

# Expert Finding for Microblog Misinformation Identification

*Chen Liang\** *Zhiyuan Liu\** *Maosong Sun*

Department of Computer Science and Technology  
State Key Lab on Intelligent Technology and Systems  
National Lab for Information Science and Technology  
Tsinghua University, Beijing 100084, China

{chenliang.harry, lzy.thu}@gmail.com, sms@tsinghua.edu.cn

## ABSTRACT

The growth of social media provides a convenient communication scheme for people, but at the same time it becomes a hotbed of misinformation. The wide spread of misinformation over social media is injurious to public interest. We design a framework, which integrates collective intelligence and machine intelligence, to help identify misinformation. The basic idea is: (1) automatically index the expertise of users according to their microblog contents; and (2) match the experts with given suspected misinformation. By sending the suspected misinformation to appropriate experts, we can collect the assessments of experts to judge the credibility of information, and help refute misinformation. In this paper, we focus on expert finding for misinformation identification. We propose a tag-based method to index the expertise of microblog users with social tags. Experiments on a real world dataset demonstrate the effectiveness of our method for expert finding with respect to misinformation identification in microblogs.

## TITLE AND ABSTRACT IN CHINESE

### 面向微博不实信息识别的专家发现

在为人们提供了便利的交流方式的同时，社会媒体也成为不实信息传播的温床。不实信息在社会媒体中的广泛传播将对公共利益造成损害。这里，我们提出综合利用群体智能和机器智能帮助识别不实信息，基本思想是：（1）根据微博用户产生内容自动分析和索引用户的专长；（2）自动将可疑的不实信息与相应的专家匹配。通过将不实信息发送给合适的专家，我们可以收集专家们对信息可信性的评估，帮助识别不实信息和辟谣。本论文将主要探讨面向不实信息识别的专家发现问题。我们采用基于标签的方法来索引微博用户的专长。在真实数据上的实验表明，我们的方法可以有效进行微博专家发现用于不实信息识别。

---

**KEYWORDS:** misinformation identification, expert finding, microblog.

**KEYWORDS IN CHINESE:** 不实信息识别, 专家发现, 微博.

---

---

\* indicates equal contributions from these authors.

# 1 Introduction

Although rumors lack a specific definition, most theories agree that a rumor is a statement of information, whose veracity is not quickly or ever confirmed, spreading from person to person and pertaining to an object, event or issue in public concern (Peterson and Gist, 1951). Rumors are regarded as a type of misinformation. In recent years, online social media is growing rapidly. Social media provides a convenient communication scheme between people. Meanwhile, the scheme enables unreliable sources to spread large amounts of unverified information among people. Rumors are thus possible to spread more quickly and widely through online social media compared to traditional offline social communities. The wide spread of misinformation may bring disorder to people especially when they are facing crises. This indicates that it is crucial for social media to identify misinformation in time so as to limit the spread of rumors.

Most existing research efforts on rumors in social media focus on their external features, such as spread and conversation patterns. It is a consensus that automatically identifying rumors via in-depth content analysis is a challenging task. Social media still lacks solutions to effectively identify and refute misinformation to stop it from wide spread. Although fully-automatic identification of misinformation is currently a mission impossible for computer programs, most rumors can be easily identified by human experts with corresponding knowledge or experiences. Due to the popularity of social media, most experts can be found in social media. Under such a scenario, we design a framework to identify misinformation with the help of experts in microblogs. We automatically route suspected misinformation to a set of experts who can assess the credibility. With the assessments from experts, we can help determine the credibility of the information, identify and refute misinformation, and eventually stop the wide spread of rumors.

The most crucial part of the framework is finding appropriate experts for suspected misinformation. A relevant task has been studied as expert finding. However finding experts for suspected misinformation is different from the traditional task in many aspects, which make it more challenging. In this paper, we focus on expert finding for misinformation identification. With the help of experts, we incorporate the power of natural language processing techniques and the knowledge of human experts. With the method, we may help social media achieve self-management and self-organization.

## 2 Empirical Analysis of Rumors

We give empirical analysis of rumors on Sina Weibo, the largest microblog service in China. Sina Weibo has more than 250 million registered users as of October 2011, over 60% network users have accounts of Sina Weibo, and 31% Weibo users post more than 3 messages per day. Due to the overflow of rumors, Sina Weibo maintains a team to refute rumors. Since the process is carried out manually, it is manpower intensive while the refutation is usually delayed and the scope is limited.

