

Cost-Effective Online Trending Topic Detection and Popularity Prediction in Microblogging

ZHONGCHEN MIAO and KAI CHEN, Shanghai Jiao Tong University

YI FANG, Santa Clara University

JIANHUA HE, Aston University

YI ZHOU and WENJUN ZHANG, Shanghai Jiao Tong University

HONGYUAN ZHA, Georgia Institute of Technology

Identifying topic trends on microblogging services such as Twitter and estimating those topics' future popularity have great academic and business value, especially when the operations can be done in real time. For any third party, however, capturing and processing such huge volumes of real-time data in microblogs are almost infeasible tasks, as there always exist API (Application Program Interface) request limits, monitoring and computing budgets, as well as timeliness requirements. To deal with these challenges, we propose a cost-effective system framework with algorithms that can automatically select a subset of representative users in microblogging networks in offline, under given cost constraints. Then the proposed system can online monitor and utilize only these selected users' real-time microposts to detect the overall trending topics and predict their future popularity among the whole microblogging network. Therefore, our proposed system framework is practical for real-time usage as it avoids the high cost in capturing and processing full real-time data, while not compromising detection and prediction performance under given cost constraints. Experiments with real microblogs dataset show that by tracking only 500 users out of 0.6 million users and processing no more than 30,000 microposts daily, about 92% trending topics could be detected and predicted by the proposed system and, on average, more than 10 hours earlier than they appear in official trends lists.

CCS Concepts: • **Information systems** → **Data stream mining**; **Document topic models**; **Retrieval efficiency**; *Content analysis and feature selection*; *Information extraction*; Social networks;

Additional Key Words and Phrases: Topic detection, prediction, microblogging, cost

ACM Reference Format:

Zhongchen Miao, Kai Chen, Yi Fang, Jianhua He, Yi Zhou, Wenjun Zhang, and Hongyuan Zha. 2016. Cost-effective online trending topic detection and popularity prediction in microblogging. *ACM Trans. Inf. Syst.* 35, 3, Article 18 (December 2016), 36 pages.

DOI: <http://dx.doi.org/10.1145/3001833>

This work was supported in part by the National Key Research and Development Program of China (2016YFB1001003), National Natural Science Foundation of China (61521062, 61527804), Shanghai Science and Technology Committees of Scientific Research Project (Grants No. 14XD1402100 and No. 15JC1401700), and the 111 Program (B07022).

Authors' addresses: Z. Miao, K. Chen (corresponding author), Y. Zhou, and W. Zhang, Department of Electronic Engineering, Shanghai Jiao Tong University, China; emails: {miaozhongchen, kchen, zy_21th, zhangwenjun}@sjtu.edu.cn; Y. Fang, Department of Computer Engineering, Santa Clara University, U.S.A.; email: yfang@scu.edu; J. He, School of Engineering and Applied Science, Aston University, U.K.; email: j.he7@aston.ac.uk; H. Zha, School of Computational Science and Engineering, Georgia Institute of Technology, U.S.A.; email: zha@cc.gatech.edu.

Permission to make digital or hard copies of part or all of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies show this notice on the first page or initial screen of a display along with the full citation. Copyrights for components of this work owned by others than ACM must be honored. Abstracting with credit is permitted. To copy otherwise, to republish, to post on servers, to redistribute to lists, or to use any component of this work in other works requires prior specific permission and/or a fee. Permissions may be requested from Publications Dept., ACM, Inc., 2 Penn Plaza, Suite 701, New York, NY 10121-0701 USA, fax +1 (212) 869-0481, or permissions@acm.org.

© 2016 ACM 1046-8188/2016/12-ART18 \$15.00

DOI: <http://dx.doi.org/10.1145/3001833>

1. INTRODUCTION

Nowadays people's daily life across the world is closely tied to online social networks. Microblogging services (e.g., Twitter¹ and Weibo²), as one of the representative online social network services, provide a more convenient approach for everybody around the world to read news, deliver messages, and exchange opinions than traditional media such as TV or newspaper. So huge quantities of users tend to post microposts in microblogging services and talk about the things they just witnessed, the news they just heard, or the ideas they just thought. In our article, the topic of a micropost refers to a group of keywords (such as the name of items, news headline, or thesis of ideas) in its content. Those semantically related microposts that talk about the same items, news, or thoughts within a given time window are the set of microposts of that topic.

Commonly, microblogging social networks are filled with a large number of varied topics all the time. However, if one topic is suddenly mentioned or discussed by an unusual amount of microposts within a relatively short time period, that topic is "trending" in microblogging services. As microblogging becomes the earliest and fastest source of information, these trending topics shown on social networks are often referred to the breaking news or events that have or will have societal impact in our real lives, such as first-hand reports of natural or man-made disasters, leaking an excellent product, unexpected sports winners, and very controversial remarks/opinions.

Because of this, identifying trending topics in microblogging networks is receiving increasing interest among academic researchers as well as industries. Moreover, it will produce even more scientific, social, and commercial value if the trending topics can be detected in real time and those topics' future popularity can be predicted at early stages. For example, the early awareness of first-hand reports of disasters can give rescuers more priceless time to reach the incident site and help more victims. Taking another example, higher predicted popularity and longer lifetime of a "leaks of a specific product" trending topic can bode well for the product's future reputation and sales, so businesspeople can be prepared to increase inventory and production.

In fact, microblogging service providers themselves, such as Twitter and Weibo, are publishing their official trending topic lists regularly. However unfortunately, these official lists are commonly delayed in publishing, small in size (Top-10 only), and not fully customized by individual user's preferences. They also do not contain topics' future popularity prediction function at all, so it is hard to tell how long a trending topic will last. More critically, there are concerns that some trending topics will never appear in these official lists that are subjected to the service provider's commercial considerations or even government censorship policy [Chen et al. 2013a]. If we rely only on the official trends lists, then some topics would most likely be missed or delayed by us. Therefore, business companies, organizations, or even individuals are in need of a reliable online real-time trending topics detection and prediction system on microblogging services and other social networks, which can produce impartial, accurate, and even customized results from a third-party perspective.

Traditionally, online trending topics detection and prediction systems for microblogging comprise three major steps: (1) retrieving microposts and related information from a microblogging website as much as possible, (2) detecting trends from the obtained microblog dataset, and (3) predicting the detected topics' future popularity using the obtained microblog dataset. In this way, the performance in steps 2 and 3 largely depends on the quantity and quality of the dataset retrieved in step 1. If at some time periods the sampled data source is biased [Morstatter et al. 2013, 2014], then an extra large-scaled

¹<http://twitter.com>., the top microblogging service worldwide.

²<http://weibo.com>., the top microblogging service in China.

dataset at those time periods will be needed in order to remove the bias and get representative topic detection results for all times. However, for any third-party analyzers, the tremendous number of users in microblogging services and the ever-growing volumes of microposts pose significant challenges to the capturing and processing of such a big scale dataset *in real time*. Although we are in an era of cloud-based services, it is still very challenging for any third-party analyzers to acquire the full real-time data stream of the whole microblogging network in time, as microblogging services companies are heavily limiting and narrowing the API (Application Program Interface) requests rate per account or per IP (Internet Protocol) address to prevent large dataset collection.³ Moreover, to get online detection results in real time, it requires a large resource budget on network bandwidth, IP address, storage, CPU (Central Processing Unit), and RAM (Random-access Memory) in cloud-based services for collecting and processing these large scaled data in time. As a result, in practical usage, the cost to obtain the fresh data and to detect and predict trending topics in real time should be seriously considered, and how to make a full use of the limited budget becomes a very important problem.

To deal with this difficulty, in this article we propose a cost-effective detection and prediction framework for trending topics in microblogging services. The core notion of the framework is to select a small subset of representative users among the whole microblog users in offline based on historical data. Then, online, the system will continuously track this small-sized subset of representative users and utilize their real-time microposts to detect the trending topics and predict these topics' future popularity. Therefore, our proposed system can run under limited resources, which sharply reduces data retrieval and computation cost and not compromise on performance.

The idea of selecting a subset of users in a microblogging network for trending topics detection and prediction is somewhat similar to the question of putting alerting sensors in city electricity or water monitoring networks that are analyzed by Leskovec et al. [2007], in which any single point power failure in the electricity monitoring network can be covered by a nearby alerting sensor. For microblogs, a topic can be viewed to be covered by a user if he posts⁴ or re-posts⁵ a micropost that is related to that topic. Thus, both problems aim to decide where to put the "sensors" in the network, given constraints on monitoring cost. However, electricity or water monitoring is a single-coverage outbreak detection system, which means any abnormal signal detected by one sensor should be reported as an issue. In contrast, when a new topic appears in the microblogging network and it is covered by one or a few users simultaneously, it should not be treated as a trending topic until a certain coverage degree is reached indicating the topic is really trendy among the whole network. Therefore, the placement of such "sensors," that is, selecting a proper subset of representative users in a microblogging network is a multi-coverage problem. And the selected users should be both effective in detecting trending topics and predicting those topic's future popularity.

It is worth pointing out that the representative users can be selected offline and fixed for a period of time in online usage. After some time, the set of selected users can be updated by running user selection algorithms again using the newly collected training data. However, as more microblog data are needed to run the user selection algorithm, in real-world usage the updating frequency need not be too high, or there will be no advantage in saving the data retrieving and processing costs.

The contributions of this article can be summarized as follows:

- (1) We treat online trending topics detection in microblogs as a multi-coverage problem: How to select a subset of users from all users in microblogging networks first, so

³See Twitter API Rate Limits as an example, <https://dev.twitter.com/rest/public/rate-limits>.

⁴A micropost is also called a "Tweet" in Twitter.

⁵Re-posting a micropost is similar to the "Forward" action in Email. It is also called a "ReTweet" in Twitter.

trending topics in the whole network can be detected by monitoring only this subset of users and utilizing their posted/re-posted microposts. In this way, the real-time monitoring and computation costs can be greatly reduced.

- (2) We formulate the subset user selection problem as a mixed-integer optimization problem with cost constraints, topic coverage, and prediction requirements. The topic coverage requirements can be customized for different individual topics or even different categories of topics, which enable the system to be more sensitive on high-priority topics that users are more interested in.
- (3) We integrate trending topics detection and their future popularity prediction into a single system. We propose efficient subset user selection algorithms for the optimization task by taking into account both detection and prediction accuracy. The experimental results show that the proposed algorithms outperform the state-of-art algorithms.
- (4) We collect nearly 1.6 million real-world microposts in Weibo as the testbed, and evaluate performance of the proposed system and algorithms from several dimensions. The real-time testing evaluations show that, using only 500 of 0.6 million users in the dataset, our proposed system can detect 92% of the trending topics that are published in Weibo official trends. Besides, it can also detect and predict the topics much earlier than they are published in the official trends. We also release our source code and the collected dataset to the public.⁶

2. RELATED WORKS

In Allan [2002], Fung et al. [2007], and Makkonen et al. [2004], a topic is defined as a coherent set of semantically related terms or documents that express a single argument. In this article, we follow the similar definitions, so microposts of one topic are those semantically related microposts/reposts that talk about the same items or news within a given time window. With the fast development of online services in recent years, detection and analysis of topics over microblogging services and other websites with user-generated contents are receiving more and more research interests.

One aspect of the research focuses on emerging topic discovery in online content, such as real-time earthquake detection over Twitter [Sakaki et al. 2010], “SigniTrend” emerging topic early detection with hashed significance thresholds [Schubert et al. 2014], real-time emergent topic detection in blogs [Alvanaki et al. 2012], “TwitterMonitor” trend detection system that treats bursting keywords as entry points [Mathioudakis and Koudas 2010], and two-level clustering methods [Petkos et al. 2014] that improve document-pivot algorithms in detection.

Some research papers track and analyse topics in longer time periods. Memes are identified on a daily basis by Leskovec et al. [2009], and their temporal variation is discussed by Yang and Leskovec [2011]. Cataldi et al. [2010] also use lifecycle models of key words to detect emerging topics. Event evolutions are mined with short-text streams by Huang et al. [2015].

Besides detecting emerging topics in social networks, some other works propose algorithms and techniques for analysing topic patterns and predicting trends online [Han et al. 2013]. Myers and Leskovec [2014] and Naaman et al. [2011] discuss factors that affect topic trends and the bursty dynamics in Twitter, and hashtags in microposts are utilized by Tsur and Rappoport [2012] for predicting topic propagation. Regression and classification algorithms are used by Asur et al. [2011] and Bandari et al. [2012] to predict news popularity in social media, temporal patterns evolution and state transition-based topic popularity prediction methods are discussed by Ahmed et al. [2013], and a Gradient Boosted Decision Tree model for microposts show counts is proposed by

⁶The dataset and source code is available at <https://github.com/zcmiao/Topic>.

Kupavskii et al. [2013]. There are also other purposes of topic analysis in social networks. For example, the event classification approach with utilization of spatio-temporal information carried by microposts is proposed by Lee et al. [2011], activity networks are used to identify interesting social events by Rozenshtein et al. [2014], and events trends are modeled with cascades of Poisson processing by Simma and Jordan [2010].

From all the above works, we see that various topic detection and analysis systems with different purposes, structures, and algorithms have been developed for social networks. However the above-reported systems need to process all data streams and extract features from it to accomplish these tasks. This will generate very heavy communication and computation loads, which requires large time and resource costs, hence its performance is restricted in real-time operations.

Our proposed online microblogging trending topics detection and popularity prediction system differs from the above reported systems in that our system tracks only a very small number of microblog users that are pre-selected by our algorithms and utilizes their real-time microposts to accomplish real-time detection and prediction tasks for trending topics in the whole microblogging network. One of the main contributions in this article is how to select most representative subset of users that are vital in both detection and prediction. The concept is somewhat similar to the influence maximization problem, which is to acquire maximum users or events cover under limited cost in social networks, which was first proposed by Domingos and Richardson [2001] and further discussed in Estevez et al. [2007], Narayanam and Narahari [2011], Pal and Counts [2011], and Weng et al. [2010]. The influence maximization problem is formulated as an optimization task by Kempe et al. [2003] and is proved to be NP-hard, and then a greedy algorithm is proposed to solve it approximately. A sub-modular property of the nodes selection problem in networks was found by Leskovec et al. [2007], and faster greedy algorithms were developed by Chen et al. [2009]. In our preliminary works [Chen et al. 2013b; Miao et al. 2015; Yang et al. 2014], the idea of selecting subset users for single tasks such as topic detection or topic prediction in microblogs was proposed. Some other greedy-based algorithms to get top-K influential users in social networks were proposed [Du et al. 2013; Wang et al. 2010], and an algorithm was proposed by Gomez-Rodriguez et al. [2012] to infer website influence in blogs. In addition, topic-specific influence and backbone structures in networks were studied [Bi et al. 2014; Bogdanov et al. 2013].

In this article, we extend the cost-effective framework and propose an integrated system for both trending topics detection as well as topic future popularity prediction in microblogs. Hence, subset users selection algorithm for joint detection and prediction are developed, and extensive experiments are carried out to evaluate joint multi-coverage and prediction performance under cost constraints.

3. OVERALL FRAMEWORK OF THE SYSTEM

In this section, we introduce the framework of our proposed system and then explain each module in detail. The overall system structure is shown as Figure 1, comprising the following five function modules:

- Module I: Training data retrieval;
- Module II: Subset microblog user selection;
- Module III: Real-time online data retrieval;
- Module IV: Online trending topic detection;
- Module V: Online trending topic popularity prediction.

