heightened conversational activities, the persistence of the cluster indicates that a community has formed. The SIL is mapped with the algorithm in 3.2.3 and has not been reconfigured. The nodes are colour coded – black nodes are retweets, blue nodes are conversational posts, the shades of red are Twitter users. The reddest nodes are ones with higher measure of betweenness centrality, a sign that they are the “gate keepers” of information. An analysis of the #FreeJahar data set collected and mapped with the methods described here is available (Ch’ng, 2015c). The paper investigates the formation, development, and decline of the #FreeJahar community and elaborates the phenomenon discovered here.

## 5. Discussion

Social media certainly increases our range of communicable activities. Social media also, collectively in the global sense, extends the human connectedness beyond the boundary of a participant’s geographical location. The time of interaction between participants and



**Notes:** The expression of the SIL is typical of a network with heightened conversational activities. Black nodes are retweets, blue nodes are conversational posts, the shades of red are Twitter users. The reddest nodes are ones with higher measure of betweenness centrality

**Figure 11.**  
Visualisation of the  
#FreeJahar SIL

communications sent now may be retrieved and responded to, much later by other participants due to Information Technology infrastructures. Thus, the distance between time and space has been greatly shortened in the age of social media. This invariably opens up a broad range of opportunities as space and time, in the eye of a user are “compressed” to within a digital display, allowing diverse communications from extremely large demographics groups. Participants interact without thinking much about the information networks that they are helping to build, and users reading social media are oblivious to the fact that hidden networks exist. This is due to the fact that the limited presentation of the sequential nature of content listings within social media, as viewed on web browsers or on mobile devices do not necessarily reveal nor make obvious an unknown nature of the medium; that every participant, from content producers, to consumers, to followers and subscribers, including the contents they produce or subscribed to, are intrinsically connected in a hidden but potentially massive contextual network. These hidden networks are a valuable source of information for scientific inquiry. Such networks when mapped and analysed using appropriate methods could reveal valuable information. These multimodal networks, together with information about the nodes, direction of edges, the polarity of interactions, the structure of the graph and secondary data captured longitudinally form a SIL. SIL have become an important aspect of contemporary scientific inquiry, and the importance will only increase in the future when all things both men and machines are digitally connected.

In this paper, the definition of SIL is set, together with a detail of the methods for mapping such landscapes. The mapping of SILs undoubtedly requires a Big Data approach, for the temporal range may be great if meaningful longitudinal data sets are to be acquired. Volume and velocity have been a major challenge in acquiring social media data sets, this may mean that manual approaches in collecting and mapping networks will no longer be feasible. To address the issue of volume and velocity, a Big Data architecture integrated with scalable open source libraries and API is provided as a solution. The solution uses server-side asynchronous I/O software developed for connecting to and streaming data from social media and the web. The system is designed so that incoming data is immediately stored in scalable and distributed NoSQL databases, and when the data collection is complete, separate algorithms are used for streaming structured data in stable, evenly spaced chunks for mapping SIL. A set of algorithms were presented in Section 3 for mapping, reconfiguring and storing SIL, together with patterns for dynamic and static data formats that could be used with graph visualisation packages. Section 4 demonstrates the feasibility of the Big Data architecture and the usefulness of mapping large multimodal and longitudinal data by presenting real-world social media data sets as examples.

The methodology presented here has laid a foundation for mapping multimodal and longitudinal networks. However, much work is required for the future, especially with regards to issues of Big Data. One of such issues is the (1) processing of extremely large graphs with hundreds of thousands of nodes and edges. In the data sets given within this paper, the largest graph reaches 200k nodes in a single data point, in actual fact however, a SIL is a combination of all the data points as the activities are continuous through time. This will bring the size of nodes to over one million with edges doubling the size of the nodes. On a Dell PowerEdge C6,100 with 2x Intel Xeon X5,660 Processor (2.80 GHz, 12 M Cache, 6.40 GT/s QPI, Turbo, HT), 1,333 MHz Max Memory, 48 GB Memory (6 × 8 GB Dual Rank RDIMMs), 200k nodes is manageable. In larger social media involving months to years of data collection, the size could easily reach tens of millions. Extremely large graph processing is therefore an issue and needed parallel