We collect 859 rumors identified by both Sina Weibo team and a public interest organization (guokr.com) for empirical analysis. All the rumors have widely spread over Sina Weibo ranging from 18 November, 2010 to 29 December, 2011. By studying the real-world rumors, according to what restricts people from identifying them as rumors, we manually categorize rumors into two classes: (1) 588 out of 859 rumors are **domain knowledge constrained (DKC) misinformation**. This type of information usually talks about some domain-specific

topics. Most people do not master the corresponding professional knowledge and thus cannot verify the correctness. For example, a rumor claimed that “A nutritionist finds that if you eat a bag of instant noodles, the liver will need 32 days for detoxification”, which is related to the knowledge of food hygiene and nutrition. (2) 271 out of 859 rumors are **time-space constrained (TSC) misinformation**. This type of misinformation is usually related to some events that occur in some places and time. Most people have not experienced these events and thus cannot check the authenticity. For example, a rumor claimed that “Some mentally retarded children in Xiangyang, Hubei were cut out tongues and genitals”. The people not living in that city will not be able to verify the reliability of the information. According to the analysis, we summarize that the two types of rumors are different in: (1) Their topics are quite different. DKC rumors focus on the topics of science and technologies, while TSC rumors focus on the topics of public security, society and politics. (2) TSC rumors usually talk about events and thus often mention specific names of persons, places or organizations; while DKC rumors seldom mention specific names.

We build a two-class classifier using the occurrences of names and topics as features to quantitatively identify the differences between the two types of misinformation. The classifier is built using LIBLINEAR (Fan et al., 2008). For each message, features are boolean values indicating the appearance of names and topics. We perform 5-cross validation for evaluation and the prediction accuracies are 0.848, 0.804, and 0.865 when using names, topics or both as features. We can see that most misinformation can be well distinguished according to these features. Based on the characteristics of DKC and TSC misinformation, we propose a unified method to find applicable experts for them.

### 3 Tag-based Method for Expert Finding

Suppose all microblog users are candidate experts denoted as a set  $E$ . These users may be either people or organizations. Given a suspected misinformation  $m$ , the probability of selecting a microblog user  $e$  from  $E$  being an expert on  $m$  can be estimated as

$$\Pr(e|m) = \frac{\Pr(m|e)\Pr(e)}{\Pr(m)} \propto \Pr(m|e)\Pr(e), \quad (1)$$

where  $\Pr(e)$  is the prior probability of an expert;  $\Pr(m)$  is the prior probability of  $m$ , which remains the same for all candidate experts and thus is ignored for expert ranking. Here,  $\Pr(e)$  is estimated as the authority of  $e$ , and  $\Pr(m|e)$  is as the expertise of the expert  $e$  on  $m$ . We adopt social tags annotated by microblog users to model expertise.

For DKC misinformation, we use Eq.(1) to find experts from  $E$  directly; while for TSC misinformation, we have to restrict the set of candidate experts with respect to the named entities that have appeared in  $m$ . We introduce our method in details in the following three aspects: (1) modeling expertise with tag-based method to compute  $\Pr(e|m)$ ; (2) computing authorities of experts, i.e.,  $\Pr(e)$ ; and (3) restricting candidate set of experts for TSC.

#### 3.1 Modeling Expertise with Tag-based Method

Social tagging is an iconic application in social media (Gupta et al., 2010). Sina Weibo allows users to annotate tags for themselves, which may attract users with similar interests to follow them. Expertise of experts can be represented in terms of social tags because tags represent the interests or characteristics of users to some extent. For example, a user whose occupation is ophthalmologist may annotate itself with the tag “ophthalmology”.

We denote the set of tags as  $T$ . Suppose  $\Pr(m|e)$  is a generation process as follows. An expert  $e$  first generates a tag  $t \in T$ , and then  $t$  generates the message  $m$ . Given the generation process, the probability  $\Pr(m|e)$  can thus be estimated as follows:

$$\Pr(m|e) = \sum_t \Pr(m|t)\Pr(t|e) = \sum_t \frac{\Pr(t|m)\Pr(m)}{\Pr(t)} \Pr(t|e) \propto \sum_t \frac{\Pr(t|m)}{\Pr(t)} \Pr(t|e), \quad (2)$$