In general, Module I and II run in offline, and they are mainly used for selecting representative users. Module I is used to obtain a historical training dataset including

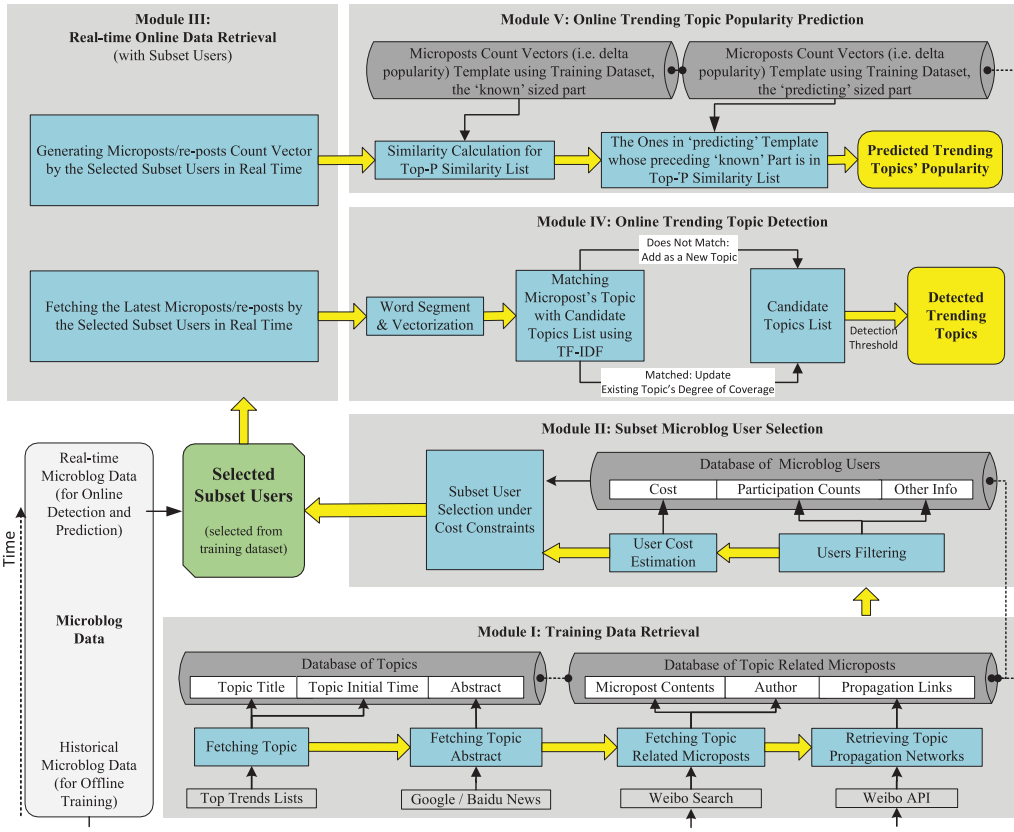


Fig. 1. Overall framework of our microblog trending topics detection and prediction system. Subset users are selected by Modules I and II using a training dataset, and the real-time microposts by these selected subset users are used for online detection and prediction in Modules III, IV, and V.

microposts, microposts' propagation (re-posting) links, and user profiles from microblogging websites. The “ground truth” trending topics are also collected in this module. Module II plays a role in selecting subset users from training dataset, and they should be optimally selected according to cost constraints and other configurable settings.

After selecting a subset of users offline, these users will be used in online modules, namely Modules III, IV, and V. Module III will continuously monitor only the selected users in real time and gather their fresh microposts/re-posts as data sources (these microposts/re-posts' further re-posting links are not gathered) for online trending topics detection in Module IV and prediction in Module V.

The above five modules work together to accomplish the overall online detection and prediction tasks under monitoring and computation cost constraints. In addition, the offline training can also be run periodically when a newly collected training dataset is ready, so the selected subset users can be updated for online operations. But please notice that the updating frequency need not be too high in order to save the cost for building up new training dataset.

3.1. Training Data Retrieval

In this subsection, we explain the components of Module I, so the building process of the training dataset, including historical microblog data gathering and several

pre-processing procedures, will be introduced. We use Weibo, the largest microblogging service in China, as a data source in our experiment. While most microblog contents in Weibo are written in Chinese, the proposed framework can be readily applied to other languages, with removal of some steps pertinent to the Chinese language such as Chinese word segmentation [Foo and Li 2004] during content vectorization.

3.1.1. Fetching Topics. Except for the microblogging service providers themselves, it is almost impossible to obtain datasets containing all the microposts and corresponding topics over the whole microblogging network. Therefore, a microblog dataset should be collected according to background knowledge of specific problem definitions and targets. As our research is focused on trending topics, the first thing we need to know is what topics are indeed popular over the Weibo microblogging network, so we could pay more emphasis on gathering these trending topic-related microposts.

Every 10 minutes we collected titles of the Top-10 Trends published officially by Weibo. To reduce the potential risk of commercial and political bias from Weibo Official Trends, we also collected titles of Top-10 Trends provided by some popular search engine companies in China, namely Baidu,⁷ Sogou,⁸ and Soso.⁹ Generally, the titles in all these Top Trends List are too short (commonly less than 20 Chinese characters) to describe the topic in detail. Therefore, we searched the titles in Google News¹⁰ and Baidu News¹¹ to get more textual information and keywords about the topic. They form an “abstract” of a trending topic, which contains around 80–160 Chinese characters or about 15–30 phrases on average. We make sure that the publishing time of these abstracts is consistent with the fetch time of the corresponding topic title. In addition, Term Frequency Inverted Document Frequency (TF-IDF) vectors of the title and the corresponding abstract are also compared, ensuring they discuss the same topic. Topics are also combined if they have large similarity on titles/abstracts within a given time period. Afterwards, the keywords (especially nouns, names, places, and verbs) from titles and abstracts can act as “descriptor” of trending topics that discriminate each topic clearly, so a topic will contain these keywords and the timestamp.

It is worth noting that the trending topics collected from these search engines may not completely eliminate the bias from Weibo Official Trends, due to the commercial considerations or government censorship policy in these sources themselves.

3.1.2. Fetching Topics Related Microposts and Their Propagation Networks. In a microblogging network, a user u_a can start a topic e' by posting a micropost, and the timestamp and keywords of the micropost content will also be the topic's timestamp and keywords. In fact, if that topic e' matches an existing topic e (i.e., e is started some times earlier by another user u_b) in keywords and in the same time period, then user u_a is actually joining the existing topic e unconsciously, even though u_a and u_b are in disjoint social relation/following networks (i.e., they do not know/follow each other). These posted (or to say non-reposted) microposts posted by u_a and u_b are then viewed as different *initializing microposts* of the same topic e , and there will be no topic e' .

Besides that, microblog user u_c can also join topic e by re-posting one of e 's existing microposts (including re-posts), and thus he also spreads the topic to his followers by his *re-posting micropost*. The re-posting action is actually one of the most effective ways that attracts users' interests on microblogs.

⁷<http://top.baidu.com>.

⁸<http://top.sogou.com>.

⁹<http://top.soso.com>, currently unavailable.

¹⁰<http://news.google.com/news?ned=cn>.

¹¹<http://news.baidu.com>.

We use different strategies in fetching these two kinds of microposts for the same topic e obtained in Section 3.1.1:

- In order to retrieve *initializing microposts* that are related to a specific topic e , keywords of the topic are used as query strings in Weibo Search API. In the returned results, a micropost is marked as related to that topic e if it meets all of the following rule: (1) it is not a re-posting micropost, (2) it matches the TF-IDF of that topic e 's keywords, (3) it is posted within a reasonable time window of topic e 's timestamp, and (4) it has a re-posts count larger than a certain threshold (e.g., 5) to speed up retrieval efficiency.
- For every *initializing microposts* fetched above, we recursively retrieved its full re-posting networks (including its re-posts and re-re-posts, etc.) using Weibo API. All *re-posting microposts* in a re-posting network discussing the same topic during a time period belong to that topic e .

Every author (microblog user) of topic e 's initializing and re-posting microposts is participating in e , and he or she can be regarded as nodes of topic e 's propagation network.

3.1.3. Topic Filtering. The trending topics fetched in Section 3.1.1 are crawled from the top trends lists provided by both Weibo and search engines. However, topic trends provided by search engines come from various kinds of information sources such as portal websites, blogs, forums, and microblogs. Therefore, if the number of participants/nodes in a topic's full propagation network is less than a threshold (e.g., 750), this topic seems to be not popular in microblogging services and will not be used. On the contrary, the rest topics with participants count bigger than the threshold are indeed trendy, and thus we regard these topics as "ground-truth" trending topics in our dataset.

Although we have tried our best to avoid the bias in the "ground-truth" trending topics and training dataset, the bias might still not be completely eliminated. So please note that there is a chance that the final detection results may also reflect such bias.

3.2. Subset Microblog User Selection

Module II is to select a suitable subset of representative users among all users in the whole trending topics propagation networks in offline, whose real-time posted/re-posted microposts can then be used to detect trending topics as well as to predict their popularity online in a cost-effective way. The whole user selection procedure comprises the following steps.

3.2.1. User Filtering. In microblogging networks, there are many inactive users and even spam users that should be excluded from selection, since efficiency is one major concern in our system. As this article is not focused on identifying spam users, we first apply some filtering rules on the domains of users. These filter rules remove the users who are highly inactive (far less than the average posting/re-posting frequency), apparently not influential (very low on followers count), or with spamlike behavior (such as repeatedly re-posting the same topic/micropost or putting many irrelevant keywords together into a single micropost). Filtering these users could reduce computation loads in later steps, and the final system accuracy will not likely to be affected, as these users are not likely to be selected anyway according to the strategies of all of the user selection algorithms mentioned in Section 6.2.

3.2.2. User Cost Estimation. When a user is selected as a representative user into the subset, the proposed system will keep monitoring and retrieving his/her microposts continuously, and then his/her microposts will be the input of real-time topic detection and prediction modules. The cost for monitoring and retrieving such real-time data

is related to the user's posting/re-posting frequency, as the number of API calls that fetches the microposts content is limited during each time window (e.g., per 15 minutes in Twitter) by the microblogging service provides. So the cost will arise if more API requests are needed to collect each selected user's data continuously. Besides that, the computational cost such as CPU and RAM needed for online detection and prediction algorithms are also related to input data scale or to, say, the selected user's posting frequency. Therefore, to quantitatively measure the system's monitoring and computational cost spent on each selected user, we define *user cost* as the average number of microposts that posted/re-posted by him/her per day during a long period of time. User cost will be taken into account during user selection for the sake of efficiency and system overall cost constraints.

Technically, we assume that user cost would not change much during a long period of time, thus *user cost* is estimated according to the time difference between the first and the last micropost among his/her latest 100 microposts (including re-posts). For example, if it takes a user 8 days to post his latest 100 microposts, then his cost is estimated as $100/8 = 12.5$. The number 100 here is enough, as we have tested the numbers much bigger than 100 and found no more than 10% difference on estimated costs. Moreover, due to the API rates limit per time window and the API limitations on the max number of one user's microposts that can be retrieved using one API call, it will cause extra consumptions on API resources to estimate user cost by retrieving much more than 100 latest microposts for each user, which is neither worthwhile nor affordable.

3.2.3. Subset User Selection. This is one of the core procedures in our system. An optimal subset of users are selected by minimizing detection and prediction loss while satisfying the system constraints. The formal problem definitions and solutions will be explained in detail in Section 4 and 5.

3.3. Real-Time Online Data Retrieval

In this module, previously selected subset users are monitored continuously online. The microposts that are posted/re-posted by these users within the latest time slot are periodically collected by our system using Weibo API, and thus selected users' microposts are gathered as real-time online dataset to be used in detection and prediction. It is notable that the further re-posting links or networks of these subset users' microposts/re-posts is not needed for the following real-time detection and prediction.¹²

3.4. Online Trending Topic Detection

Real-time microposts by the selected users that are collected in the previous module are fed into this module as the input dataset for trends detection. Generally, we can use almost any text-mining and trend identification methods with these data. Nevertheless, many research works focus on extracting features using a *huge* amount of data, and this is not suitable here since the input microposts data of this module is already downsized. Therefore, in order to meet the intention of our cost-effective framework and demonstrate the power of proposed subset user selection algorithms, we just apply a simple content matching-based single-pass clustering algorithm [Papka and Allan 1998; Yang et al. 1998] in this online detection module.

The online trending topic detection steps are outlined as follows, while the mathematical definition will be stated later, in Section 4.2.

¹²With the intention of evaluating the detection and prediction performance of our system, we still retrieved these microposts propagation links as ground truth. For the same reason, the "ground truth" trending topics of the real-time microblog dataset are also collected using similar methods mentioned in Module I.

- (1) Microposts that posted or re-posted by the subset users within the latest time slot are fetched periodically using Module III;
- (2) Word segmentation,¹³ stop-words filtering, and text vectorization are applied to micropost contents;
- (3) Each micropost is compared with the topic list that has been specified in the latest N_h (a configurable threshold) time slots using TF-IDF:
 - If a micropost is matched with an existing topic with high similarity, then mark the micropost to be related with that topic;
 - Otherwise, a new topic is created and added to the topic list whose timestamp and keywords are based on that micropost’s timestamp and its content keywords.
- (4) Update detection coverage of all the topics. If one topic’s detection coverage goes beyond a predefined threshold, then it is regarded to be detected as a trending topic.

It is worth pointing out that this detection module can be updated with more advanced text mining or any other types of detection methods that are compatible with our framework in accomplishing online trending topics detection task.

3.5. Online Trending Topic Popularity Prediction

After a trending topic is detected, our system can predict its future popularity. In this article, we define a topic’s *popularity* at a given time point as the total number of microposts and re-posts of that topic since the topic begun to that time point. Similarly to the considerations in choosing detection methods, we again propose a simple algorithm in terms of topic popularity prediction, whose formal definition and detailed method will be explained in Section 4.3 and 5.3. The basic idea is to calculate weighted average over template vectors as prediction results: At first, we calculate similarity between “known” part of a detected trending topic’ delta popularity vector (its size is τ) among selected users and each τ sized part of the template vector taken from training dataset. Then the “predicting” part of a trending topic’s delta popularity vector among all users can be predicted by weighted majority voting of the succeeding part of the top- \mathcal{P} most similar templates’ “known” part and other factors.

It is also worth noting that this prediction module can be updated with any other prediction algorithms that are compatible with our framework that uses subset users’ microposts to predict the topic’s future popularity among all users.

4. PROBLEM STATEMENT FOR SUBSET USER SELECTION

4.1. Basic Settings

Given a set of trending topics \mathcal{E} , the users who have posted/re-posted microposts for at least one trending topic in \mathcal{E} can be seen as the nodes of topic set \mathcal{E} ’s propagation network \mathcal{G} . Let \mathcal{V} denote the whole nodes set in the network; the goal is to select a suitable subset of representative nodes S from \mathcal{V} ($S \subseteq \mathcal{V}$), so trending topics \mathcal{E} among users \mathcal{V} can still be detected and their future popularity can be predicted using only the microposts from S .

There are two basic but necessary constraints when selecting subset nodes S : The maximum number of nodes (K) in the subset and the maximum total cost of all nodes (M) in the subset. The purpose of constraints K and M is to keep the real monitoring, data retrieving, and processing cost within budgets when solving practical problems.

¹³Chinese word segmentation system ICTCLAS is used, available at <http://ictclas.nipir.org>.

Denoting m_v as node v 's cost (defined in Section 3.2.2), the above two constraints can be represented by Equation (1),

$$|\mathcal{S}| \leq K, \quad \sum_{v \in \mathcal{S}} m_v \leq M. \quad (1)$$

4.2. Loss Function for Detection

This subsection formulates the loss function of trending topic detection in microblogs by selected subset users \mathcal{S} .

A node v ($v \in \mathcal{V}$) is regarded as a participant of a topic e ($e \in \mathcal{E}$) by posting or re-posting topic e -related microposts within a given time period T_M since topic e was initiated by its earliest micropost.¹⁴ And topic e is viewed to be covered for one time by node v if v participates in topic e . If node v participates in e for multiple times, then topic e is still viewed to be covered once by v . So binary variable $a_{v,e}$ is used to indicate this status, where $a_{v,e} = 1$ if and only if v participates in topic e at least once. Otherwise, the value of $a_{v,e}$ is 0.

As mentioned in Section 1, selecting subset users for trending topic detection is a multi-coverage ‘‘sensor placement’’ problem in a microblog propagation network. Therefore, we define a concept called *Degree of Coverage (DoC)*, denoted as $D_e(\mathcal{S})$, to measure the degree that a topic e has been covered by a subset of users \mathcal{S} ($\mathcal{S} \subseteq \mathcal{V}$). In the simplest form, $D_e(\mathcal{S})$ can be calculated by e 's participants count in \mathcal{S} , shown in Equation (2),

$$D_e(\mathcal{S}) = \sum_{v \in \mathcal{S}} a_{v,e}. \quad (2)$$

Given a threshold X_e , topic e is said to be multi-covered (or detected as a trending topic) by user set \mathcal{S} if and only if $D_e(\mathcal{S}) \geq X_e$. This detection threshold can be set accordingly for different training datasets and different cost constraints. Furthermore, the threshold for each topic X_e ($e \in \mathcal{E}$) or the threshold for topics in different categories can be customized according to the system user's preferences. For example, one topic containing a specific keyword can be set to have smaller detection threshold than the other topics, so it is easier for this topic to be covered as less users are needed.

The loss function for detecting trending topics \mathcal{E} using subset \mathcal{S} is shown as Equation (3). The value of function $\mathbf{1}(x)$ is equal to 1 if x is logical True, and it is equal to 0 if x is False. So there is no loss for a topic if its *DoC* reaches the detection threshold,

$$\mathcal{L}_{detect}(\mathcal{E}, \mathcal{S}) = \sum_{e \in \mathcal{E}} \mathcal{L}_{detect}(e, \mathcal{S}) = \sum_{e \in \mathcal{E}} \mathbf{1}(X_e > D_e(\mathcal{S})). \quad (3)$$

4.3. Loss Function for Prediction

Besides identifying e as a trending topic with subset user \mathcal{S} , we also would like to predict e 's future popularity among all users \mathcal{V} , using only the existing observed microblog data from subset user \mathcal{S} . With the predicted future popularity among all users, analyzers can understand the importance of the topic in advance, as well as how long will this trending topic last.

