

# iASK: A Distributed Q&A System Incorporating Social Community and Global Collective Intelligence

Guoxin Liu, *Student Member, IEEE*, and Haiying Shen\*, *Senior Member, IEEE*

**Abstract**—Traditional web-based Question and Answer (Q&A) websites cannot easily solve non-factual questions to match askers' preference. Recent research efforts begin to study social-based Q&A systems that rely on an asker's social friends to provide answers. However, this method cannot find answerers for a question not belonging to the asker's interests. To solve this problem, we propose a distributed Q&A system incorporating both social community intelligence and global collective intelligence, named as iASK. iASK improves the response latency and answer quality in both the social domain and global domain. It uses a neural network based friend ranking method to identify answerer candidates by considering social closeness and Q&A activities. To efficiently identify answerers in the global user base, iASK builds a virtual server tree that embeds the hierarchical structure of interests, and also maps users to the tree based on user interests. To accurately locate the cooperative experts, iASK has a fine-grained reputation system to evaluate user reputation based on their cooperativeness and expertise, and uses a reputation based reward strategy to encourage users to be cooperative. To further improve the performance of iASK, we propose a weak tie assisted social based potential answerer location algorithm and an interest coefficient based uncategorized question forwarding algorithm. To further improve the response quality and cooperativeness, we propose a reputation based reward strategy that motivates users to answer questions from unknown users. Experimental results from large-scale trace-driven simulation and real-world daily usages of the iASK prototype show the superior performance of iASK. It achieves high answer quality with 24% higher accuracy, short response latency with 53% less delay and effective cooperative incentives with 16% more answers compared to other social-based Q&A systems. The results also show the effectiveness of the enhancement algorithms in improving the performance of iASK.

**Keywords:** Distributed systems, Question and answer systems, Social networks, Information search

## 1 INTRODUCTION

Question and Answer (Q&A) systems play a vital role in our daily life as one of the most important information sources. Q&A websites such as Ask.com [1], Answers.com [2], Yahoo! Answers [3], stackoverflow [4] and Quora [5] publish the questions on the web, making them available to all users to answer. These Q&A websites may allow users to build directed relationships, such as follower-followee. However, they cannot easily solve non-factual questions [6], because followers are unaware of their followees' personnel preferences. The non-factual questions here mean the questions without specific correct answers, such as questions about opinion or suggestion. Also, due to the anonymous global users, a question may not receive answers or the response delay may be long, and the provided answers may not be trustable (such as spam) or accurate [7]. To address these problems, more and more research efforts begin to study social-based Q&A systems [6], [7], [8], [9], [10], [11], [12]. Since social friends always share common-interests and they trust and like to help each other, the social-based Q&A systems rely on an asker's social friends to provide answers.

However, users sometimes may be more likely to seek the information not related to their social community. For instance, a researcher in "distributed systems" may ask questions on "social networks"; a football fan at New York may already know much information about the football sports in New York, but needs suggestions when he decides to watch a melodrama in New York. Then, it may be difficult to find the best answerers from an asker's social community for questions irrelevant to this social community. Indeed, previous social network studies show that weak ties play a more dominant role in the dissemination of information online than strong ties in social network [13], [14]. By limiting the search scope to a user's strong ties, it confines the Q&A activities within individual social communities and prevents the knowledge sharing between different social communities. Therefore, neither a pure social-based Q&A system nor a global Q&A website suffices as a both comprehensive and personalized Q&A system. Thus, we face a challenge of *connecting different social communities to fully utilize the cohesive power of weak ties for users to efficiently receive answers outside of their social communities*.

To solve this challenge, we propose a unified system that incorporates social community intelligence and global collective intelligence into a single distributed Q&A system, named as iASK. Compared to other social-based Q&A systems, iASK is the first work that uses the global collective intelligence to complement the social community intelligence in order to efficiently and accurately locate potential answerers outside the asker's social communities. When an answer

- \* Corresponding Author. Email: [hs6ms@virginia.edu](mailto:hs6ms@virginia.edu); Phone: (434) 924-8271; Fax: (434) 982-2214.
- Haiying Shen is with the Department of Computer Science, University of Virginia, Charlottesville, VA, 22904. E-mail: [hs6ms@virginia.edu](mailto:{hs6ms}@virginia.edu) [guoxinl@clemsun.edu](mailto:{guoxinl}@clemsun.edu)