and distributed processing. The second issue is (2) real-time mapping of networks. Unlike the data management covered in this paper, which stores all data sets prior to processing them in a stable stream, real-time processing may need to discard unwanted or redundant raw data as they may be too big for storage, keeping only data which has been structured. This invariably requires real-time mapping and configuration of graphs, implying that parallel and distributed processing is required. The final issue is (3) large SIL visualisation, which involves (1) and (2) will require distributed and parallel processing, particularly with GPGPUs. A single GPGPU with thousands of processing cores within a single computer is a suitable platform for developing algorithms that manage all three issues discussed above. The Big Data and Visual Analytics Lab (Ch'ng, 2015a) is at present working with Nvidia in various research projects in developing algorithms in all three issues as part of our future work. One of the propositions for visualisation is described in Sun *et al.*'s (2015) article.

The methodology presented in this paper has practical implications not only in the academia, but also in marketing, business applications, and security. The mapped SILs can be used for tracking the outreach and impact of topical marketing activities on social networks, via keywords and hash tags for example. Influencing factors such as profile attributes of social media users in marketing strategies can be analysed together with marketing impacts in order to see if the factors correlate with profits in the longer term, if company contents are diffused via social networks using different media types (text, audio, video, and image contents). SILs do have important implications on mapping contextual networks for monitoring terrorist activities, outreach, connections, and context for security purposes (e.g. ISIS has online presence), including terrorist attacks and their psychological effects on civilians and the wider society. Coupled with machine learning techniques, SILs could be a powerful tool that can be used for classification and prediction purposes.

Data relationality is an important issue, and this paper has tackled a very important aspect of the broader issues in Big Data research.

## References

- Anderson, C. (2008), "The end of theory: the data deluge makes the scientific method obsolete", *Wired Magazine*, 16 July, p. 16.
- Bontcheva, K. (2014a), "EU project to build lie detector for social media", available at: <https://www.sheffield.ac.uk/news/nr/lie-detector-social-media-sheffield-twitter-facebook-1.354715>
- Bontcheva, K. (2014b), "Lie detector on the way to test social media rumours", *BBC Technology*, available at: [www.bbc.co.uk/news/technology-26263510](http://www.bbc.co.uk/news/technology-26263510) (accessed 21 April 2015).
- Bourdieu, P. (1985), "The social space and the genesis of groups", *Theory and Society*, Vol. 14 No. 6, pp. 723-744.
- Boyd, D. and Crawford, K. (2011), "Six provocations for big data: a decade in internet time", *Symposium on the Dynamics of the Internet and Society, Social Science Research Network, New York*.
- Ch'ng, E. (2014), "The value of using big data technology in computational social science", in Alvin, C., Zhan, J., Ding, W. and Wu, J. (Eds), *The 3rd ASE Big Data Science 2014, Tsinghua University 4-7 August*, ACM, Beijing and New York, NY, pp. 1-4, available at: <http://dx.doi.org/10.1145/2640087.2644162>
- Ch'ng, E. (2015a), "Big Data and visual analytics lab", available at: [www.nottingham.edu.cn/en/science-engineering/departments/computer-science/research/groups/bdva/big-data-and-visual-analytics.aspx](http://www.nottingham.edu.cn/en/science-engineering/departments/computer-science/research/groups/bdva/big-data-and-visual-analytics.aspx) (accessed 19 February 2015).