In this article, the popularity of a topic is measured by its total micropost (including re-posts) count at a given time point since the topic begun. For convenience, we segment a topic's whole lifetime T_M from when it is initiated until it ends into discrete time slots. These time slots can be indexed as $\{T_s^{(1)}, \dots, T_s^{(i)}, \dots, \}$, and the time points right

¹⁴ T_M is uniformly set to 3 days in our experiment settings.

between every time slot are denoted as $\{t_1, \dots, t_i, \dots | t_0 = 0\}$. Comparing with a topic's whole lifetime window T_M , each time slot length L_s should be set relatively small.¹⁵

We use a time series $\{y_e(1, \mathcal{V}), y_e(2, \mathcal{V}), \dots\}$ to represent topic e 's microposts (including re-posts) count that are posted/re-posted by users in \mathcal{V} during each time slot $T_s^{(1)}$, $T_s^{(2)}$ Thus, denoting the counting time series among all users \mathcal{V} since the first time slot $T_s^{(1)}$ until $T_s^{(\tau)}$ as $\vec{y}_e([1, \tau], \mathcal{V})$, the popularity of topic e at time point t_τ (i.e., right after time slot $T_s^{(\tau)}$) can be calculated by summing up its elements using Equation (4). As micropost counts are always non-negative, the sum can also be denoted as an L1 norm of \vec{y} ,

$$\text{pop}_e(t_\tau, \mathcal{V}) = \sum_{i=1}^{\tau} y_e(i, \mathcal{V}) = \|\vec{y}_e([1, \tau], \mathcal{V})\|_1. \quad (4)$$

Following the above definitions and design philosophy of our system, in real-time online prediction the actual microposts we observe are the ones posted or re-posted by subset user S from the beginning until t_τ , and the observed counting time series known to us can be denoted as $\vec{y}_e([1, \tau], S)$.¹⁶

We denote a prediction function Ψ in Equation (5) that can predict the future microposts counting time series $\vec{y}_e([\tau + 1, \kappa], \mathcal{V})$ among the whole users \mathcal{V} from time slot $T_s^{(\tau+1)}$ until $T_s^{(\kappa)}$, using input time series $\vec{y}_e([1, \tau], S)$. The value of κ indicates the longest time that can be predicted by function Ψ . Then the topic's future popularity at time point t_κ can be predicted using Equation (6) by summing the known counts until t_τ (the first term) as well as the predicted micropost counts from $t_{\tau+1}$ until t_κ (the second term) at each time slot,

$$\Psi(\kappa, \vec{y}_e([1, \tau], S)) = \hat{y}_e([\tau + 1, \kappa], \mathcal{V}), \quad (5)$$

$$\begin{aligned} \widehat{\text{pop}}_e(t_\kappa, \mathcal{V} | \Psi, \tau, S) &= \|\vec{y}_e([1, \tau], \mathcal{V})\|_1 + \|\hat{y}_e([\tau + 1, \kappa], \mathcal{V})\|_1 \\ &= \text{pop}_e(t_\tau, \mathcal{V}) + \|\Psi(\kappa, \vec{y}_e([1, \tau], S))\|_1. \end{aligned} \quad (6)$$

Having all the definitions above, the loss of popularity prediction on trending topics \mathcal{E} by a subset user S and a prediction function Ψ can be defined as the absolute popularity prediction error at time point t_κ ($\kappa > \tau$), shown in Equation (7).

$$\begin{aligned} \mathcal{L}_{\text{predict}}(\mathcal{E}, S) &= \sum_{e \in \mathcal{E}} \mathcal{L}_{\text{predict}}(e, S) \\ &= \sum_{e \in \mathcal{E}} |\widehat{\text{pop}}_e(t_\kappa, \mathcal{V} | \Psi, \tau, S) - \text{pop}_e(t_\kappa, \mathcal{V})|. \end{aligned} \quad (7)$$

Substituting the predicted popularity term in Equation (7) by Equation (5) and Equation (6), the loss can then be calculated by the sum of absolute micropost count prediction error in each time slot from time point $t_{\tau+1}$ until t_κ . The deduction is demonstrated

¹⁵Length of a time slot L_s is set to 6 minutes in our experiment settings.

¹⁶We would like to give an example for better illustration: Suppose there are 10 users in \mathcal{V} for a topic e . The first half of them each posted one micropost at the first time slot, and the other half of them each posted one micropost during the second time slot. Then the time series $\vec{y}_e([1 : 2], \mathcal{V})$ will be (5,5), and topic e 's popularity $\text{pop}_e(t_2, \mathcal{V})$ at time point t_2 is $5 + 5 = 10$. If 2 users in the first half are selected as subset users S , then the $\vec{y}_e([1 : 2], S)$ observed by the system will be (2,0), and $\text{pop}_e(t_2, S)$ is 2.

in Equation (8),

$$\begin{aligned}
 \mathcal{L}_{predict}(\mathcal{E}, \mathcal{S}) &= \sum_{e \in \mathcal{E}} |\widehat{\text{pop}}_e(t_\kappa, \mathcal{V} | \Psi, \tau, \mathcal{S}) - \text{pop}_e(t_\kappa, \mathcal{V})| \\
 &= \sum_{e \in \mathcal{E}} |\widehat{\text{pop}}_e(t_\kappa, \mathcal{V} | \Psi, \tau, \mathcal{S}) - (\text{pop}_e(t_\tau, \mathcal{V}) + \|\bar{y}_e([\tau + 1, \kappa], \mathcal{V})\|_1)| \\
 &= \sum_{e \in \mathcal{E}} |\|\hat{y}_e([\tau + 1, \kappa], \mathcal{V})\|_1 - \|\bar{y}_e([\tau + 1, \kappa], \mathcal{V})\|_1| \\
 &= \sum_{e \in \mathcal{E}} |\|\Psi(\kappa, \bar{y}_e([1, \tau], \mathcal{S}))\|_1 - \|\bar{y}_e([\tau + 1, \kappa], \mathcal{V})\|_1|.
 \end{aligned} \tag{8}$$

It should be pointed out that the time point of prediction function Ψ 's output is $t_{\tau+1}$ until t_κ given the input from t_1 until t_τ . If the input of Ψ are more recent observations such as $\bar{y}_e([1 + k, \tau + k], \mathcal{S})$ ($k > 0$), then it can produce prediction results $\hat{y}_e([\tau + 1 + k, \kappa + k], \mathcal{V})$ at further time points. In this way, the prediction results at any future time points can be recursively predicted.

4.4. Combined Objective Function

Based on the above loss functions, we formulate an optimization task for selecting a subset of nodes \mathcal{S} from the whole node set \mathcal{V} in network \mathcal{G} under resource constraints. In the optimization, argument is \mathcal{S} , and the target objective function is minimizing both detection loss $\mathcal{L}_{detect}(e, \mathcal{S})$ and prediction loss $\mathcal{L}_{predict}(e, \mathcal{S})$ for all topics $e \in \mathcal{E}$.

Let b_v be a binary variable where $b_v = 1$ indicates that node $v \in \mathcal{V}$ is selected as one of the subset users and $b_v = 0$ otherwise. Overall optimization objective function can then be represented in Equation (9) by mixing up Equation (1), Equation (3), and Equation (7). In the equation, λ is a coefficient that indicates the weight of prediction loss when selecting subset users \mathcal{S} . When $\lambda = 0$, the prediction loss will not be considered during user selection. The effect of λ in experiments will be discussed in Section 6.4.3,

$$\begin{aligned}
 \underset{\mathcal{S}}{\text{argmin}} \left(\sum_{e \in \mathcal{E}} \mathcal{L}_{detect}(e, \mathcal{S}) + \lambda \cdot \sum_{e \in \mathcal{E}} \mathcal{L}_{predict}(e, \mathcal{S}) \right) \\
 \text{s.t. } \sum_{v \in \mathcal{V}} m_v b_v \leq M, \quad \sum_{v \in \mathcal{V}} b_v \leq K, \quad \lambda \geq 0 \\
 \mathcal{S} = \{v | b_v = 1, v \in \mathcal{V}\}.
 \end{aligned} \tag{9}$$

5. EFFICIENT ALGORITHMS

Generally speaking, the original problem formulated in Section 4 is mixed-integer programming [Bertacco 2006], and we propose efficient algorithms to find a feasible solution that satisfies all constraints. For our joint detection and prediction system, we define a ‘‘reward’’ function $R(\Lambda, \Theta)$, which maps a subset Λ ($\Lambda \subseteq \mathcal{V}$) of nodes and a subset of topics Θ ($\Theta \subseteq \mathcal{E}$) into a real number. The value of this number shows the current detection and prediction ‘‘reward’’ on the topics set Θ using selected user subset Λ . Therefore, different ways of selecting subset users will lead to different sets of detected trending topics, and thus the rewards differ. We define the total joint reward in Equation (10), in which detection and prediction rewards are the derived and opposite

of loss function \mathcal{L}_{detect} and $\mathcal{L}_{predict}$, respectively.

$$\begin{aligned}
 R(\Lambda, \Theta) &= R_{detect}(\Lambda, \Theta) + \lambda \cdot R_{predict}(\Lambda, \Theta) \\
 R_{detect}(\Lambda, \Theta) &= \sum_{e \in \Theta} \mathbb{1}(D_e(\Lambda) \geq X_e) \\
 R_{predict}(\Lambda, \Theta) &= - \sum_{e \in \Theta} \left| \|\Psi(\kappa, \bar{y}_e([1, \tau], \Lambda))\|_1 - \|\bar{y}_e([\tau + 1, \kappa], \mathcal{V})\|_1 \right|.
 \end{aligned} \tag{10}$$

With the help of function R , various ways of utilizing reward values can be developed, that is, different heuristic strategies in selecting subset users. In following Section 5.1, we first introduce a straightforward user selection algorithm *SWC*, and then, in Section 5.2, a more effective user selection method *JNT* is proposed. Section 5.3 describes the popularity prediction algorithm and the prediction reward calculation in detail.

5.1. Algorithm SWC

In single-coverage problems where the objective is maximizing node placement coverage with nodes having equal or unequal cost, a widely used heuristic is the greedy algorithm described in Leskovec et al. [2007]. In that article, the node with maximized ratio of reward to cost is chosen iteratively in each round of selection. Based on the idea of maximizing ratio in that greedy algorithm, we adapted it to be compatible for solving subset user selection problems with multiple coverage requirements. This algorithm runs in a Stage-Wise Covering manner and thus is called algorithm *SWC*.

At first, the algorithm is initiated with an empty set of selected nodes $S = \emptyset$ and an empty set of topics $\mathcal{E}_c = \emptyset$ that includes the topics with $DoC \geq X$ (i.e., the trending topics that are multi-covered by S). Then the multi-covering problem can split into looping single-covering stages. During each single-coverage stage, every uncovered topic e ($e \in \mathcal{E} \setminus \mathcal{E}_c$) needs only to be covered once. In subsequent single-coverage stages, topic e still needs to be covered one time in each stage until its overall DoC reaches X_e and is moved into \mathcal{E}_c . In total, there will be at most $\max(X_e), e \in \mathcal{E}$ single-coverage stages.

More specifically, at the initiation step of the i th single-coverage stage, $\mathcal{E}_c^{(i)}$ (denoting the topics that has been single-covered in i th stage) is set to be empty; the detection threshold $X_e^{(i)}$ of each not-yet-multi-covered topic e 's ($e \in \mathcal{E} \setminus \mathcal{E}_c$) is set to 1 in this stage, and the threshold $X^{(i)}$ of the rest topics in $\mathcal{E}_c \cup \mathcal{E}_c^{(i)}$ are set to $+\infty$, indicating that the reward for these topics is not considered. The optimization target of this stage is to find a subset of nodes that can single-cover topics set as $\mathcal{E} \setminus \mathcal{E}_c$.

For each single-coverage stage, users are iteratively selected in rounds. In each user selection round, marginal detection reward/cost ratio of each user v ($v \in \mathcal{V} \setminus S$) is calculated with Equation (11). A user v_{\max} with the largest marginal detection reward per unit cost is then selected and added to subset S . Afterwards, the topics covered by user v_{\max} are added to $\mathcal{E}_c^{(i)}$, and the marginal detection reward/cost ratio is recalculated using Equation (11) again. Then, in next the round, another user with the largest marginal reward/cost ratio is selected. In this way, users can be iteratively selected for each single-coverage stage. Each single-coverage stage stops when all topics are single covered (i.e., $\mathcal{E}_c^{(i)} = \mathcal{E} \setminus \mathcal{E}_c$) or when the overall cost constraints are reached. At the end of i th stage, \mathcal{E}_c is updated. If the overall cost constraints are not reached, then the $i + 1$ th stage will then begin. In case there are more than one user maximizing Equation (11), we can select the user who participates in more topics to break the tie,

$$v_{\max} = \operatorname{argmax}_{v \in \mathcal{V} \setminus S} \frac{R_{detect}(\mathcal{S}^{(i)} \cup \{v\}, \mathcal{E} \setminus (\mathcal{E}_c \cup \mathcal{E}_c^{(i)})) - R_{detect}(\mathcal{S}^{(i)}, \mathcal{E} \setminus (\mathcal{E}_c \cup \mathcal{E}_c^{(i)}))}{m_v}. \tag{11}$$

ALGORITHM 1: Algorithm *SWC* for Subset User Selection

Require: Full nodes set \mathcal{V} , nodes cost $m_v (v \in \mathcal{V})$, trending topics \mathcal{E} , constraints M and K

Ensure: A set of optimal selected nodes $\mathcal{S} \subseteq \mathcal{V}$

```

1:  $\mathcal{S} \leftarrow \emptyset, \mathcal{E}_c \leftarrow \emptyset, M_{curr} \leftarrow 0, i \leftarrow 1$ 
2: while  $|\mathcal{S}| < K$  and  $M_{curr} < M$  and  $i \leq \max\{X_e | e \in \mathcal{E}\}$  do
3:    $\mathcal{E}_c^{(i)} \leftarrow \emptyset, \mathcal{S}^{(i)} \leftarrow \emptyset, X_e^{(i)} = 1 (e \in \mathcal{E} \setminus \mathcal{E}_c), X_e^{(i)} = +\infty (e \in \mathcal{E}_c)$ 
4:   while  $\mathcal{E}_c^{(i)} \neq \mathcal{E} \setminus \mathcal{E}_c$  do
5:     Calculate current reward  $R_{detect}(\mathcal{S}^{(i)}, \mathcal{E} \setminus (\mathcal{E}_c \cup \mathcal{E}_c^{(i)}))$  by Equation (10)
6:     Find a node  $v_{max} \in \mathcal{V} \setminus \mathcal{S}$  with max reward/cost ratio by Equation (11)
7:     if  $M_{curr} + m_{v_{max}} \leq M$  and  $|\mathcal{S}| + 1 \leq K$  then
8:        $\mathcal{E}_c^{(i)} \leftarrow \mathcal{E}_c^{(i)} \cup \{e | a_{v_{max}, e} = 1, e \in \mathcal{E} \setminus \mathcal{E}_c\}$ 
9:        $\mathcal{S}^{(i)} \leftarrow \mathcal{S}^{(i)} \cup \{v_{max}\}, \mathcal{S} \leftarrow \mathcal{S} \cup \{v_{max}\}$ 
10:       $M_{curr} \leftarrow M_{curr} + m_{v_{max}}$ 
11:     else
12:       Abort user selection, return  $\mathcal{S}$ 
13:     end if
14:   end while
15:    $\mathcal{E}_c \leftarrow \mathcal{E}_c \cup \{e | D_e(\mathcal{S}) \geq X_e, e \in \mathcal{E}^{(i)}\}, i \leftarrow i + 1$ 
16: end while
17: return  $\mathcal{S}$ 

```

After running all the stages in algorithm *SWC*, a subset of users \mathcal{S} are finally selected, and then those real-time microposts of \mathcal{S} will be retrieved and used in subsequent real-time detection and prediction procedures. Pseudo-code of the whole algorithm *SWC* is listed in Algorithm 1.

5.2. Algorithm *JNT*

In the user selection algorithm *SWC*, the target of topic coverage is set to 1 per single covering stage. It is not efficient, as it needs many loops when overall detection threshold X is large. As a matter of fact, when solving multi-coverage problems, it is more efficient to cover a topic more than once by different users during one user selection iteration. Additionally, algorithm *SWC* does not take the selected user's prediction performance into consideration at all, which might not be appropriate for the joint task. Therefore, we propose an efficient algorithm to solve the multi-coverage problem that takes into account the *JoiNT* detection and prediction accuracy of selected users. Thus we name it algorithm *JNT*, and it contains three major improvements compared with *SWC*.