where  $\Pr(t)$  indicates the prior probability for the tag  $t$ ,  $\Pr(t|m)$  measures the probability of  $t$  given the message  $m$ , and  $\Pr(t|e)$  computes the probability of expertise  $t$  given the expert  $e$ . The prior  $\Pr(t)$  can be estimated using the number of microblog users who annotate themselves with the tag  $t$ ; while  $\Pr(t|m)$  and  $\Pr(t|e)$  are modeled as a problem of social tag suggestion task.  $\Pr(t|e)$  can be decomposed into two parts: one is the suggestion score of  $t$  given the messages posted by  $e$ , and the other is whether  $e$  has annotated itself with  $t$ . The two parts are combined with a smoothing factor  $\gamma$  ranging from 0 to 1,  $\Pr(t|e) = \gamma \sum_{m \in M_e} \Pr(t|m)\Pr(m|e) + (1-\gamma)\mathbf{1}_{t \in T_e}$ , where  $M_e$  is the set of messages that are posted by  $e$ ,  $\Pr(t|m)$  is the ranking score of  $t$  given the message  $m$ ,  $\Pr(m|e)$  indicates the weight of message  $m$  within all messages posted by  $e$ , and  $\mathbf{1}_{t \in T_e}$  equals 1 if the tag set  $T_e$  annotated by  $e$  contains  $t$  and 0 otherwise. In this paper, we simply set the weights of all messages  $\Pr(m|e)$  for  $e$  to be equal, and set  $\gamma = 0.5$ .

Based on the above analysis to  $\Pr(t|m)$  and  $\Pr(t|e)$ , the essential task is to suggest social tags for a message  $m$ . As the rapid growth of social media, social tag suggestion has been well studied (Gupta et al., 2010). There are two approaches for social tag suggestion: graph-based approach and content-based approach. Since we have to suggest tags according to the content of  $m$ , we follow the content-based approach. The specialty of our problem compared to previous problems lies in: (1) the method should be robust to noise and informal format of microblog messages; and (2)  $m$  is short with no more than 140 Chinese characters in Sina Weibo.

Taking the specialty in consideration, we propose to use word alignment model (WAM) in statistical machine translation (Brown et al., 1993) for social tag suggestion, which has been verified to outperform other existing content-based methods (Liu et al., 2011, 2012). Here we give a brief introduction to WAM, and introduce some important extensions to make the method appropriate to suggest social tags for microblog messages.

**WAM for Social Tag Suggestion.** Given a message  $m$ , WAM ranks candidate tags by computing their likelihood  $\Pr_{\text{WAM}}(t|m) = \sum_{w \in m} \Pr(t|w)\Pr(w|m)$ , where  $\Pr(w|m)$  is the weight of the word  $w$  in  $m$ , and  $\Pr(t|w)$  is the translation probability from  $w$  to  $t$  obtained from the translation models.  $\Pr(w|m)$  is estimated using term-frequency and inverse message frequency (TFIMF), which is similar to TFIDF. According to the ranking scores, we suggest the top- $M$  as tags for  $m$ . WAM can avoid the problem caused by noise and informal format of microblogs. Moreover, WAM can suggest tags that have not appeared in the given message. However, a tag that appears in the given message may be more important. Therefore, we improve WAM by combining WAM with frequency-based methods. A simple and effective frequency-based method is using TFIMF to rank candidate tags in a given message. We thus compute the ranking score using improved WAM (IWAM) for a candidate tag as follows,  $\Pr_{\text{IWAM}}(t|m) = \alpha \Pr_{\text{WAM}}(t|m) + (1-\alpha) \Pr_{\text{TFIMF}}(t|m)$ , where  $\alpha$  is a smoothing factor with range  $\alpha \in [0.0, 1.0]$ . In experiments we set  $\alpha = 0.5$  which achieves the best performance.

**Training Translation Models for WAM.** Training WAM for tag suggestion consists of two steps: preparing translation pairs and training translation models. The training set for traditional WAM consists of a number of translation pairs written in two languages. In our task, we have to collect sufficient translation pairs of microblog messages and their tags to capture the semantic relationship between them. However, microblogs are usually not annotated with tags. We thus propose to prepare translation pairs by automatically extracting tags using a simple and effective method for each message. The basic idea is that most results of the simple method are correct, while the errors can be filtered out by WAM. The preparation process is as follows. We first collect all tags annotated by microblog users in Sina Weibo. For each tag  $t$ , we record the users that annotate tag  $t$  as  $E_t$ . We group all tags with  $|E_t| > 10$  as a tag list. After that, we collect a large set of microblog messages. For each message  $m$ , we extract several tags according to the score of tag-frequency and inverse expert-frequency  $\text{TFIEF}_{(t,m)} = \text{TF}_{(t,m)}|E|/|E_t|$ . Similar to TFIDE,  $\text{TF}_{(t,m)}$  indicates the significance of the tag  $t$  in  $m$ , and  $|E|/|E_t|$  indicates the discriminative ability of the tag  $t$ . Using messages and their corresponding extracted tags, we build the translation pairs for WAM training.