- Ch'ng, E. (2015b), "Local interactions and the emergence and maintenance of a Twitter small-world network", *Social Networking*, Vol. 4 No. 2, pp. 33-40.
- Ch'ng, E. (2015c), "The bottom-up formation and maintenance of a Twitter community: analysis of the #FreeJahar Twitter community", *Industrial Management & Data Systems*, Vol. 115 No. 4, pp. 612-624.
- Gaffney, V.L., Fitch, S., Ramsey, E., Yorston, R., Ch'ng, E., Baldwin, E., Bates, R., Gaffney, C., Ruggles, C., Sparrow, T., McMillan, A., Cowley, D., Fraser, S., Murray, C., Murray, H., Hopla, E. and Howard, A. (2013), "Time and a place: a luni-solar "time-reckoner" from 8th millennium BC Scotland", *Internet Archaeology*, p. 34.
- Huberman, B., Romero, D. and Wu, F. (2008), "Social networks that matter: Twitter under the microscope", available at: SSRN 1313405.
- Jacobs, A. (2009), "The pathologies of big data", *Communications of the ACM*, Vol. 52 No. 8, pp. 36-44.
- Jansen, B.J., Zhang, M., Sobel, K. and Chowdury, A. (2009), "Twitter power: tweets as electronic word of mouth", *Journal of the American Society for Information Science and Technology*, Vol. 60 No. 11, pp. 2169-2188.
- Java, A., Song, X., Finin, T. and Tseng, B. (2007), "Why we Twitter: understanding microblogging usage and communities", *Proceedings of the 9th WebKDD and 1st SNA-KDD 2007 Workshop on Web Mining and Social Network Analysis. ACM 12 August*, pp. 56-65.
- Kwak, H., Lee, C., Park, H. and Moon, S. (2010), "What is Twitter, a social network or a news media?", *Proceedings of the 19th International Conference on World Wide Web. ACM 26-30 April*, pp. 591-600.
- Latour, B. (2010), "10 Tarde's idea of quantification", *The Social after Gabriel Tarde: Debates and Assessments*, pp. 145-162.
- Lazer, D., Pentland, A.S., Adamic, L., Aral, S., Barabasi, A.L., Brewer, D., Christakis, N., Contractor, N., Fowler, J. and Gutmann, M. (2009), "Life in the network: the coming age of computational social science", *Science*, Vol. 323 No. 5915, pp. 721-723.
- Manovich, L. (2011), "Trending: the promises and the challenges of big social data", *Debates in the Digital Humanities*, pp. 460-475.
- Ngai, E.W.T., Moon, K.K., S.S., L., Chin, E.S.K. and Tao, S.S.C. (2015), "Social media models, technologies, and applications: an academic review and case study", *Industrial Management & Data Systems*, Vol. 115 No. 5, pp. 769-802.
- Sakaki, T., Okazaki, M. and Matsuo, Y. (2010), "Earthquake shakes Twitter users: real-time event detection by social sensors", *Proceedings of the 19th International Conference on World Wide Web, ACM*, pp. 851-860.
- Sun, H., Ch'ng, E. and See, S. (2015), "A provision for visualisation of evolving virtual communities", *The Ninth International Conference on Operations and Supply Chain Management, University of Nottingham Ningbo China, Ningbo, 13-15 July*.

#### About the author

Dr Eugene Ch'ng is an Associate Professor in Computer Science and Deputy Director for the International Doctoral Innovation Centre (IDIC) Digital Economy Strand at the University of Nottingham Ningbo China campus. He leads the Big Data and Visual Analytics Research Lab at Nottingham's China campus. Eugene holds a Visiting Professorship position at the Centre for Creative Content and Digital Innovation, University of Malaya. Dr Ch'ng has previously served as Innovations Director at the IBM Visual and Spatial Technology Centre and the Digital Humanities Hub a £3.5 m strategic investment bid at the University of Birmingham where he led research in the development and application of cutting-edge technology in digital heritage and culture. Eugene's research has an overarching theme in Complex Systems Science related to the

---

reconstruction and modelling of terrestrial, social, political, and virtual landscapes. These topics naturally involve the collection and generation of massive multimodal and longitudinal data sets and therefore the need for Big Data research, an area that he is currently researching. Eugene's particular interest is in the scalability of software-hardware architecture, data structures and relationality, data mining and real-time visualisation. On the modelling aspects of complex systems, Eugene's expertise is in advanced interactive systems, enhanced virtual environments, agent-based modelling and multi-agent systems that require large computing clusters for processing of agent-interaction and computer graphics. Dr Ch'ng is involved in editorial boards, technical, and programme committees in international journals and conferences in his field. Dr Ch'ng is a Council Member for the Complex Systems Society. Dr Eugene Ch'ng can be contacted at: [eugenecc@gmail.com](mailto:eugenecc@gmail.com)