The first improvement in algorithm *JNT* is that dynamic detection reward is used for different topics in user selection, based on the gap between each topic's current *DoC* and its detection threshold X . In the original reward function (Equation (10)), subset users' detection reward for each topic is binary valued depending on whether its current *DoC* reaches detection threshold X , which is fine in single-coverage situations. However, in multi-coverage problems, the reward should be measured more precisely as the gap between a topic's current *DoC* and X could differ considerably among various topics and users. For example, suppose the detection threshold X is 10 in a user selection process, the current *DoC* of trending topics e_1 and e_2 are 2 and 8 respectively, user u_1 can cover e_1 once while u_2 can cover e_2 once, and one of them is to be selected. In this situation, other things being equal, topic e_1 is more urgent to be covered than e_2 , because e_1 needs 8 more coverages to reach threshold X while e_2 needs only 2. Thus, the reward for covering e_1 by u_1 should be higher than covering e_2 by u_2 so u_1 can be selected. However, the binary valued detection reward cannot handle this case as threshold X is not reached and the reward for u_1 and u_2 would be both 0.

Therefore, to improve the overall topic coverage of S , a dynamic reward function is defined according to the difference between a topic's current DoC and its threshold X . If topic e 's current DoC has not reached detection threshold yet, then it is urgent to cover e by the selected user, so the reward for covering e can be defined to be proportional (linear) to the difference between X and its DoC (i.e., $X_e - D_e(S)$). Consequently, topics with lower coverage degree are prioritized to have higher reward, and thus they are stimulated to be covered by the selected users in subsequent user selections. In contrast, when topic e 's DoC has reached X , it is not urgent to cover e any more, so the reward is set to be inversely proportional to its DoC to discourage further covers on this topic in subsequent subset user selection operations.

Denoting Λ as the set of already selected subset users, the dynamic reward $r_e(\Lambda)$ for topic e is denoted by Equation (12), based on the above settings. In Equation (12), reward is commonly no larger than 1, and α ($0 \leq \alpha \leq 1$) is a configurable number to control sensitivity level in dynamic reward calculation. α is set to 0.01 in our experiments, and it is discussed further in Section 6.4.4,

$$r_e(\Lambda) = \begin{cases} \frac{(1-\alpha)[X_e - D_e(\Lambda)]}{X_e}, & D_e(\Lambda) < X_e \\ \frac{\alpha}{D_e(\Lambda)}, & D_e(\Lambda) \geq X_e \end{cases}. \quad (12)$$

Hence, the definition of detection reward for algorithm JNT should be updated accordingly, shown in Equation (13). Theoretically speaking, using dynamic reward in user selection can be helpful in improving the overall recall rate of trending topics detection,

$$R_{detect}^{JNT}(\Lambda, \Theta) = \sum_{e \in \Theta} r_e(\Lambda). \quad (13)$$

The second improvement in algorithm JNT is that we apply a dynamic user cost boundary in user selection, so the users whose cost is beyond boundary are excluded from selection. Afterwards, the users having maximum reward and who are within cost boundaries are selected iteratively as subset users. At the beginning of each iteration, the cost boundary is dynamically updated according to current system spare cost and current subset users' size. Comparing to the strategy of just selecting the user with the highest marginal reward/cost ratio in algorithm SWC , the aforementioned operation is a more flexible user selection strategy that can make full use of the remaining available cost budget, especially when the total cost constraints is not so tight.

Concretely, when there are K_l available nodes (maximum is K) and M_l available microposts monitoring and processing cost (maximum is M), the cost boundary $M_b(K_l, M_l)$ is proposed in Equation (14). In the equation, γ is a configurable value to control boundary size. The cost boundary will always be bigger than the current average available cost per user (M_l/K_l), so the users with better coverages but relatively larger cost are allowed to be selected; the boundary will be no larger than the current total available cost M_l in order to meet the system cost constraints. γ is set to 0.7 in our experiments and is discussed in Section 6.4.4,

$$M_b(K_l, M_l) = \min\left(\frac{M_l}{K_l^\gamma} + \frac{M_l}{K_l}, M_l\right). \quad (14)$$

Third, in algorithm JNT , we consider the selected users' prediction reward during user selection. For algorithm SWC , users who have the best marginal detection reward per unit cost are selected. However, the fact that a user is doing well in detection does not necessarily mean that he or she will also be the best choice in prediction. For example, the detection result will not be affected if most of the selected users prefer to

ALGORITHM 2: Algorithm *JNT* for Subset User Selection**Require:** Full nodes set \mathcal{V} , nodes cost $m_v (v \in \mathcal{V})$, trending topics \mathcal{E} , constraints M and K **Ensure:** A set of optimal selected nodes $\mathcal{S} \subseteq \mathcal{V}$

```

1:  $\mathcal{S} \leftarrow \emptyset, M_l \leftarrow M, K_l \leftarrow K$ 
2: while  $K_l > 0$  and  $M_l > 0$  do
3:   Calculate cost boundary  $M_b(K_l, M_l)$  by Equation (14)
4:   Calculate current joint reward  $R^{JNT}(\mathcal{S}, \mathcal{E})$  by Equation (15)
5:   Find a node  $v_{\max} \in \mathcal{V} \setminus \mathcal{S}$  with max joint reward increment by Equation (16)
6:   if  $m_{v_{\max}} \leq M_l$  then
7:      $\mathcal{S} \leftarrow \mathcal{S} \cup \{v_{\max}\}, K_l \leftarrow K_l - 1, M_l \leftarrow M_l - m_{v_{\max}}$ 
8:   else
9:     Abort user selection, return  $\mathcal{S}$ 
10:  end if
11: end while
12: return  $\mathcal{S}$ 

```

attend trending topics relatively later than the other users, but their microposts might be too late to be used as prediction input and the prediction result may not be so ideal.

Therefore, prediction reward is added into the total reward function R^{JNT} for algorithm *JNT* shown in Equation (15). In the equation, coefficient λ controls the prediction reward weight.¹⁷ By default, λ is bigger than 0 in *JNT*, so both the detection and prediction will be taken into account when selecting users, and it is discussed in Section 6.4.3,

$$R^{JNT}(\Lambda, \Theta) = R_{detect}^{JNT}(\Lambda, \Theta) + \lambda \cdot R_{predict}(\Lambda, \Theta). \quad (15)$$

Combining the above three improvements and modifications together, algorithm *JNT* selects the subset users in an iterative manner: In each selecting iteration, cost boundary M_b is updated based on current available budget, then reward of each user in $\mathcal{V} \setminus \mathcal{S}$ whose cost is within current cost boundary is calculated, and then user v_{\max} with maximized total reward among them is selected using Equation (16). After adding v_{\max} into \mathcal{S} and updating K_l and M_l , cost boundary M_b is re-calculated and then the next user selection iteration begins. The user selection process will stop when any cost constraints is met. The full procedures of algorithm *JNT* is summarized in Algorithm 2,

$$v_{\max} = \underset{v \in \mathcal{V} \setminus \mathcal{S}, m_v \leq M_b(K_l, M_l)}{\operatorname{argmax}} (R^{JNT}(\mathcal{S} \cup \{v\}, \mathcal{E}) - R^{JNT}(\mathcal{S}, \mathcal{E})). \quad (16)$$

5.3. Prediction Algorithm

In the above two subsections, we introduced different user selection algorithms and corresponding reward calculation methods. In this subsection, we explain the algorithm of utilizing selected subset users' microposts for predicting topic's future popularity among whole users. After that, selected users' prediction reward as well as topic popularity prediction result can be calculated.

The intention and definition of prediction function Ψ is already listed in Equation (5). To mathematically describe the detailed algorithm for the prediction function, we follow the setting introduced in Section 4.3: When a subset of users \mathcal{S} is selected, what we can observe by monitoring their microposts are the beginning *known part* of each topic e 's microposts/re-posts counting time series $\bar{y}_e([1, \tau], \mathcal{S})$ until time point t_τ . The prediction target is the succeeding *unknown part* $\bar{y}_e([\tau + 1, \kappa], \mathcal{V})$ among all users \mathcal{V} , from time

¹⁷If the detection and prediction reward are not normalized, then the weight coefficient should be denoted as $\lambda_s \cdot \lambda$, where λ_s is the scale factor. In this article, we assume $\lambda_s = 1$.

ALGORITHM 3: Prediction Algorithm with Selected Subset Users

Require: Training dataset Ω_{tra} , user subset \mathcal{S} , observed $YG_e^{\mathcal{S}}$ with length τ
Ensure: Predicted time series $YU_e^{\mathcal{V}}$ and $YU_e^{\mathcal{S}}$ from $t_{\tau+1}$ until t_κ

- 1: **for all** topic's full counting time series in Ω_{tra} **do**
- 2: Generate template vectors $\mathbf{PG}^{\mathcal{V}}$ and $\mathbf{PU}^{\mathcal{V}}$ using sliding window
- 3: Generate template vectors $\mathbf{PG}^{\mathcal{S}}$ and $\mathbf{PU}^{\mathcal{S}}$ using sliding window
- 4: **end for**
- 5: **for all** template vector $PG_j^{\mathcal{S}} \in \mathbf{PG}^{\mathcal{S}}$ **do**
- 6: Calculate similarity between $PG_j^{\mathcal{S}}$ and $YG_e^{\mathcal{S}}$
- 7: **end for**
- 8: Get index i_ρ ($\rho \in [1, \mathcal{P}]$) of Top- \mathcal{P} most similar template vector with $YG_e^{\mathcal{S}}$ in $\mathbf{PG}^{\mathcal{S}}$
- 9: Predict $YU_e^{\mathcal{V}}$ with Equations (17), (18), and (19)
- 10: Predict $YU_e^{\mathcal{S}}$ with Equation (20)

point $t_{\tau+1}$ to t_κ . That is to say, we will use current subset users' microposts of a topic to predict the topic's future popularity among all users. For simplicity, the aforementioned known part and predicting part are denoted as $YG_e^{\mathcal{S}}$ and $YU_e^{\mathcal{V}}$, respectively.

5.3.1. Template Vectors. In our prediction algorithm, there is an assumption that if the first part of two time-series vectors are high in similarity, their succeeding part within a small time period will also likely to be similar, especially when the two vectors represent the same group of users' posting behavior on trending topics. Therefore, besides the observed *known* vectors $YG_e^{\mathcal{S}}$, additional *template* time-series vectors that can reflect subset users \mathcal{S} 's posting/re-posting counts on historical trending topics are needed. Each template vector consists of two parts: the τ sized known part used for similarity calculation and the succeeding $\kappa - \tau$ -sized *predicting* part used for prediction.

For a given selected user subset \mathcal{S} , a set of template vectors $\mathbf{P}^{\mathcal{S}}$ can be extracted from the training dataset. In the training dataset, each trending topic's full counting time series has $L_M = \lfloor T_M/L_s \rfloor$ (Max lifetime of a topic/length of a time slot) time slots in total, and it is commonly much bigger than κ . Thus, we use a sliding window with size = κ and step = 1 to extract every κ sized template vector from each trending topic's full counting time series in the training dataset. After that, each template vector $P_j^{\mathcal{S}} \in \mathbf{P}^{\mathcal{S}}$ is segmented into known and succeeding predicting parts as $\langle PG_j^{\mathcal{S}}, PU_j^{\mathcal{S}} \rangle$, with size $(\tau, \kappa - \tau)$.

Concretely, denoting $\vec{y}_A([1, L_M], \mathcal{S})$ as topic A 's full counting time series among users set \mathcal{S} in the training dataset, the first extracted template will be $PG_1^{\mathcal{S}} = \vec{y}_A([1, \tau], \mathcal{S})$, $PU_1^{\mathcal{S}} = \vec{y}_A([\tau + 1, \kappa], \mathcal{S})$. Then, the window will slide one step and the second extracted template will be $PG_2^{\mathcal{S}} = \vec{y}_A([2, \tau + 1], \mathcal{S})$, $PU_2^{\mathcal{S}} = \vec{y}_A([\tau + 2, \kappa + 1], \mathcal{S})$, and so on. The extraction window will gradually slide for $L_M - \kappa$ times until the other side of the window reaches the last time slot of \vec{y}_A , and thus $L_M - \kappa + 1$ templates are extracted. Afterwards, extraction of the next trending topic B 's full counting time series $\vec{y}_B([1, L_M], \mathcal{S})$ begins, and so on. In the end, the set $\mathbf{P}^{\mathcal{S}}$ will include all the template vectors extracted from all trending topics in the training dataset.

Using similar operations, template vectors set $\mathbf{P}^{\mathcal{V}}$ that contains the microposts counting time series among whole user set \mathcal{V} is also built. Moreover, if $\mathbf{P}^{\mathcal{V}}$ and $\mathbf{P}^{\mathcal{S}}$ are extracted from the training dataset exactly in the same order, they have synchronized indexes on topics and time points as users set $\mathcal{S} \subseteq \mathcal{V}$.

5.3.2. Similarity Calculation and Popularity Prediction. After building up template vectors with the offline training dataset, it is time to predict $YU_e^{\mathcal{V}}$ for topic e . At first,

similarities between YG_e^S and every template $PG_j^S \in \mathbf{PG}^S$ are calculated using Pearson's correlation coefficient. If all the elements in YG_e^S or PG_j^S are the same value, then the Pearson's coefficient cannot be calculated. In this case, we use a cosine correlation coefficient instead to represent the similarity. The template whose similarity with YG_e^S is less than a threshold (0.5 in our experiments) will be skipped in following steps.

Denoting i_ρ ($\rho \in [1, \mathcal{P}]$) as indexes of the top- \mathcal{P} most similar template vectors $PG_{i_\rho}^S$ to YG_e^S , then those top- \mathcal{P} template vectors' succeeding part $PU_{i_\rho}^S$ and their same-indexed template vectors $PU_{i_\rho}^\mathcal{V}$ among all users \mathcal{V} can be used to predict the time series $YU_e^\mathcal{V}$ among all users \mathcal{V} . The prediction calculation is done by the weighted average of the template vectors $PU_{i_\rho}^\mathcal{V}$ ($\rho \in [1, \mathcal{P}]$) shown in Equation (17), in which w_{i_ρ} is the weight of the i_ρ th template $PU_{i_\rho}^\mathcal{V}$ and is measured by the similarity between $YU_e^\mathcal{V}$ and $PU_{i_\rho}^\mathcal{V}$. Also, η_{i_ρ} in Equation (17) is a scale coefficient that measures the scale ratio between template vector $PU_{i_\rho}^\mathcal{V}$ and predicting target $YU_e^\mathcal{V}$, so the scale of the predicted result will not be affected by the scale of template vectors, which commonly differ.

According to the assumption that similarity of two vectors would not change much within a short time period, as well as the fact that subset users will be selected according to their representativeness among all users, we can estimate the similarity between $YU_e^\mathcal{V}$ and $PU_{i_\rho}^\mathcal{V}$ by using the previously calculated similarity between YG_e^S and $PG_{i_\rho}^S$, shown in Equation (18). Besides that, as YG^S and \mathbf{PG}^S are observed by the same group of users, their scale ratio can also be utilized to estimate the scale ratio η of their succeeding parts among all users \mathcal{V} . This estimation is stated in Equation (19),

$$\widehat{YU}_e^\mathcal{V} = \frac{\sum_{\rho \in [1, \mathcal{P}]} w_{i_\rho} \cdot \eta_{i_\rho} \cdot PU_{i_\rho}^\mathcal{V}}{\sum_{\rho \in [1, \mathcal{P}]} w_{i_\rho}}, \quad (17)$$

$$w_{i_\rho} = \text{sim}(YU_e^\mathcal{V}, PU_{i_\rho}^\mathcal{V}) \approx \text{sim}(YG_e^\mathcal{V}, PG_{i_\rho}^\mathcal{V}) \approx \text{sim}(YG_e^S, PG_{i_\rho}^S), \quad (18)$$

$$\eta_{i_\rho} = \eta(YU_e^\mathcal{V}, PU_{i_\rho}^\mathcal{V}) \approx \eta(YG_e^\mathcal{V}, PG_{i_\rho}^\mathcal{V}) \approx \eta(YG_e^S, PG_{i_\rho}^S) \approx \frac{\sum YG_e^S}{\sum PG_{i_\rho}^S}. \quad (19)$$

Additionally, the counting time series YU_e^S among subset users can also be predicted similarly using Equation (20),

$$\begin{aligned} \widehat{YU}_e^S &= \frac{\sum_{\rho \in [1, \mathcal{P}]} w_{i_\rho} \cdot \eta_{i_\rho} \cdot PU_{i_\rho}^S}{\sum_{\rho \in [1, \mathcal{P}]} w_{i_\rho}} \\ w_{i_\rho} &= \text{sim}(YU_e^S, PU_{i_\rho}^S) \approx \text{sim}(YG_e^S, PG_{i_\rho}^S) \\ \eta_{i_\rho} &= \eta(YU_e^S, PU_{i_\rho}^S) \approx \eta(YG_e^S, PG_{i_\rho}^S) \approx \frac{\sum YG_e^S}{\sum PG_{i_\rho}^S}. \end{aligned} \quad (20)$$

The overall prediction algorithm is summarized in Algorithm 3. Using this algorithm, $\{\hat{y}_e(\tau+1, \mathcal{V}), \hat{y}_e(\tau+2, \mathcal{V}), \dots, \hat{y}_e(\kappa, \mathcal{V})\}$ can be predicted, so the overall predicted popularity $\widehat{\text{pop}}_e(t_k, \mathcal{V})$ at any time point t_k , $k \in [\tau+1, \kappa]$ can be calculated with Equation (6).