cannot be found within the social network of an asker, it is forwarded to the global user base. iASK does not simply combine the previously proposed social-based Q&A system and global Q&A website platform. It improves the response latency and answer quality (trust and accuracy) in both the social domain and global domain. In the social domain, by using neural network, iASK considers multiple factors (e.g., response delay, quality, social closeness) in answerer candidate identification, and also gives users options to set different priorities on the factors. In the global domain, there exist three challenges. First, the system must identify potential answerers in an efficient and scalable manner. Second, it is important to identify potential answerers that can provide accurate and trustable answers and are willing to answer the question. Third, it is critical to encourage users to cooperatively answer questions. To handle the first challenge, iASK builds central servers into a virtual server tree that embeds the hierarchical structure of interests (i.e., categories). In iASK, interests not only includes long term interests (i.e., music, book, movie), but also includes short term activities (i.e., job hunting, falling in love). It also classifies the global user base based on user interests and maps the user groups to the virtual servers, so that the potential answerers in a specific interest can be efficiently located along the tree. To handle the second and third challenges, iASK has a fine-grained reputation system to evaluate user reputation based on their cooperativeness and expertise, and a reputation-based reward strategy. In iASK, question forwarding failure may happen due to node dynamism such as node off-line or delivery failure. To improve fault tolerance, we let receivers send an acknowledge to the sender when it receives the message. If a sender does not receive an acknowledge from a node in TTL (time-to-live), it assumes that the node fails and then choose another node to send the message.

Our contributions can be summarized as follows:

- (1) A Q&A system structure that incorporates both social community and global collective intelligences, which complement each other in potential answerer search.
- (2) A neural network based friend ranking method that considers multi-factors to identify answerer candidates in the social network that can provide quick and accurate response. It further provides users the flexibility to choose candidates based on their preference priorities on different factors.
- (3) A virtual server tree in the central servers to efficiently locate answerer candidates in the global user base. Each virtual server manages users in a fine-grained interest and is responsible for locating the answerer candidates in this interest.
- (4) A fine-grained reputation system that accurately locates cooperative global experts to answer questions.
- (5) A weak tie assisted social based potential answerer location algorithm that finds the social community of the interest of a question to locate potential answerers when an asker's question is not in his/her interest (i.e., social community).
- (6) An interest coefficient based uncategorized question forwarding algorithm that forwards a question to the users that can better categorize the interest of a question when the question's asker has limited knowledge to identify the interest of the question.
- (7) A reputation-based reward strategy that encourages

global experts to be cooperative. A higher reward for answerers with a higher reputation motivates users to be cooperative in question responding.

(8) Experimental results from large-scale trace-driven simulation and real-world daily usages of the iASK prototype confirm iASK's superior performance. It achieves high answer quality with 24% higher accuracy, short response latency with 53% less delay and effective cooperative incentives with 16% more answers compared to other social-based Q&A systems.

The rest of the paper is organized as follows. Section 2 presents related work. Section 3 presents the design of the iASK system and describes our strategies. Section 4 shows the trace-driven simulation results of iASK compared to other systems. Section 5 presents a real implementation of the iASK prototype, and demonstrates iASK's performance in the wild testing. We conclude this paper with remarks on the future work in Section 6.

## 2 RELATED WORK

Recently, many research efforts began to study social-based Q&A systems [6], [7], [8], [9], [10], [11], [12]. The systems in [6], [7] are based on broadcasting. Morris *et al.* [6] studied the answer quality and response speed of questions asked through status messages in an online social network as well as how to format questions in order to improve the performance. By posting questions on the status wall, a user can broadcast the questions to all of his/her friends. Harper *et al.* [7] investigated the question quality predictors, and found that the reward strategy and community networks lead to better answer quality. The works in [8], [9] are centralized based systems that identify the most appropriate friends of a user to answer his question. These works and [11] also studied the influence of different factors (e.g., users' profiles, system interactions and community size) in the social networks on Q&A performance. The study results lay the foundation of social-based Q&A systems to leverage social network properties in the design. However, a broadcasting method generates high overhead and a large number of received questions make users hard to find what they can answer. Centralized methods have problems of single point of failure, higher bandwidth and server maintenance costs [10]. Zhang *et al.* [12] proposed an expert finding mechanism coupling with profile matching and social acquaintance prediction methods in order to forward referral requests through social links to experts. SOS [10] is a distributed Q&A system based on a social network that forwards questions in a distributed manner in an asker's social network, and uses knowledge engineering techniques to find the potential answerers of questions in the social network. Different from SOS and all the previous Q&A systems, iASK focuses on incorporating social community intelligence and global collective intelligence to find answerer candidates for higher user quality of service (QoS) (i.e., lower response latency, higher accurate and trustable answers).