We use IBM Model 1 (Brown et al., 1993) for WAM training. IBM Model 1 is a widely used word alignment algorithm which does not require linguistic knowledge for two languages. We have also tested more sophisticated word alignment algorithms such as IBM Model 3 for tag suggestion. However, these methods do not achieve better performance than IBM Model 1. Therefore, in this paper we only demonstrate the experimental results using IBM Model 1. In experiments, we select GIZA++ (Och and Ney, 2003) to train IBM Model 1.

### 3.2 Measuring Authority

Some works have been devoted to authority analysis of social media (Pal and Counts, 2011). The basic conclusion is that a microblog user has more authority if it has more followers and posts more original messages. Therefore, in this paper, we simply compute authority of a user  $e$  as:

$$\text{Pr}(e) = \frac{\log\left(\frac{|F_e|}{|A_e|}\right) \times \log(|M_e|)}{\sum_{e \in E_m} \log\left(\frac{|F_e|}{|A_e|}\right) \times \log(|M_e|)}, \quad (3)$$

where  $F_e$  is the follower set of  $e$  and  $A_e$  is the user set followed by  $e$ . The score is normalized over all experts in  $E_m$ .

### 3.3 Restricting Candidate Expert Set for TSC

We denote the names of each user  $e \in E$  as  $N_e$  and the names in  $m$  as  $N_m$ . We perform named entity disambiguation for  $N_m$  according to microblog users, and link each name  $n$  in  $N_m$  to all relevant microblog users that  $n$  really mentions. We denote the restricted candidate experts as a set  $E_m$ . To restrict candidate expert set for TSC, we perform the following three steps.

**Extracting Names for Microblog Experts.** We extract and index the names  $N_e$  of each microblog user  $e \in E$  according to its nickname, introduction and authentication reason. This problem is addressed as a sequence labeling task solved by conditional random fields (CRF) (Lafferty et al., 2001)<sup>1</sup>. Since nicknames, introductions and authentication reasons

<sup>1</sup>We use CRF++ for implementation, which can be obtained in <http://crfpp.sourceforge.net/>.

have obvious patterns, we can obtain the labeling accuracy of above 90% by training on a set of 500 manually annotated users.

**Named Entity Recognition for  $m$ .** Since Sina Weibo is in Chinese, we first perform Chinese word segmentation (CWS) and part-of-speech (POS) tagging for messages using the algorithm originally proposed in (Jiang et al., 2008). After that, we perform named entity recognition (NER). Since misinformation always pretends to be credible by written in a formal style, we can thus achieve high accuracy using the algorithm CRF (Nadeau and Sekine, 2007) based on the CWS and POS tagging output.

**Named Entity Disambiguation.** We disambiguate the names in  $m$  with respect to microblog users, and thus restrict candidate expert set from  $E$  to  $E_m = \{e | N_e \cap N_m \neq \emptyset\}$ . First, we find all microblog users by substring match between the names of microblog users and the names extracted from the message, denoted as  $E_s$ . These users are not all relevant to the names in  $m$ . The following task is to disambiguate the names in  $m$  to microblog users in  $E_s$  according to the relevance of these users with  $m$ . We follow the state-of-the-art algorithm in (Zheng et al., 2010), and use list-wise learning to rank (L2R) framework to address the problem. After investigating various combinations of features, we use the following effective features for L2R: (1) Follow-attention ratio which indicates the popularity of  $e$ . (2) The number of original messages that  $e$  has posted which indicates the vitality of  $e$ . (3) The numbers of comments and reposts for recent 100 messages which also indicates the recent vitality of  $e$ . (4) The number of microblog user names that appear in both recent 100 messages of  $e$  and  $m$ . This measures the semantic relatedness between  $e$  and  $m$ .