Table I. Statistics of Training and Testing Dataset

Dataset	Time Period	Trending Topics	Microposts	Users
Ω_{tra}	10th Sept.–25th Sept., 2012	75	753,486	585,640
Ω_{test}	26th Sept.–10th Oct., 2012	93	840,572	634,840

For offline training (user selection) process, prediction reward $R_{predict}$ using the given selected subset users \mathcal{S} can also be calculated with Equation (10). Additionally, time series further than $\hat{y}_e(\kappa, \mathcal{V})$ can also be recursively estimated by inputting newer YG_e^S (either observed or previously predicted) at more recent time points into the prediction system.

It is worth noting that, during the whole prediction process, we never need or use YG^V , that is, the known part of trending topics' microposts count among *all* users.

6. EXPERIMENTS

6.1. Data Collections

We use Weibo as the microblog data source in our experiments. Weibo is the dominant microblogging service provider in China, which has more than 222 million monthly active users and 100 million daily active users¹⁸ as of September 2015.

Based on the procedures described in Section 3.1, we crawled the titles and abstracts of “ground truth” trending topics. We used the mentioned methods to retrieve microblog data in the period from September 10, 2012, to October 10, 2012. Meanwhile, we collected each topic's initializing microposts and their full reposting network in Weibo, where each topic was tracked for 3 days since it started. In total, there are 168 trending topics in the dataset containing 1,594,058 microposts/re-posts and 1,104,960 nodes (distinct users). We then split the topics and corresponding microposts into two disjoint parts for different purposes, whose statistics are listed in Table I.

- The first part contains the topics that were initiated during the first 15 days, denoted as \mathcal{E}_{tra} . These topics' corresponding microposts are treated as training dataset Ω_{tra} . Given cost constraints K and M , the operations of Modules I and II in our framework (see Section 3) and the proposed efficient algorithms are applied to Ω_{tra} , and thus the subset users \mathcal{S} will be selected from all users in set \mathcal{V}_{tra} who participated in \mathcal{E}_{tra} .
- The rest of the trending topics \mathcal{E}_{test} initiated during the last 15 days are used for testing, and all the microposts of topics \mathcal{E}_{test} are regarded as a full testing dataset Ω_{test} to simulate the real-time microposts exist in Weibo network. For our system, only the microposts/re-posts $\Omega_{test}^S \subseteq \Omega_{test}$ that are generated by the selected subset users $\mathcal{S} \subseteq \mathcal{V}_{tra}$ will be used in real-time testing. The rest of the dataset $\Omega_{test} \setminus \Omega_{test}^S$ is kept untouched during online detection and prediction, and it is only used as ground truth in the final prediction performance evaluation.

It is worth noting that in a real online environment, it is almost impossible for any third-party analyzers to crawl and collect all of the newly generated microposts Ω_{test} *in real time* because of the fact that microblog service providers are limiting API usage, as well as the high expense incurred in gathering and processing the full-sized fresh data to fulfill the time requirements. However, in our system, the needed testing dataset Ω_{test}^S is small in size, which can be easily picked up by Module III of our system with a small amount of Weibo API requests. Then the small-sized testing dataset can be used to conduct the detection and prediction tasks described in Modules IV and V of our system framework.

¹⁸Weibo Official Reports: <http://ir.weibo.com/phoenix.zhtml?c=253076&p=irol-newsArticle&ID=2113781>.

Table II. Statistics of Followers For Microblog Users in Ω_{tra}

Followers No.	>2M	500k–2M	50k–500k	5k–50k	1k–5k	<1k
User Amount	17	338	1,580	7,363	32,536	the rest

Table III. Statistics of Trending Topic Participation Counts per Microblog User in Ω_{tra}

Participated Topics Count	1	2	3	4–10	≥ 10
User Amount	487,776	65,033	18,032	14,057	742
Percentage	83.29%	11.10%	3.08%	2.40%	0.13%

Table IV. Parameters Configuration in Offline Training

Parameters Set	M	K	X
Set I	8,000	100	10
Set II	15,000	200	20
Set III	30,000	500	40
Set IV	50,000	800	60

To illustrate the basic characteristics of microblog dataset, some statistical analysis on training dataset Ω_{tra} are carried out, and it can be helpful in deciding some threshold configurations in data pre-processing and user selection algorithms. Distributions of microblog users' followers counts are shown in Table II. From the table, we can see that only 1.6% of users have more than 5,000 followers, and nearly 93% users have fewer than 1,000 followers. In terms of per-user's participation counts for trending topics in the training dataset, Table III shows that only 2.53% of users were observed participating in ≥ 4 different trending topics.

From the above statistics, we can observe a long-tail phenomenon in the microblogging network. Therefore, before running subset user selection algorithms (including the proposed algorithm *SWC*, *JNT*, and all the other baseline user selection algorithms that are introduced in the next section) that compare all users in the training dataset, we use pre-processing filtering, mentioned in Section 3.2.1, to remove the inactive users and users who exhibit spamlike behaviours. Thus the user selection process can be more efficient.

6.2. Evaluation Criteria

In this section, we will first introduce the methods that are used for comparisons with the proposed algorithms, and then explain the criteria for performance evaluation.

According to statistical analysis in the previous section, there are two straightforward strategies for subset user selection problem. One strategy is iteratively picking the user who has the most followers in the training dataset, which can be denoted as algorithm *FM*. The other is called algorithm *ECM*, which iteratively picks a user who has the highest topic participation count. Besides that, we also use PageRank [Brin and Page 2012] as another baseline method *PR* in selecting subset users among all users in the training dataset. In algorithm *PR*, the users who were involved in topics \mathcal{E}_{tra} and the re-posting actions between those users are treated as nodes and edges of a directed multi-graph. Then the nodes with the highest PageRank values in the graph are selected as subset users. User selection operations in the above three methods will stop once the system cost constraints are reached.

In terms of system training parameter configurations (including cost constraints settings), for each offline user selection algorithm *FM*, *ECM*, *PR*, *SWC*, and *JNT*, we run four sets of parameters I through IV, listed in Table IV. Constraints of maximal

microposts monitoring and processing cost M and maximal selected subset users size K are applied to all five user selection algorithms. Detection threshold X is applied when running algorithms SWC and JNT , and the same threshold is used for each topic. In other words, under identical cost constraints and parameters, our experiments will use the training dataset Ω_{tra} to select five different subsets of users \mathcal{S} from \mathcal{V}_{tra} using algorithms FM , ECM , PR , SWC , and JNT , respectively. After selecting a subset users \mathcal{S} by each algorithm offline, real-time detection and prediction performances on topics \mathcal{E}_{test} are evaluated using the corresponding real-time testing dataset Ω_{test}^S .

In general, the value of training parameter X can be set according to subset users size K as well as the desired quality of the selected users, since it can be used to control the least-desired average DoC per subset user (denoted as d) on \mathcal{E}_{tra} . For example, in parameter set II, if we want to have averagely at least $d = 8$ trending topics covered per subset user on \mathcal{E}_{tra} , the corresponding X can be estimated by $K * d / |\mathcal{E}_{tra}| = 200 * 8 / 75 = 21.33 \approx 20$. In experiments, the value of d should be set based on the dataset characteristics, especially the statistics of each user's participation counts on \mathcal{E}_{tra} shown in Table III. If d is set too small (e.g., $d < 4$), then a huge number of users who have lower trending topic participation counts (see Table III) could be selected into the subset, and thus the subset users' overall coverage on trending topics and the training quality will not be ideal. If d is too big (e.g., $d \geq 10$), then the amount of users fulfilling the coverage requirement could be quite small (also See Table III), and thus the proposed algorithm will be somewhat similar to the strategy used in algorithm ECM , that is, only selecting the users with largest participation counts.

In addition to the methods that use only subset users' microposts as data sources for real-time online trending topic detection, we also run an experiment using a state-of-art detection method called "Two-level clustering" [Petkos et al. 2014] on the Weibo testing dataset. The algorithm is a document-pivot algorithm and is denoted as TLC . It scans the contents and other features of all microposts in the full dataset Ω_{test} , puts them into different clusters, and then extracts the top-ranked topics from the clusters. Algorithm TLC has no training or user selection procedures (all users are selected, that is, $\mathcal{S} = \mathcal{V}$ and $\Omega_{test}^S = \Omega_{test}$ in this case), so in order to detect real-time trending topics online with algorithm TLC , the full micropost dataset by all users must be obtained and processed in real time, which is quite expensive in the online environment.

In online evaluation, a trending topic e is viewed as detected if its DoC reaches or exceeds an online detecting threshold \tilde{X}_e using microposts posted by subset user \mathcal{S} . It is noted that this online detection threshold \tilde{X} in real-time testing is generally not equal to the threshold X used in offline subset user selection, because the scales of the datasets Ω_{tra} and Ω_{test}^S differ substantially. \tilde{X} could also be set differently for each topic according to user preferences on its content, so a topic with lower \tilde{X} can be detected more easily in online usage. Generally speaking, during a real-time trending topic detection process, any topic whose DoC reaches the threshold \tilde{X} will be identified as a trending topic by our system, so some of the topics that are not in the "ground truth" topic list \mathcal{E}_{test} may also be detected. Denoting $\hat{\mathcal{E}}_{test}$ as all the trending topics detected by Ω_{test}^S with $DoC \geq \tilde{X}$, the recall and precision that quantitatively measure trending topics detection performance can be defined to benchmark the results.

Detection recall rate is calculated by Equation (21), that is, the ratio of unique correctly matched trending topics' size to the ground-truth trending topics' size. In the equation, function $\mathbb{1}(x)$ equals 1 if x is logical True or 0 otherwise,

$$\text{recall}(\Omega_{test}^S) = \frac{\sum_{e \in \mathcal{E}_{test}} \mathbb{1}(e \in \hat{\mathcal{E}}_{test})}{|\mathcal{E}_{test}|}. \quad (21)$$

The detection precision rate is calculated in Equation (22), that is, the ratio of the number of total correctly matched trending topics to the number of detected trending topics. So if several trending topics in $\hat{\mathcal{E}}_{test}$ match the same trending topic in \mathcal{E}_{test} , they will be counted multiple times in a precision calculation as follows:

$$\text{precision}(\Omega_{test}^S) = \frac{\sum_{e \in \hat{\mathcal{E}}_{test}} \mathbb{1}(e \in \mathcal{E}_{test})}{|\hat{\mathcal{E}}_{test}|}. \quad (22)$$

The above precision value can be denoted as precision@All (or P@All), as it is calculated based on all detected trending topics in $\hat{\mathcal{E}}_{test}$ with Equation (22). Sometimes, however, in order to put emphasis on the most trending topics in microblogs, only the Top- N topics (ranked by their *DoC*) in $\hat{\mathcal{E}}_{test}$ will be regarded as the detected trending topics. In this case, precision@ N (or P@ N) is reported by treating only Top- N topics as $\hat{\mathcal{E}}_{test}$ in Equation (22).

It should be pointed out that the above recall and precision are based on the “ground-truth” topics \mathcal{E}_{test} gathered from Weibo and other search engines’ Top-10 Trends. But, actually, there are always some microposts discussing topics other than \mathcal{E}_{test} in the dataset Ω_{test}^S , especially in users’ comments when re-posting. That is to say, due to the incompleteness of real ground-truth topic lists for Ω_{test}^S , a recall metric might be more convincing than precision in this evaluation scenario. Therefore, both F1-score ($\beta = 1$) and F2-score ($\beta = 2$) are calculated with Equation (23) to benchmark the detection performance, in which $\beta > 1$ means the F-score relies more on recall than precision,

$$F_\beta = (1 + \beta^2) \cdot \frac{\text{precision} \cdot \text{recall}}{\beta^2 \cdot \text{precision} + \text{recall}}. \quad (23)$$

In terms of prediction, the predicted trending topic’s popularity is measured by the commonly used Root Mean Square Error (RMSE) in a topicwise manner. Let us denote $\hat{y}_e(k, \mathcal{V})$ as the predicted microposts count of topic e during a future time slot $T_s^{(k)}$ among all users \mathcal{V} , which is predicted by using real-time microposts dataset Ω_{test}^S from selected users \mathcal{S} . Prediction RMSE of this time slot for all topics can then be calculated by Equation (24) on the topic $e \in \mathcal{E}_{test} \cap \hat{\mathcal{E}}_{test}$ that belongs to the ground-truth topics and is detected by our system,

$$\text{RMSE}(T_s^{(k)}, \mathcal{V}|\mathcal{S}) = \sqrt{\frac{\sum_e (\hat{y}_e(k, \mathcal{V}) - y_e(k, \mathcal{V}))^2}{|\mathcal{E}_{test} \cap \hat{\mathcal{E}}_{test}|}}. \quad (24)$$

In addition, in order to compare prediction result during a larger period of time between time slot $T_s^{(k_a)}$ and $T_s^{(k_b)}$, the mean RMSE per time slot is commonly used in later experiments using Equation (25),

$$\overline{\text{RMSE}}(T_s^{(k_a)}, T_s^{(k_b)}) = \frac{\sum_{k=k_a}^{k_b} \text{RMSE}(T_s^{(k)}, \mathcal{V}|\mathcal{S})}{|k_b - k_a + 1|}. \quad (25)$$

6.3. Topic Coverage Evaluation on Training Dataset

Before evaluating system performance on a real-time testing dataset, in this section we will first exhibit the selected subset users’ multi-covering performance with different user selection algorithms on topics \mathcal{E}_{tra} in the training dataset. That is to say, right

Table V. Topic Coverage Comparison of Training Dataset with Selected Users

Set ^a	Alg.	K_S	$Cost_S$	Covered Topics Ratio				A-DoC	Run-Time
				$\bar{X} = 1$	$\bar{X} = 3$	$\bar{X} = 5$	$\bar{X} = 8$		
I	FM	100	3,631	50/75	27/75	15/75	11/75	3.2	0.5
	ECM	43	8,291	66/75	55/75	44/75	32/75	10.6	0.7
	PR	100	3,180	53/75	32/75	23/75	14/75	4.4	1.1
	SWC	100	32.9	74/75	20/75	8/75	2/75	2.5	32
	JNT ($\lambda = 0$)	100	6,246	73/75	72/75	70/75	67/75	16.5	16
	JNT ($\lambda = 0.5$)	100	4,925	73/75	73/75	71/75	62/75	15.5	233
	JNT ($\lambda = 1.0$)	100	5,119	73/75	72/75	69/75	58/75	14.0	265
II	FM	200	13,392	54/75	36/75	25/75	19/75	5.8	0.4
	ECM	116	15,030	67/75	63/75	57/75	52/75	22.9	0.8
	PR	200	6,181	63/75	48/75	38/75	25/75	9.1	0.6
	SWC	200	104.8	74/75	74/75	26/75	11/75	5.3	63
	JNT ($\lambda = 0$)	200	11,788	73/75	73/75	72/75	72/75	29.9	15
	JNT ($\lambda = 0.5$)	200	12,140	73/75	73/75	72/75	70/75	29.5	534
	JNT ($\lambda = 1.0$)	200	10,454	73/75	73/75	72/75	70/75	27.8	539
III	FM	500	29,752	67/75	53/75	45/75	37/75	14.9	0.4
	ECM	380	30,029	73/75	69/75	67/75	64/75	53.6	0.8
	PR	500	14,718	73/75	67/75	58/75	50/75	23.6	0.6
	SWC	500	456.4	74/75	74/75	74/75	74/75	12.9	152
	JNT ($\lambda = 0$)	500	23,979	73/75	73/75	73/75	72/75	56.4	36
	JNT ($\lambda = 0.5$)	500	23,358	73/75	73/75	73/75	72/75	56.0	1607
	JNT ($\lambda = 1.0$)	500	22,891	73/75	73/75	72/75	72/75	54.9	1539
IV	FM	800	39,209	72/75	64/75	59/75	46/75	22.8	0.5
	ECM	800	49,200	74/75	70/75	67/75	66/75	88.4	0.8
	PR	800	25,491	73/75	70/75	66/75	61/75	38.8	0.6
	SWC	800	920.5	74/75	74/75	74/75	74/75	20.6	224
	SWC ($K = 1500$)	1500	2,452	74/75	74/75	74/75	74/75	38.4	398
	SWC ($K = 2000$)	2000	4,050	74/75	74/75	74/75	74/75	52.7	574
	SWC ($K = 3000$)	3000	7,051	74/75	74/75	74/75	74/75	80.6	826
	JNT ($\lambda = 0$)	800	34,288	74/75	74/75	73/75	72/75	79.1	116
	JNT ($\lambda = 0.5$)	800	34,458	74/75	74/75	73/75	72/75	78.8	4311
JNT ($\lambda = 1.0$)	800	34,281	74/75	73/75	73/75	72/75	79.8	4248	

^aThe parameters and constraints used in each set is listed in Table IV.

after selecting subset users S by each offline user selection algorithm, we exhibit the detection performance on trending topics set \mathcal{E}_{tra} in the training dataset using the corresponding subset microposts Ω_{tra}^S .