The works in [15], [16], [17], [18] focus on locating experts and authoritative users as potential answers for Q&A systems. Zhang *et al.* [15] measured the performance of a set of network-based algorithms for finding experts

on a large-size social network, and found several structural characteristics in the social networks that affect the algorithms' performance for online communities. In [16], the reputation of answerers is calculated in Q&A systems to increase the credibility of answers. In [17], authoritative users for specific question subjects are discovered in order to improve the quality of answerers and answer ranking. In [18], an Opinion-based Cascading (OC) model is proposed to identify the user with positive opinions of a product promotion, and by spreading promotions to these users, OC maximizes the spread of positive influence. Different from these works, iASK's fine-grained reputation system considers more factors for more accurate reputation evaluation, and it further uses the reputation system in its reward strategy to incentivize users to respond to non-friends.

Social networks also have been leveraged for search engines [19], [20], [21], [22], [23], [24] in recent years. Amityay *et al.* [19] used enhanced faceted search engine to query social entities (including friends and social bookmarks), the scores of which are evaluated based on the relevance of their documents to the query. David *et al.* [20] re-ranked the search results by calculating their relevance with individuals in the requester's social network. Kolay *et al.* [21] studied the usefulness of social bookmarked URLs for search engines to find qualified contents on the Web. Bao *et al.* [22] observed that social annotations can represent the popularities of the corresponding web pages. Accordingly, they proposed a SocialPageRank algorithm to calculate the popularity of web pages based on social annotations. Evans *et al.* [23] identified searching as a social activity and demonstrated that social interactions can help improve the search results. Carretero *et al.* [24] proposed Geology, which leverages a gossip protocol to gather neighbor information (activities, friends and interests) from social networks in order to recommend locations where a user might be interested in. Quora [25] as a commercial Q&A system forwards questions based on user's follower and interests. It hosts its service in Amazon Clouds. While iASK incorporates both social community intelligence and global collective intelligence to locate the expertise to forward the questions, and is built on decentralized social networks. These social-based search engine works focus on finding most relevant contents for a Web search, while our work focuses on finding potential answerers that are most likely to provide accurate answers quickly.

Social networks have been used for efficient and cooperative file sharing and distribution in peer-to-peer (P2P) networks [26], [27], [28], [29]. Cheng *et al.* [26] proposed NetTube, which clusters peers to social communities and leverages them for efficient short video sharing. Wang *et al.* [27] studied the propagation patterns of social video contents in social networks, and based on these patterns, they proposed a social-aware content dissemination method with a hybrid edge-cloud and a peer-assisted architecture to facilitate the video sharing. Shen *et al.* [28] proposed a social network-aided efficient live streaming system, which leverages social friends to connect to new video channels in order to release the load of centralized servers. Zhang *et al.* [29] identified transient connected components in a social graph, and leveraged them to disseminate data in mobile

social networks. These works cannot be directly used to locate the best cooperative potential answerers for questions, because they do not consider the users' expertise to answer a specific question, while iASK does.

### 3 iASK: INCORPORATING SOCIAL COMMUNITY AND GLOBAL COLLECTIVE INTELLIGENCE

#### 3.1 Design Rationality

The QoS of a Q&A system depends on whether an asker receives answers, the response latency, answer quality and whether the answers match the asker's needs. The QoS of Q&A systems can be improved by leveraging social networks due to social friend properties. It can improve the answer quality [7] since the friendship is altruistic and trustable [30]. Also, friends in an online social network tend to share similar interests, and be clustered based on their interests [31]. Friends inside the same community may know the asker well so they can provide with satisfied answers. Thus, the friends are better potential answerers for non-factual questions to match askers' personnel preferences and personalized needs. For example, in real life, the persons a student resorts to for answers of questions such as "Is the computer organization qualify exam in our ECE Department difficult?" are usually those in his social community in the ECE Department at his university. Therefore, we can leverage social community intelligence to solve the questions based on interest topics.

It is critical to identify potential answerers in an asker's social community that can provide high-quality answers. Inside the social community, the interaction frequencies between a user and his friends are largely different and vary over time [32], which means that the willingness, availability and trust of a user's different friends to answer his questions need to be evaluated individually and updated over time. iASK considers the dynamic social interactions, which represent friend social closeness, and other Q&A activity factors (e.g., response rate, response delay, answer quality) to identify friends who are willing and trustable to provide answers [23]. iASK also allows users to set different weights on these multi-factors based on their preference to rank friends for potential answerer identification.