Since the number of TSC rumors are limited for training and testing, we instead manually annotate 6394 names in news articles ranging from June to December, 2011 as dataset, with each name linked to a microblog user. By training on 3,985 instances and testing on 2,409 instances, we obtain accuracy of 96.3%, which indicates the effectiveness of L2R for named entity disambiguation to microblog users. With the trained model on the dataset, we perform named entity disambiguation to names in given message  $m$  and restrict the candidate expert set to  $E_m$ .

## 4 Experiments and Analysis

We perform experiments on 859 rumors manually collected from Sina Weibo. We also collect 5 million the most active microblog users with their profiles and messages to build expert database. For each rumor, we recommend 10 experts from microblog users. We ask two editors to manually annotate the correctness of the results, who discussed and finally achieved final agreement on annotation. For the inconsistent annotations, the two editors discuss to achieve agreement. We use P@N for evaluation where  $N$  ranges from 1 to 10.

**Evaluation Results on DKC Rumors.** To investigate the effectiveness of tag-based method, we compare our method with language model, the state-of-the-art method for expert finding, on 588 DKC rumors. For each DKC rumor, we suggest maximum 10 microblog experts.

In language model, a candidate expert  $e$  is represented by a multinomial probability distribution over the vocabulary of words, i.e.,  $\Pr(w|\theta_e)$ . A message  $m$  is represented by a bag of words with each word generated independently. Therefore, the probability of  $m$  being generated by the language model  $\theta_e$  can be obtained by taking the product across all words in  $m$ :  $\Pr(m|e) = \prod_{w \in m} \Pr(w|\theta_e)^{n(w,m)}$ , where  $\Pr(w|\theta_e)$  is the probability of a word  $w$  given  $\theta_e$ , and  $n(w,m)$  is the number of times word  $w$  appears in  $m$ . The language model

of  $e$ ,  $\Pr(w|\theta_e)$ , is estimated as  $\Pr(w|\theta_e) = \sum_{m \in M_e} \Pr(w|m) \Pr(m|e)$ , where  $\Pr(m|e)$  is the weight of a message posted by  $e$ , and  $\Pr(w|m)$  is the generation probability by the message  $m$ . We set  $\Pr(m|e)$  equal for all messages posted by  $e$ ; while  $\Pr(w|m)$  is estimated using the TFIMF score of  $w$ . We also apply the Jelinek-Mercer method to smooth language model (Zhou et al., 2009), which is not introduced in detail for space limit.

We show the evaluation results in Fig. 1. From Fig. 1 we observe that: (1) The tag-based method consistently and significantly outperforms language model for expert finding. This indicates the effectiveness of the tag-based method. The reason is that tags are annotated by microblog users and provide sufficient information. (2) Although the performance of expert finding is far from perfection, it can help find experts and reduce manual work to a great extent. Moreover, the performance of expert finding can be further improved as more knowledge are taken into consideration, which will be our future work.

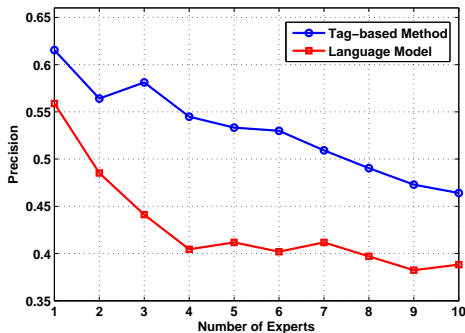


Figure 1: Evaluation results on expert finding for DKC rumors.

We also analyze the sample errors in expert finding for DKC rumors. The main reasons for these errors are: (1) Some tags are ambiguous and we may find experts that are irrelevant to  $m$ . For example, the tag “apple” may refer to either a type of fruits or an IT company. Although the tag-based method takes the topics of  $m$  into consideration, it still cannot thoroughly solve the problem. In future work, we may adopt tag disambiguation (Yeung et al., 2007) in our method. (2) We treat all tags annotated by experts equally. However, microblog users may annotate tags arbitrarily, which will thus import noise to our method. In future work, we will estimate confidence scores to the tags annotated by microblog users.