During the subset user selection process, algorithms *FM*, *ECM*, *PR*, and *SWC* do not consider each user's prediction reward at all. In terms of *JNT*, it will consider prediction loss unless its coefficient λ is set to 0. So performances with some different λ values ranging from 0 to 1 are evaluated for algorithm *JNT* for more detailed comparison. Additionally, as algorithm *SWC* tends to select users with the highest reward/cost ratio, its total cost of the selected users (denoted as $Cost_S$) using parameter sets I through IV is too low to be comparable with the other algorithms. Therefore, we run additional experiments for *SWC* by enlarging the subset user size constraints K to 1500, 2000, and 3000 in parameter set IV, thus more users can be selected and included in S . Algorithm *TLC* does not contain training or a user selection process, so it is not included in the subset users' detection performance comparison on the training dataset.

Table V shows the selected users performance on the training dataset with different user selection algorithms, using the four sets of parameters mentioned in Table IV. In

the table, column “ K_S ” shows the size of subset users that are actually selected, and column “ Cost_S ” shows the total estimated cost $\sum_{v \in S} m_v$ of these selected subset users. The values of these two columns are related to the system cost constraints K and M ; Column “Covered Topics” represents the proportion of topics in \mathcal{E}_{tra} whose DoC reaches \bar{X}_c using Ω_{tra}^S , and the results with several different threshold \bar{X}_c values are reported; Column “A-DoC” shows the actual average degree of coverage of topics in \mathcal{E}_{tra} , that is, $\sum_{e \in \mathcal{E}_{tra}} D_e(S)/|\mathcal{E}_{tra}|$. Generally, higher average DoC means better overall trending topic coverage performance by the selected subset users S ; Column “Run-Time” lists the running time of the offline user selection procedures in minutes.¹⁹

Let us first focus on the comparison on topic coverage and selected user’s cost. From Table V, it is easy to find out that algorithm FM has the lowest number of covered topics as well as the second lowest average DoC among the five user selection algorithms under the same parameter settings. It suggests that when some cost constraints are considered, following only the microblog users with largest followers is not a good strategy in covering more trending topics. This result may be a little bit beyond one’s intuition, since in the real world we are more likely to follow the users with more followers, whom are often the celebrities or famous ones, to receive fresh news and information. In terms of ECM , as its philosophy is to select users who participate more in trending topics, apparently it has a higher average DoC and relatively larger covered topic count than FM , PR , and SWC . However, ECM is also the highest in Cost_S and lowest in K_S among all algorithms, which means that ECM is selecting the users with larger average cost per user. Therefore, it is probably not as cost-effective as other algorithms, especially when the total cost budget is tight and thus fewer users are selected. In contrast, per-user cost for algorithm SWC is the lowest, as it is designed to be, but its average DoC is also the lowest and its topic coverage is not as ideal as other algorithms when K_S is small. To improve topic coverage performance, SWC has to select two times or more users than the other algorithms, which diverges our initial intention of small sized but representative S . Moreover, continuously monitoring too many users will also cause extra consumption on API requests, which is strictly restricted by the microblogging company. For algorithm PR , its topic coverage performance is worse than ECM but better than FM , while per-user cost is relatively fair. For algorithm JNT with $\lambda = 0$ or $\lambda > 0$, it outperforms all the other algorithms in topic coverage ratio and average DoC under identical constraints conditions, and its selected users’ cost is moderately small. The covered topics and average DoC of JNT ($\lambda > 0$) is a little bit smaller than JNT ($\lambda = 0$), but the latter has slightly less Cost_S . More discussions on λ value are covered in Sections 6.4.3.

Next is the discussion on the running time needed for selecting subset users. It can be found from column “Run-Time” in Table V that algorithm JNT with $\lambda > 0$ runs is slower than all the other algorithms. But, in our opinion, the longer training time in algorithm JNT ($\lambda > 0$) is acceptable for the following three reasons: (1) The whole user selection process is an offline running procedure, so the time requirement is much less urgent, and overall better detection performance is more preferred. Besides, it takes more time to run algorithm JNT ($\lambda > 0$), as it takes into account prediction reward while other algorithms do not. (2) During each user selection iteration, each user’s reward on detection and prediction will be computed, and then all the results are compared. Thus distributed processing techniques such as MapReduce [Dean and Ghemawat 2008] can be further applied to speed up the current training time. (3) Due to the rate limits on API usage by the microblog service providers, it is too expensive for third-party users to collect all of the newly generated full training dataset within

¹⁹We run the training experiments in a single virtual machine with 16 cores and 60G RAM from Google Computing Engine. Multiple cores are used when running algorithm JNT .

a short period of time, so, in practice, the training dataset itself would not be updated very frequently. As a consequence, it is not necessary to re-run offline user selection algorithms and to update the subset users very often if there is no new training dataset. Based on the above three reasons, we think it is fine to take longer time in running the user selection algorithm, so better detection and prediction accuracy can be achieved.

Last in this subsection, the detection threshold \tilde{X} is discussed. In Table V, it can be observed that trending topics' coverage changes with respect to detection threshold \tilde{X} . Generally speaking, the detection threshold should be determined by the detectability of trending topics. For our system, this lies in finding a proper value of online detection threshold \tilde{X} , which is important in deciding whether a topic can be regarded as a trending topic, as well as evaluating recall and precision with Equation (21) and Equation (22) in the later real-time experiments with the testing dataset. Thus, empirically in the article, we set $\tilde{X} = 3$ as default online detection threshold value based on the above evaluations over training dataset, which could allow at least 95% ($\geq 72/75$) topics in training dataset to be marked as trending topics, using JNT ($\lambda \geq 0$) with the lowest cost constraints parameter set I. It should be pointed out that in later online evaluations, $\tilde{X} = 3$ is also globally used for parameters sets I through IV for comparison convenience across different constraints settings. However, in practical usage, \tilde{X} can be set larger than 3 accordingly when a selected users' size is bigger.

6.4. Evaluation on Real-Time Testing Dataset

After benchmarking topic coverage with selected subset users on the training dataset, the proposed system performance is then evaluated on the real-time online testing dataset. With the subset users \mathcal{S} that are selected from the offline training dataset Ω_{tra} by different user selection algorithms, we use only their microposts in the testing dataset $\Omega_{test}^{\mathcal{S}} \subseteq \Omega_{test}$ to detect and predict the trending topics \mathcal{E}_{test} online. For prediction, we set t_{τ} and t_{κ} to be 6 hours and 30 hours after each topic is initiated. This means we observe the first $t_{\tau} = 6$ hours of a topic's counting time series using only the selected users microposts $\Omega_{test}^{\mathcal{S}}$ and then predict its popularity among all users \mathcal{V} in the next $t_{\kappa} - t_{\tau} = 24$ hours.

Table VI shows the real-time online testing performance results of FM , ECM , PR , SWC , and JNT using the corresponding dataset $\Omega_{test}^{\mathcal{S}}$ by their selected subset users \mathcal{S} . In the table, column "Recall" reports the recall rate using Equation (21). In order to exhibit more detailed results, recall rates with various online detection thresholds \tilde{X} are reported. In the columns "Precision," "F1-Score," and "F2-Score," as the purpose of using the same online detection threshold \tilde{X} value for different parameters set in these evaluations is already explained at the end of the previous subsection, the evaluation results of Equation (22) and Equation (23) with $\tilde{X} = 3$ are listed. In some cases, people may prefer to compare the precision rate based on the top- N detected trending topics in $\hat{\mathcal{E}}_{test}$. Thus the precision@ N ($N = 30$ and 50) as well as corresponding F1-/F2-scores ($N = 50$) are also listed in the performance comparison table. Column "A-D" in Table VI is the average degree of coverages for topics in \mathcal{E}_{test} , and a higher value means that the trending topics are still likely to be detected even if threshold \tilde{X} is set to be higher; Column "T-G" shows the average time difference *in hours* of trending topics detected by our system *before* they are published by the Weibo and Baidu/Sogou/Soso search engines' Top Trend Lists, in which the posting time of the last micropost that makes topic e 's *Degree of Coverage* reach detection threshold \tilde{X}_e is regarded as the time point that topic e is detected as a trending topic by our system. In other words, column "T-G" reflects the time gained by using our system than relying on the official trends list to get trending topics. Column "RMSE" in the table reports the average popularity predicting

Table VI. Real-Time Online Performance over Testing Dataset Ω_{Test}^S with Selected Users

Set	Alg ^c	Recall					Precision ($\bar{X}=3$)					F1 ($\bar{X}=3$)		A-D	T-G ^b	RMSE	R-T ^v
		X = 1	X = 3	X = 5	X = 8	P@All	P@30	P@50	P@All	P@50	P@All	P@50					
I	FM	27/93	27/93	14/93	3/93	28/31	27/30	28/31	0.4394	0.4394	0.3359	1.5	1.0	58.33	3.5		
		66/93	66/93	61/93	42/93	88/108	27/30	44/50	0.7857	0.7857	0.7383	9.7	12.6	58.18	3.7		
	ECM	23/93	23/93	15/93	5/93	27/32	26/30	27/32	0.3825	0.3825	0.2880	1.7	1.0	57.67	3.2		
		4/93	4/93	2/93	0/93	4/4	4/4	4/4	0.0825	0.0825	0.0532	0.2	-0.4	60.71	3.2		
	SWC	75/93	75/93	66/93	52/93	103/121	28/30	46/50	0.8282	0.8282	0.8150	11.1	19.7	41.59	3.6		
		JNT ($\lambda = 0$)	71/93	71/93	61/93	37/93	93/112	27/30	41/50	0.7955	0.7907	0.7741	8.3	15.7	40.90	3.6	
	JNT ($\lambda = 1.0$)	71/93	71/93	57/93	39/93	89/105	28/30	45/50	0.8033	0.8261	0.7873	7.9	14.0	40.94	3.6		
		FM	33/93	33/93	22/93	9/93	36/41	26/30	36/41	0.5054	0.5028	0.4028	2.7	1.5	47.08	3.6	
	ECM	79/93	79/93	75/93	61/93	138/171	27/30	44/50	0.8277	0.8645	0.8406	20.8	19.4	47.54	4.1		
		PR	47/93	47/93	31/93	16/93	56/64	27/30	43/50	0.6407	0.6366	0.5508	4.1	7.1	56.12	3.4	
II	SWC	14/93	14/93	5/93	1/93	15/16	15/16	15/16	0.2594	0.2594	0.1809	0.8	0.8	45.55	3.3		
		80/93	80/93	75/93	63/93	125/154	28/30	45/50	0.8352	0.8797	0.8501	18.2	22.4	42.52	3.9		
JNT ($\lambda = 0$)	83/93	83/93	76/93	61/93	128/159	29/30	44/50	0.8465	0.8862	0.8735	18.0	23.3	42.94	3.8			
	JNT ($\lambda = 0.5$)	80/93	80/93	71/93	59/93	118/144	27/30	45/50	0.8393	0.8797	0.8517	15.9	21.2	43.10	3.8		
JNT ($\lambda = 1.0$)	FM	62/93	62/93	42/93	31/93	77/89	27/30	44/50	0.7531	0.7586	0.6987	7.8	9.4	40.57	3.7		
	ECM	83/93	83/93	82/93	78/93	176/236	28/30	45/50	0.8125	0.8962	0.8587	46.8	29.7	52.70	4.6		
III	PR	67/93	67/93	53/93	43/93	105/133	26/30	43/50	0.7534	0.7841	0.7333	11.9	19.0	41.66	3.6		
		36/93	36/93	17/93	8/93	42/44	29/30	42/44	0.5508	0.5508	0.4393	2.4	0.6	42.68	3.4		
SWC	85/93	85/93	83/93	77/93	161/207	29/30	44/50	0.8404	0.8967	0.8831	33.8	26.3	47.75	4.1			
	JNT ($\lambda = 0$)	85/93	85/93	84/93	78/93	162/209	28/30	45/50	0.8388	0.8824	0.9111	33.5	27.9	46.79	4.1		
JNT ($\lambda = 0.5$)	86/93	86/93	82/93	76/93	158/203	27/30	44/50	0.8452	0.9018	0.8912	32.2	30.1	44.88	4.1			
	JNT ($\lambda = 1.0$)	74/93	74/93	55/93	42/93	104/128	28/30	43/50	0.8040	0.8266	0.7990	13.0	14.8	42.44	3.9		
IV	ECM	87/93	87/93	85/93	84/93	198/266	28/30	45/50	0.8290	0.9174	0.8898	76.4	39.0	65.85	5.0		
		PR	75/93	75/93	66/93	55/93	135/174	28/30	45/50	0.7909	0.8507	0.8001	20.2	25.2	41.07	3.9	
SWC	50/93	50/93	32/93	14/93	63/69	29/30	44/50	0.6768	0.6675	0.5858	4.2	7.4	42.18	4.4			
	64/93	64/93	47/93	30/93	82/98	29/30	44/50	0.7724	0.7724	0.7135	7.9	9.3	41.31	3.7			
SWC ($K = 1.5k$)	74/93	74/93	59/93	38/93	99/122	28/30	44/50	0.8035	0.8357	0.7988	11.2	13.8	40.64	3.7			
	80/93	80/93	68/93	53/93	121/154	28/30	44/50	0.8213	0.8700	0.8442	18.3	17.4	98.51	3.9			
SWC ($K = 2k$)	87/93	87/93	85/93	80/93	191/241	29/30	46/50	0.8581	0.9277	0.9029	47.0	37.3	49.51	4.4			
	JNT ($\lambda = 0$)	86/93	86/93	83/93	78/93	192/241	28/30	45/50	0.8560	0.9122	0.8959	47.0	33.2	43.53	4.4		
JNT ($\lambda = 0.5$)	88/93	88/93	85/93	77/93	189/242	28/30	44/50	0.8557	0.9119	0.9078	46.5	31.3	47.80	4.5			
	JNT ($\lambda = 1.0$)	88/93	88/93	85/93	77/93	189/242	28/30	44/50	0.8557	0.9119	0.9078	46.5	31.3	47.80	4.5		

^aPlease also refer to the corresponding selected users' cost ($Cost_S$) listed in Table V as well as Figure 2 to compare the estimated online cost of each algorithm.

^bThe unit is in hours. It shows the time gained by using our system than relying on the official trends list to get trending topics.

^cThe unit is in minutes. It shows the total running time of both online detection and prediction procedures for algorithms with user selection mechanism, which is much faster than the running time of algorithm *TLC* listed in Table VII.

Table VII. Real-time Online Performance over Testing Dataset Ω_{test} using Algorithm *TLC*

Threshold	$\bar{X} = 50$	$\bar{X} = 100$	$\bar{X} = 300$	$\bar{X} = 500$	$\bar{X} = 800$	$\bar{X} = 1000$
Recall	90/93	88/93	83/93	78/93	64/93	59/93
Precision@All	999/1058	755/796	404/430	281/297	200/210	170/177
Precision@30	30/30	30/30	30/30	30/30	30/30	30/30
Precision@50	50/50	50/50	50/50	50/50	50/50	50/50
F1-Score (P@All)	0.9559	0.9474	0.9154	0.8892	0.7990	0.7641
F1-Score (P@50)	0.9836	0.9724	0.9432	0.9123	0.8153	0.7763
F2-Score (P@All)	0.9629	0.9467	0.9015	0.8582	0.7286	0.6806
F2-Score (P@50)	0.9740	0.9565	0.9121	0.8667	0.7339	0.6845
Average DoC	7,890	7,697	7,024	6,507	5,963	5,673
Running Time	527.5					

RMSE per time slot for the next 24 hours using Equation (25). Column “R-T” shows the total running time of both online detection and prediction procedures in minutes.²⁰

Besides evaluating online testing performance of the above five algorithms that uses a small subset of dataset Ω_{test}^S , we also evaluate detection performance of algorithm *TLC* that needs to use and process the full testing dataset Ω_{test} during online detection. That is to say, in order to detect trending topics with algorithm *TLC* in real-time, all users’ microposts in the microblog website must be gathered quickly and continuously, so all these microposts can then be processed for clustering and topic extraction in near real time. In a practical online environment, the size of microblog users and their newly generated microposts are extremely huge, thus the cost of collecting and processing such full-sized dataset in real time is prohibitive. The detection performance of algorithm *TLC* using Ω_{test} and its running time (in minutes, without the data collecting time) is listed in Table VII, in which the detected topics containing no fewer than \bar{X} microposts are treated as trending topics. As the size of input dataset Ω_{test} (shown in Table I) is much bigger than any Ω_{test}^S , detection performance with various online detection thresholds \bar{X} much bigger than 3 are reported.

In the following subsections, performances by different algorithms with different parameter settings are compared and discussed in detail.

6.4.1. Discussions on Performance of Different Algorithms. Viewing the online performance of all five user selection algorithms *FM*, *ECM*, *PR*, *SWC*, and *JNT* under each cost parameter Set I through IV as a whole in Table VI, the F-scores, average *DoC*, and detection time gain apparently increase as value of cost constraints M and K increase from Set I to Set IV. This shows that system cost budget settings are indeed affecting online testing cost as well as the overall system online performance.

At first, we discuss the detection performance with algorithm *TLC*, for which the system cost is not limited at all and all the testing dataset are used in online testing. According to the performance shown in Table VII, the F-scores of algorithm *TLC* are indeed better than *JNT* and other algorithms (shown in Table VI) when *TLC*’s $\bar{X} < 300$. But please keep in mind that the former algorithm uses microposts from 0.6 million users, and the latter uses only microposts from no more than 800 selected users. More importantly, the running time of online testing for *TLC* is also increasing as the size of its input dataset needed is apparently larger, and it takes more than 115 times (527.5 min vs. 4.5 min) longer than any other algorithms with the user selection mechanism. That is to say, although the detection performance of algorithm *TLC* is the best, it is not practically suitable to accomplish a real-time online trending topics detection task

²⁰In all online testing experiments, a commodity computer with 2.0GHz CPU and 16G RAM is used.

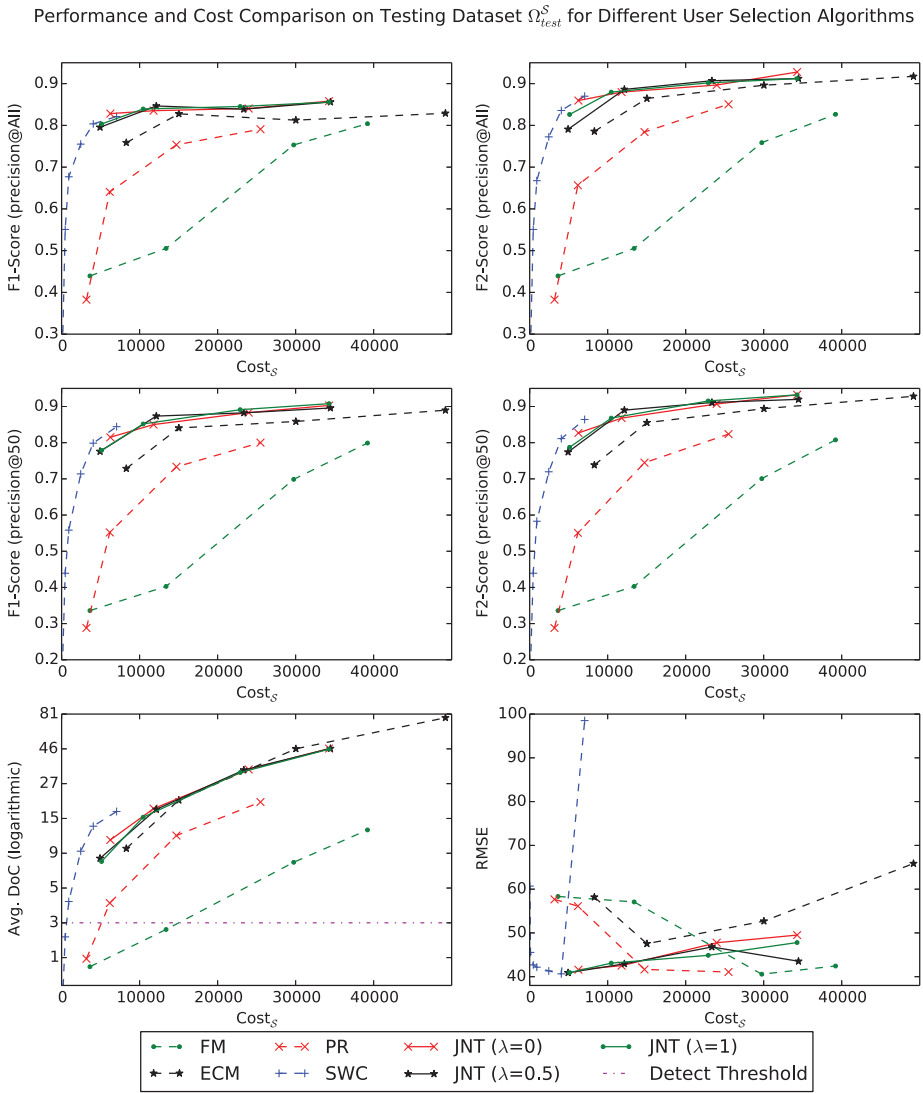


Fig. 2. Performance and cost comparison on testing dataset Ω_{test}^S with selected users. The x -axis in the figure shows the total cost of selected users. Under the same cost constraints, the proposed algorithm JNT shows better performance in F-scores, average degree of coverage, and lower RMSE than other algorithms.

by using the full dataset directly, in addition to the cost and time needed to capture such full dataset in real time, especially for third-party analyzers who need to crawl the dataset on their own. As a conclusion, the *cost-effectiveness* of an algorithm should be considered seriously in an online environment since there are always some kinds of cost constraints in practice.

Therefore, based on performance shown in Table VI and $Cost_S$ shown in Table V, we draw a figure showing performance vs. total cost of selected users ($Cost_S$) in Figure 2 to compare the cost-effectiveness of all the algorithms with user selection procedures. Algorithm TLC is excluded in this figure, as it uses all users' real-time microposts,

and thus its total cost is too high to compare. In Figure 2, the x -axis is the total cost of selected users (Cost_S) listed in Table V; the y -axes of the sub-figures are F1-score, F2-score, average DoC (logarithmic scaled), and RMSE, respectively. The dash-dotted purple horizon line in the bottom left sub-figure indicates the detection threshold \tilde{X} , which is set to default value 3 for all topics in our experiments. All the other solid lines (representing JNT) and dashed lines (representing the other algorithms, including SWC with user size K bigger than 800) in the figure represent results using the microposts of the subset users that are selected by different user selection algorithms.

In the first place, we compare the recall and precision rates of different algorithms shown in Table VI as well as the F-scores shown in Figure 2, for parameter sets I through IV. Similarly to the topic coverage performance in the training dataset, in real-time evaluations, algorithms FM and SWC have lower recall rates and lower average degree of coverage among all algorithms. The low average DoC suggests that their detection performance is too sensitive on cost constraints and online detection threshold. When cost constraint M decreases and their average DoC become lower than the detection threshold, their detection performance will be too low to be competitive. In contrast, average DoC of algorithm JNT ($\lambda \geq 0$) under various cost constraints are always beyond the detection threshold, so its recall is better all the time. The detection performance of algorithm PR is a little bit better than FM and SWC but worse than ECM and JNT . For algorithm SWC with larger user size constraints $K > 800$, its detection performance can be comparable with PR , but its selected users size is several times bigger than all the other algorithms, which contrast with the intention for selecting small-sized but representative subset users, and its performance is still worse than JNT . In terms of algorithm ECM , it can be found from Figure 2 that its F-scores are higher than PR , SWC , and FM , but its detection performance is still worse than JNT in most cases while its cost is even larger. As JNT ($\lambda \geq 0$) has the highest F-scores under the same cost constraints, it is the most cost-effective one in detection performance among all algorithms.

Next, prediction performance is compared using column “RMSE” in Table VI and the bottom right sub-figure of Figure 2. The RMSE of ECM becomes higher than the other algorithms when its selected users size increases, and the RMSE of algorithm SWC with larger user size constraints $K > 800$ is also very high. Due to the fact that algorithm ECM tends to select users with higher cost and the selected users’ size of SWC with larger K is much bigger, the most reasonable explanation for their RMSE increment is that the prediction accuracy is affected by the “noise” microposts in their selected users’ microposts, since the two algorithms do not consider the users’ prediction accuracy during subset user selection. In contrast, the RMSEs of FM and PR are higher when cost constraints are low, as in these cases they are short of valuable users’ microposts for prediction. In other words, the prediction performance of the above four methods are too sensitive on cost constraints and thus not cost-effective. In general, RMSE of algorithm JNT ($\lambda > 0$) is quite stable and relatively low among the four sets of parameters.

In light of the above, the proposed algorithm JNT has the overall best joint online detection and prediction performance over testing dataset within cost constraints.

6.4.2. Discussions on Early Detection and Prediction. In Table VI, column “T-G” shows the average detection time advantage of our system in hours, which is always positive using the proposed algorithm JNT ($\lambda \geq 0$). It means that our system can detect the trending topics much earlier than they appear in the official Trends Lists of Weibo and search engines. In our experiments, the observation time t_t needed for future popularity prediction is set to 6 hours after the trending topic is initiated and detected by our system. Removing the 6 hours from column “T-G” in Table VI, the result is still

Table VIII. Performance Comparison for Algorithm *JNT* with Different λ

Set	λ	Recall			Precision	F1 ($\bar{X}=3$)	F2 ($\bar{X}=3$)	A-D	T-G	RMSE
		$\bar{X}=3$	$\bar{X}=5$	$\bar{X}=8$	P@All	P@All	P@All			
III	0	85/93	83/93	77/93	161/207	0.8404	0.8831	33.8	26.3	47.75
	0.1	86/93	84/93	77/93	162/209	0.8433	0.8904	33.2	25.8	48.21
	0.5	85/93	84/93	78/93	162/209	0.8388	0.8824	33.5	27.9	46.79
	1	86/93	82/93	76/93	158/203	0.8452	0.8912	32.2	30.1	44.88
	2	86/93	82/93	76/93	154/198	0.8449	0.8911	32.5	27.1	42.96
	5	84/93	82/93	74/93	141/175	0.8517	0.8819	27.8	25.1	43.87

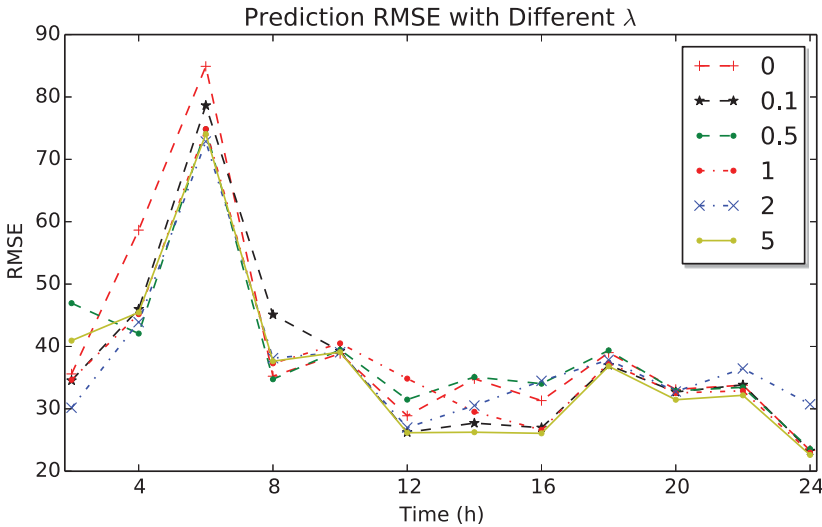


Fig. 3. Popularity prediction RMSE comparison for Algorithm *JNT* with different λ using parameter Set III. It shows the average RMSE per time slot within each hour for the next 24 hours after the prediction began.

decent, and thus we can accomplish the joint tasks of trending topic detection and prediction several hours in advance of the official lists. This reveals another advantage of our proposed framework: It is a third-party system that is very practical in both early trending topic detection and early prediction for real microblogging services, using a relatively small budget on cost.

6.4.3. *Discussions on λ* . During the user selection procedure for algorithm *JNT* with λ greater than 0, the reward of a user consists of both detection reward and prediction reward. During user selection, the system will consider more about a selected user’s contributions on prediction accuracy when the value of λ increases; and the system will focus more on user’s topic detection ability when λ drops.

To exhibit the effect of λ , we run additional experiments with various λ values. Taking experiments using parameter set III as an example, detection and prediction performance with different λ values are shown in Table VIII. The average RMSE per time slot within every hour for the next 24 hours are also shown in Figure 3.

From Table VIII, it can be seen that recall rate drops a little bit as λ increases from 0. In the meantime, corresponding prediction performance improves as expected, which can be observed in Figure 3 and column “RMSE” in Table VIII. Based on this trend, if

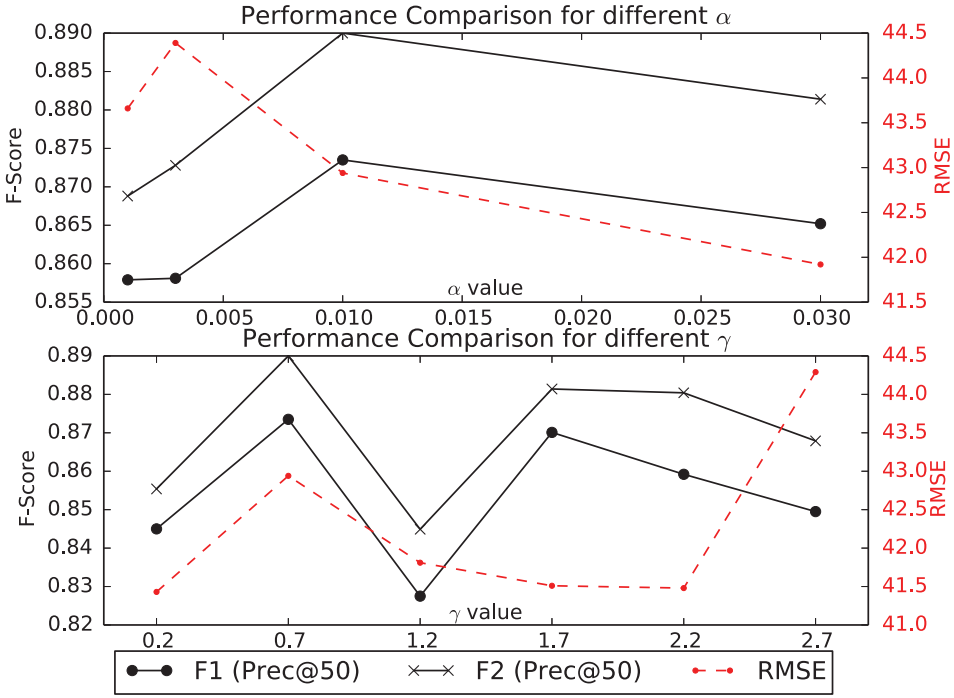


Fig. 4. Performance comparison for Algorithm *JNT* with different values of α and γ . The two coefficients α and γ are defined in Equation (12) and Equation (14). Parameter Set II and $\lambda = 0.5$ are used in the comparison.

λ is too high (e.g., >5), then the detection performance, average *DoC*, and “T-G” (time gained) will drop a lot, and thus the joint performance will not be ideal. Therefore, we should pay attention to the weight of prediction during user selection to maintain good detection and prediction accuracy, as well as timeliness to ensure the time gained is still enough to ensure early detection and prediction. In our datasets, it is desirable to set λ between 0.5 and 2.

6.4.4. Discussions on Other Coefficients in Algorithm *JNT*. Besides using λ to indicate the weights of prediction performance in user selection, algorithm *JNT* also uses the concept of *dynamic reward* and *dynamic cost boundary* to improve trending topic coverage, which is explained in Section 5.2. Here, we exhibit the impact of different α and γ values in Equation (12) and Equation (14) by experiments. The result comparison is shown in Figure 4, in which parameters set II and $\lambda = 0.5$ are used. The x-axes in the upper and lower sub-figures show the varying α values and γ values, respectively. The y-axes are their corresponding F1-score, F2-score, and RMSE with corresponding testing dataset Ω_{test}^S .

When α is smaller, the detection reward for covering topics with lower *DoC* will become a little bit larger, and thus these topics must be covered in user selection. In terms of γ , if its value is too small, then the cost boundary will become quite large and users with quite large costs will be selected first, and it might not be so cost-effective. According to the results shown in Figure 4, $\alpha = 0.01$ and $\gamma = 0.7$ are chosen as the default values in all experiments with algorithm *JNT*, as their F1- and F2-scores are the best and the RMSE are relatively small with these coefficient values.