**Evaluation Results on TSC Rumors.** Different from DKC rumors, for TSC rumors our method will identify named entities in the suspected message to restrict the candidate expert set. In experiments, we set a person name may correspond to only one microblog user, while a place/organization name may refer to multiple microblog users. We evaluate the disambiguation for person names to demonstrate the performance of named entity disambiguation. The accuracy achieves 0.818 for all person names occurred in 271 TSC rumors. The precisions of expert finding for TSC rumors are 0.760 and 0.592 when suggesting 1 and 10 experts. The performance is slightly better than DKC rumors due to the impact of restricting candidate expert set, and is also much better than language model.

Take the rumor “Some mentally retarded children in Xiangyang, Hubei were cut out tongues and genitals” for example. We extract the named entities “Xiangyang, Hubei”. By substring

matching, we find microblog users such as “Unicom Xiangyang” (a mobile company), “PSB Xiangyang”, and “News Broadcast Xiangyang”. According to the relatedness between the given message and microblog experts, we rank “PSB Xiangyang” and “News Broadcast Xiangyang” higher than other microblog experts, which are more probable to refute the rumor.

**Discussion.** From the above evaluation and analysis, we validate the effectiveness of our tag-based method for expert finding from microblog users. This will greatly improve the efficiency of refuting misinformation and further prevent rumors from wide spread.

## 5 Related Work

Rumors have been extensively studied in sociology (Pendleton, 1998). However, quantitative studies of rumors have just begun, and microblog services provide a chance. Recently, researchers have developed different approaches to study rumors or misinformation. Some researchers devoted to finding information diffusion patterns over social networks (Kempe et al., 2003; Gruhl et al., 2004; Leskovec et al., 2009; Romero et al., 2011) and limiting the spread of misinformation by means of network structure (Budak et al., 2011). The spread patterns of rumors with respect to the content and conversations were also studied (Ennals et al., 2010; Mendoza et al., 2010; Qazvinian et al., 2011; Castillo et al., 2011). On one hand, most of these methods all focused on *external* features of rumors, which cannot ultimately determine whether a message is misinformation. On the other hand, the features can be obtained only after the information has spread over social networks.

Existing methods find experts based on either people relations (graph-based approach) or people meta-data (content-based approach). In the graph-based approach, users are ranked according to their authority scores computed by the algorithms such as HITS and PageRank (Zhang et al., 2007; Jurczyk and Agichtein, 2007). In the content-based approach, topic models (Mimno and McCallum, 2007) and language models (Balog et al., 2006; Petkova and Croft, 2006; Zhou et al., 2009; Li et al., 2011) have been explored. Due to sound foundations in statistical theory and sufficient performance, language models have been dominating techniques for expert finding (Balog, 2012). Different from existing methods, this paper proposes a tag-based method to find experts for suspected misinformation.

## Conclusion and Future Work

This paper proposes a novel framework for microblog misinformation identification with the favor of experts. We focus on the task of finding experts for suspected misinformation. By categorizing rumors into two types, i.e. domain-knowledge constrained and time-space constrained, we propose a unified tag-based method to find experts from microblog users and match suspected misinformation to appropriate experts. Experiments on the real-world dataset indicate the effectiveness of our method.

This is an initial step to fight against microblog misinformation. We plan the following future work. (1) Build a real-world system to fight against rumors and evaluate the effectiveness of our method. (2) Extend the work by considering more factors, such as the spread patterns (Budak et al., 2011) and conversation patterns (Ennals et al., 2010) of rumors. (3) Improve our method by considering social networks and tag disambiguation.

## Acknowledgments

This work is supported by the National Natural Science Foundation of China (NSFC) under the grant No. 61170196 and 61202140.



## References

- Balog, K. (2012). Expertise retrieval. *Foundations and Trends in Information Retrieval*, 6(2-3):127–256.
- Balog, K., Azzopardi, L., and De Rijke, M. (2006). Formal models for expert finding in enterprise corpora. In *Proceedings of SIGIR*, pages 43–50.
- Brown, P., Pietra, V., Pietra, S., and Mercer, R. (1993). The mathematics of statistical machine translation: Parameter estimation. *Computational linguistics*, 19(2):263–311.
- Budak, C., Agrawal, D., and El Abbadi, A. (2011). Limiting the spread of misinformation in social networks. In *Proceedings of WWW*, pages 665–674.
- Castillo, C., Mendoza, M., and Poblete, B. (2011). Information credibility on twitter. In *Proceedings of WWW*, pages 675–684.
- Ennals, R., Byler, D., Agosta, J., and Rosario, B. (2010). What is disputed on the web? In *Proceedings of the 4th workshop on Information credibility*, pages 67–74.
- Fan, R., Chang, K., Hsieh, C., Wang, X., and Lin, C. (2008). Liblinear: A library for large linear classification. *The Journal of Machine Learning Research*, 9:1871–1874.
- Gruhl, D., Guha, R., Liben-Nowell, D., and Tomkins, A. (2004). Information diffusion through blogspace. In *Proceedings of WWW*, pages 491–501.
- Gupta, M., Li, R., Yin, Z., and Han, J. (2010). Survey on social tagging techniques. *ACM SIGKDD Explorations Newsletter*, 12(1):58–72.
- Jiang, W., Mi, H., and Liu, Q. (2008). Word lattice reranking for chinese word segmentation and part-of-speech tagging. In *Proceedings of COLING*, pages 385–392.
- Jurczyk, P. and Agichtein, E. (2007). Discovering authorities in question answer communities by using link analysis. In *Proceedings of CIKM*, pages 919–922.
- Kempe, D., Kleinberg, J., and Tardos, É. (2003). Maximizing the spread of influence through a social network. In *Proceedings of SIGKDD*, pages 137–146.
- Lafferty, J., McCallum, A., and Pereira, F. (2001). Conditional random fields: Probabilistic models for segmenting and labeling sequence data. In *Proceedings of ICML*, pages 282–289.
- Leskovec, J., Backstrom, L., and Kleinberg, J. (2009). Meme-tracking and the dynamics of the news cycle. In *Proceedings of SIGKDD*, pages 497–506. ACM.
- Li, B., King, I., and Lyu, M. (2011). Question routing in community question answering: putting category in its place. In *Proceedings of CIKM*, pages 2041–2044.
- Liu, Z., Chen, X., and Sun, M. (2011). A simple word trigger method for social tag suggestion. In *Proceedings of EMNLP*, pages 1577–1588.
- Liu, Z., Chen, X., and Sun, M. (2012). Mining the interests of chinese microbloggers via keyword extraction. *Frontiers of Computer Science*, 6(1):76–87.

- Mendoza, M., Poblete, B., and Castillo, C. (2010). Twitter under crisis: can we trust what we rt? In *Proceedings of the 1st workshop on social media analytics*, pages 71–79.
- Mimno, D. and McCallum, A. (2007). Expertise modeling for matching papers with reviewers. In *Proceedings of SIGKDD*, pages 500–509.
- Nadeau, D. and Sekine, S. (2007). A survey of named entity recognition and classification. *Linguisticae Investigationes*, 30(1):3–26.
- Och, F. and Ney, H. (2003). A systematic comparison of various statistical alignment models. *Computational linguistics*, 29(1):19–51.
- Pal, A. and Counts, S. (2011). Identifying topical authorities in microblogs. In *Proceedings of WSDM*, pages 45–54.
- Pendleton, S. (1998). Rumor research revisited and expanded. *Language & Communication*, pages 69–86.
- Peterson, W. and Gist, N. (1951). Rumor and public opinion. *American Journal of Sociology*, pages 159–167.
- Petkova, D. and Croft, W. (2006). Hierarchical language models for expert finding in enterprise corpora. In *Proceedings of ICTAI*, pages 599–608.
- Qazvinian, V., Rosengren, E., Radev, D., and Mei, Q. (2011). Rumor has it: Identifying misinformation in microblogs. In *Proceedings of EMNLP*, pages 1589–1599.
- Romero, D., Meeder, B., and Kleinberg, J. (2011). Differences in the mechanics of information diffusion across topics: Idioms, political hashtags, and complex contagion on twitter. In *Proceedings of WWW*, pages 695–704. ACM.
- Yeung, C., Gibbins, N., and Shadbolt, N. (2007). Tag meaning disambiguation through analysis of tripartite structure of folksonomies. In *Proceedings of WI*, pages 3–6. IEEE Computer Society.
- Zhang, J., Ackerman, M., and Adamic, L. (2007). Expertise networks in online communities: structure and algorithms. In *Proceedings of WWW*, pages 221–230.
- Zheng, Z., Li, F., Huang, M., and Zhu, X. (2010). Learning to link entities with knowledge base. In *Proceedings of HLT-NAACL*, pages 483–491.
- Zhou, Y., Cong, G., Cui, B., Jensen, C., and Yao, J. (2009). Routing questions to the right users in online communities. In *Proceedings of ICDE*, pages 700–711.