7. CONCLUSIONS AND FUTURE WORK

In this article, we present a cost-effective online trending topic detection and prediction system for microblogging services from a third-party perspective. The proposed system can run under strict resource constraints while not compromising on the performance in detection and prediction. In order to satisfy resource budget, online trending topic multi-coverage requirements, as well as popularity prediction accuracy, we propose the notion of utilizing a subset of selected users to accomplish the task. We formulate the subset user selection problem as optimization tasks, and propose efficient algorithms to solve the problem.

To evaluate the online performance of joint detection and prediction system, we collect the experiment data from real microblogging service networks, and utilize them into offline dataset and real-time testing dataset that are used differently in our experiment settings. The performance comparison results prove that the proposed algorithm *JNT* outperforms the state-of-the-art algorithms in detection and prediction accuracy while being cost-effective. Experiments show that by tracking only 500 users of 0.6 million microblog users and processing at most 30,000 microposts daily, about 92% of the trending topics among all users could be detected and then predicted by the proposed system. Moreover, the trending topics and their future popularity can be detected and predicted by our system much earlier than when they are published by official trends lists in microblogging services. As the proposed system is cost-effective, it is very practically applicable to real-world usage.

In future works, we plan to extend the system, algorithm, and experiments on different categories of microposts, so users with different interests can be selected and utilized for topic analysis. Distributed computing technology can be applied to the user selection algorithm to speed up the training. More factors in the dataset can also be used in the algorithms, for example, the time factors that a user tends to participate in trending topics. In addition, a new mechanism, such as dynamically updating selected users according to overall performance or time factors, is another interesting area.

ACKNOWLEDGMENTS

The authors thank all colleagues, reviewers, and editors who contributed to this article. The authors also thank all the creators and maintainers of the tools we used in experiments.

REFERENCES

- Mohamed Ahmed, Stella Spagna, Felipe Huici, and Saverio Niccolini. 2013. A peek into the future: Predicting the evolution of popularity in user generated content. In *Proceedings of the 6th ACM International Conference on Web Search and Data Mining (WSDM'13)*. ACM, New York, NY, 607–616. DOI : <http://dx.doi.org/10.1145/2433396.2433473>
- James Allan (Ed.). 2002. *Topic Detection and Tracking: Event-based Information Organization*. Kluwer Academic Publishers, Norwell, MA.
- Foteini Alvanaki, Sebastian Michel, Krithi Ramamritham, and Gerhard Weikum. 2012. See what's en-Blogue: Real-time emergent topic identification in social media. In *Proceedings of the 15th International Conference on Extending Database Technology (EDBT'12)*. ACM, New York, NY, 336–347. DOI : <http://dx.doi.org/10.1145/2247596.2247636>
- Sitaram Asur, Bernardo A. Huberman, Gabor Szabo, and Chunyan Wang. 2011. Trends in social media: Persistence and decay. *SSRN Electr. J.* (Feb. 2011). DOI : <http://dx.doi.org/10.2139/ssrn.1755748>
- Roja Bandari, Sitaram Asur, and Bernardo A. Huberman. 2012. The pulse of news in social media: Forecasting popularity. In *Proceedings of the 6th International Conference on Weblogs and Social Media (ICWSM'12)*. The AAAI Press, 26–33.
- Livio Bertacco. 2006. *Exact and Heuristic Methods for Mixed Integer Linear Programs*. Ph.D. Dissertation. Ph. D. thesis, Università degli Studi di Padova.
- Bin Bi, Yuanyuan Tian, Yannis Sismanis, Andrey Balmin, and Junghoo Cho. 2014. Scalable topic-specific influence analysis on microblogs. In *Proceedings of the 7th ACM International Conference*

- on *Web Search and Data Mining (WSDM'14)*. ACM, New York, NY, 513–522. DOI: <http://dx.doi.org/10.1145/2556195.2556229>
- Petko Bogdanov, Michael Busch, Jeff Moehlis, Ambuj K. Singh, and Boleslaw K. Szymanski. 2013. The social media genome: Modeling individual topic-specific behavior in social media. In *Proceedings of the 2013 IEEE/ACM International Conference on Advances in Social Networks Analysis and Mining (ASONAM'13)*. ACM, New York, NY, 236–242. DOI: <http://dx.doi.org/10.1145/2492517.2492621>
- Sergey Brin and Lawrence Page. 2012. Reprint of: The anatomy of a large-scale hypertextual web search engine. *Comput. Netw.* 56, 18 (2012), 3825–3833. DOI: <http://dx.doi.org/10.1016/j.comnet.2012.10.007>
- Mario Cataldi, Luigi Di Caro, and Claudio Schifanella. 2010. Emerging topic detection on twitter based on temporal and social terms evaluation. In *Proceedings of the 10th International Workshop on Multimedia Data Mining (MDMKDD'10)*. ACM, New York, NY, Article 4, 10 pages. DOI: <http://dx.doi.org/10.1145/1814245.1814249>
- Kai Chen, Yi Zhou, Hongyuan Zha, Jianhua He, Pei Shen, and Xiaokang Yang. 2013b. Cost-effective node monitoring for online hot event detection in Sina Weibo microblogging. In *Proceedings of the 22nd International Conference on World Wide Web (WWW'13 Companion)*. ACM, New York, NY, 107–108. DOI: <http://dx.doi.org/10.1145/2487788.2487837>
- Le Chen, Chi Zhang, and Christo Wilson. 2013a. Tweeting under pressure: Analyzing trending topics and evolving word choice on Sina Weibo. In *Proceedings of the 1st ACM Conference on Online Social Networks (COSN'13)*. ACM, New York, NY, 89–100. DOI: <http://dx.doi.org/10.1145/2512938.2512940>
- Wei Chen, Yajun Wang, and Siyu Yang. 2009. Efficient influence maximization in social networks. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'09)*. ACM, New York, NY, 199–208. DOI: <http://dx.doi.org/10.1145/1557019.1557047>
- Jeffrey Dean and Sanjay Ghemawat. 2008. MapReduce: Simplified data processing on large clusters. *Commun. ACM* 51, 1 (Jan. 2008), 107–113. DOI: <http://dx.doi.org/10.1145/1327452.1327492>
- Pedro Domingos and Matt Richardson. 2001. Mining the network value of customers. In *Proceedings of the 7th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'01)*. ACM, New York, NY, 57–66. DOI: <http://dx.doi.org/10.1145/502512.502525>
- Nan Du, Le Song, Manuel Gomez-Rodriguez, and Hongyuan Zha. 2013. Scalable influence estimation in continuous-time diffusion networks. In *Advances in Neural Information Processing Systems 26*. Curran Associates, Inc., 3147–3155.
- Pablo A. Estevez, Pablo Vera, and Kazumi Saito. 2007. Selecting the most influential nodes in social networks. In *Proceedings of the 2007 International Joint Conference on Neural Networks*. IEEE, 2397–2402. DOI: <http://dx.doi.org/10.1109/IJCNN.2007.4371333>
- Schubert Foo and Hui Li. 2004. Chinese word segmentation and its effect on information retrieval. *Inform. Process. Manag.* 40, 1 (Jan. 2004), 161–190. DOI: [http://dx.doi.org/10.1016/S0306-4573\(02\)00079-1](http://dx.doi.org/10.1016/S0306-4573(02)00079-1)
- Gabriel Pui Cheong Fung, Jeffrey Xu Yu, Huan Liu, and Philip S. Yu. 2007. Time-dependent event hierarchy construction. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'07)*. ACM, New York, NY, 300–309. DOI: <http://dx.doi.org/10.1145/1281192.1281227>
- Manuel Gomez-Rodriguez, Jure Leskovec, and Andreas Krause. 2012. Inferring networks of diffusion and influence. *ACM Trans. Knowl. Discov. Data* 5, 4, Article 21 (Feb. 2012), 37 pages. DOI: <http://dx.doi.org/10.1145/2086737.2086741>
- Yi Han, Lei Deng, Binying Xu, Lumin Zhang, Bin Zhou, and Yan Jia. 2013. Predicting the social influence of upcoming contents in large social networks. In *Proceedings of the 5th International Conference on Internet Multimedia Computing and Service (ICIMCS'13)*. ACM, New York, NY, 17–22. DOI: <http://dx.doi.org/10.1145/2499788.2499834>
- Guangyan Huang, Jing He, Yanchun Zhang, Wanlei Zhou, Hai Liu, Peng Zhang, Zhiming Ding, Yue You, and Jian Cao. 2015. Mining streams of short text for analysis of world-wide event evolutions. *World Wide Web* 18, 5 (2015), 1201–1217. DOI: <http://dx.doi.org/10.1007/s11280-014-0293-1>
- David Kempe, Jon Kleinberg, and Éva Tardos. 2003. Maximizing the spread of influence through a social network. In *Proceedings of the 9th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'03)*. ACM, New York, NY, 137–146. DOI: <http://dx.doi.org/10.1145/956750.956769>
- Andrey Kupavskii, Alexey Umnov, Gleb Gusev, and Pavel Serdyukov. 2013. Predicting the audience size of a tweet. In *Proceedings of the 7th International Conference on Weblogs and Social Media (ICWSM'13)*. The AAAI Press, 693–696.
- Chung-Hong Lee, Hsin-Chang Yang, Tzan-Feng Chien, and Wei-Shiang Wen. 2011. A novel approach for event detection by mining spatio-temporal information on microblogs. In *Proceedings of the 2011 International Conference on Advances in Social Networks Analysis and Mining*. IEEE, 254–259. DOI: <http://dx.doi.org/10.1109/ASONAM.2011.74>

- Jure Leskovec, Lars Backstrom, and Jon Kleinberg. 2009. Meme-tracking and the dynamics of the news cycle. In *Proceedings of the 15th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'09)*. ACM, New York, NY, 497–506. DOI: <http://dx.doi.org/10.1145/1557019.1557077>
- Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne VanBriesen, and Natalie Glance. 2007. Cost-effective outbreak detection in networks. In *Proceedings of the 13th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'07)*. ACM, New York, NY, 420–429. DOI: <http://dx.doi.org/10.1145/1281192.1281239>
- Juha Makkonen, Helena Ahonen-Myka, and Marko Salmenkivi. 2004. Simple semantics in topic detection and tracking. *Inform. Retrieval*, 7, 3–4 (Sep. 2004), 347–368. DOI: <http://dx.doi.org/10.1023/B:INRT.0000011210.12953.86>
- Michael Mathioudakis and Nick Koudas. 2010. TwitterMonitor: Trend detection over the twitter stream. In *Proceedings of the 2010 ACM SIGMOD International Conference on Management of Data (SIGMOD'10)*. ACM, New York, NY, 1155–1158. DOI: <http://dx.doi.org/10.1145/1807167.1807306>
- Zhongchen Miao, Kai Chen, Yi Zhou, Hongyuan Zha, Jianhua He, Xiaokang Yang, and Wenjun Zhang. 2015. Online trendy topics detection in microblogs with selective user monitoring under cost constraints. In *Proceedings of the 2015 IEEE International Conference on Communications (ICC'15)*. 1194–1200. DOI: <http://dx.doi.org/10.1109/ICC.2015.7248485>
- Fred Morstatter, Jürgen Pfeffer, and Huan Liu. 2014. When is it biased? Assessing the representativeness of twitter's streaming API. In *Proceedings of the 23rd International Conference on World Wide Web (WWW'14 Companion)*. ACM, New York, NY, 555–556. DOI: <http://dx.doi.org/10.1145/2567948.2576952>
- Fred Morstatter, Jürgen Pfeffer, Huan Liu, and Kathleen M. Carley. 2013. Is the sample good enough? Comparing data from twitter's streaming API with twitter's firehose. In *Proceedings of the 7th International Conference on Weblogs and Social Media (ICWSM'13)*. 400–408.
- Seth A. Myers and Jure Leskovec. 2014. The bursty dynamics of the twitter information network. In *Proceedings of the 23rd International Conference on World Wide Web (WWW'14)*. ACM, New York, NY, 913–924. DOI: <http://dx.doi.org/10.1145/2566486.2568043>
- Mor Naaman, Hila Becker, and Luis Gravano. 2011. Hip and trendy: Characterizing emerging trends on twitter. *J. Am. Soc. Inform. Sci. Technol.* 62, 5 (May 2011), 902–918. DOI: <http://dx.doi.org/10.1002/asi.21489>
- Ramasuri Narayanam and Yadati Narahari. 2011. A shapley value-based approach to discover influential nodes in social networks. *IEEE Trans. Autom. Sci. Eng.* 8, 1 (Jan. 2011), 130–147. DOI: <http://dx.doi.org/10.1109/TASE.2010.2052042>
- Aditya Pal and Scott Counts. 2011. Identifying topical authorities in microblogs. In *Proceedings of the 4th ACM International Conference on Web Search and Data Mining (WSDM'11)*. ACM, New York, NY, 45–54. DOI: <http://dx.doi.org/10.1145/1935826.1935843>
- R. Papka and J. Allan. 1998. *On-Line New Event Detection Using Single Pass Clustering*. Technical Report. University of Massachusetts, Amherst, MA, USA.
- Georgios Petkos, Symeon Papadopoulos, and Yiannis Kompatsiaris. 2014. Two-level message clustering for topic detection in twitter. In *Proceedings of the SNOW 2014 Data Challenge*. 49–56.
- Polina Rozenshtein, Aris Anagnostopoulos, Aristides Gionis, and Nikolaj Tatti. 2014. Event detection in activity networks. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'14)*. ACM, New York, NY, 1176–1185. DOI: <http://dx.doi.org/10.1145/2623330.2623674>
- Takeshi Sakaki, Makoto Okazaki, and Yutaka Matsuo. 2010. Earthquake shakes twitter users: Real-time event detection by social sensors. In *Proceedings of the 19th International Conference on World Wide Web (WWW'10)*. ACM, New York, NY, 851–860. DOI: <http://dx.doi.org/10.1145/1772690.1772777>
- Erich Schubert, Michael Weiler, and Hans-Peter Kriegel. 2014. SigniTrend: Scalable detection of emerging topics in textual streams by hashed significance thresholds. In *Proceedings of the 20th ACM SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'14)*. ACM, New York, NY, 871–880. DOI: <http://dx.doi.org/10.1145/2623330.2623740>
- Aleksandr Simma and Michael I. Jordan. 2010. Modeling events with cascades of poisson processes. In *Proceedings of the 26th Conference on Uncertainty in Artificial Intelligence (UAI'10)*. AUAI Press, 546–555. Retrieved from https://dslpitt.org/uai/displayArticleDetails.jsp?mmnu=1&smnu=2&article_id=2139&proceeding_id=26
- Oren Tsur and Ari Rappoport. 2012. What's in a hashtag? Content based prediction of the spread of ideas in microblogging communities. In *Proceedings of the 5th ACM International Conference on Web Search and Data Mining (WSDM'12)*. ACM, New York, NY, 643–652. DOI: <http://dx.doi.org/10.1145/2124295.2124320>
- Yu Wang, Gao Cong, Guojie Song, and Kunqing Xie. 2010. Community-based greedy algorithm for mining top-k influential nodes in mobile social networks. In *Proceedings of the 16th ACM SIGKDD*

- International Conference on Knowledge Discovery and Data Mining (KDD'10)*. ACM, New York, NY, 1039–1048. DOI : <http://dx.doi.org/10.1145/1835804.1835935>
- Jianshu Weng, Ee-Peng Lim, Jing Jiang, and Qi He. 2010. TwitterRank: Finding topic-sensitive influential twitterers. In *Proceedings of the 3rd ACM International Conference on Web Search and Data Mining (WSDM'10)*. ACM, New York, NY, 261–270. DOI : <http://dx.doi.org/10.1145/1718487.1718520>
- Jaewon Yang and Jure Leskovec. 2011. Patterns of temporal variation in online media. In *Proceedings of the 4th ACM International Conference on Web Search and Data Mining (WSDM'11)*. ACM, New York, NY, 177–186. DOI : <http://dx.doi.org/10.1145/1935826.1935863>
- Mengmeng Yang, Kai Chen, Zhongchen Miao, and Xiaokang Yang. 2014. Cost-effective user monitoring for popularity prediction of online user-generated content. In *Proceedings of the 2014 IEEE International Conference on Data Mining Workshop*. IEEE, 944–951. DOI : <http://dx.doi.org/10.1109/ICDMW.2014.72>
- Yiming Yang, Tom Pierce, and Jaime Carbonell. 1998. 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 (SIGIR'98)*. ACM, New York, NY, 28–36. DOI : <http://dx.doi.org/10.1145/290941.290953>

Received January 2016; revised September 2016; accepted September 2016