In real life, users also ask questions outside of their social communities, so the questions may not be answered within a user's social community, as indicated in Section 1. Posting a question to the web and passively waiting for answers as in current web-based Q&A websites (e.g., Ask.com, Answers.com and Yahoo! Answers) cannot guarantee timely and high-quality answers. In order to pro-actively find appropriate answerers, users need to forward questions to the global experts of these subjects. iASK fulfills this task by incorporating the global collective intelligence to complement the social community intelligence in order to increase the probability that a question is successfully resolved. By a "resolved question", we mean that the asker have received best answers with respect to this question. Unlike the trustable and altruistic social community, the global domain needs another strategy to locate and motivate users who are willing and able to answer questions to unknown unfamiliar users. iASK has a fine-grained reputation system to evaluate user reputation based on their cooperativeness



and expertise, and a reputation based reward strategy to provide more effective cooperative incentives.

### 3.2 System Architecture

iASK is an online social-based Q&A system. Users register their profiles as in online social networks, including interests, education and so on, and build their social networks. There are two types of social networks in iASK: i) friend network as in Facebook, and ii) contact-fan (i.e., followee-follower) network as in Yahoo! Answers. The friendship is bidirectional, being used to locate potential answerer by leveraging social community intelligence; the contact-fan relationship is unidirectional, being used for leveraging global collective intelligence. In our real-world software development, iASK has 5 predefined categories (i.e., music, book, movie, television and research) and 40 subcategories (e.g., pop music, data mining). These first 4 categories and their sub-categories are collected from the Yahoo! questions/answers trace [11]. The research and its sub-categories in the real implementation are collected from all topics of all participate users. From the lists, users choose their interests and question topics. Users also can enter any new category and subcategory under a category as their interests and question topics, and the redundant user-defined categories will be combined based on synonyms. The user can belong with multiple community or interest categories. For example, if one user asked/answered questions in  $m$  categories, then the user belongs to  $m$  social communities.

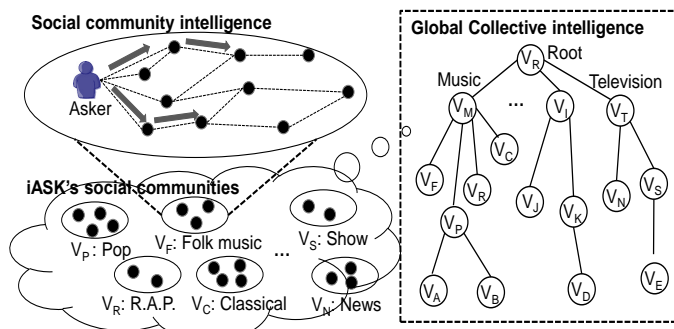


Fig. 1: The architecture of iASK.

iASK incorporates potential answerer location strategies in both the social community intelligence (within an asker's social communities) domain and the global collective intelligence domain (outside an asker's social communities) that are likely to provide high-quality answers in time. Figure 1 shows the high-level architecture of iASK based on the two domains. If a question cannot be solved within an asker's social communities, the question is forwarded to global collective intelligence. In the social community intelligence domain, it has a neural network based friend ranking method to identify potential answerers to forward a question in a distributed manner. In the global collective intelligence domain, it has a virtual server tree that helps to locate potential answerers with the interest of the question. We adopt the concept of virtual server from [33]. All virtual servers form a tree that mirrors the filiation among categories and subcategories. Therefore, each virtual server represents a group of all users with a specific interest category or subcategory, and is hosted by a physical server. That is, a virtual server's jobs, including user join and leave

management and expert location, are executed by its host physical server. To avoid user redundant efforts to forward or answer the same questions and hence reduce network traffic in both social community and global collective intelligence domains, a duplicated received question from the same asker is dropped. In order to choose answerers that will provide high-quality answers quickly and to encourage users to provide high-quality answers quickly, iASK has fine-grained reputation evaluation and reputation based reward strategies. We introduce the details of each component of iASK below.

### 3.3 Integrated Social and Global based Answerer Location

When a user asks a question, he specifies the question's topic by selecting or entering an interest. If the interest is not within the asker's interests, it is directly forwarded to the central servers. Otherwise, it is forwarded to the best  $K$  answerer candidates among his friends having this interest. Section 3.3.1 introduces how to select the answerer candidates. When a user receives a question, if he cannot answer it, he further forwards it to his friends. After the question is forwarded by TTL hops, the receiver forwards the question to the central servers. After the central servers receive a question, based on the virtual server tree, the question is then efficiently forwarded to the virtual server which manages the group of all users in the system with this interest. The responsible virtual server chooses  $K$  experts based on their reputations in this interest. The details of the global answerer candidate identification are presented in Section 3.3.2. If the answer is still not answered satisfactorily, the question is posted to the question forum as in Yahoo! Answers.

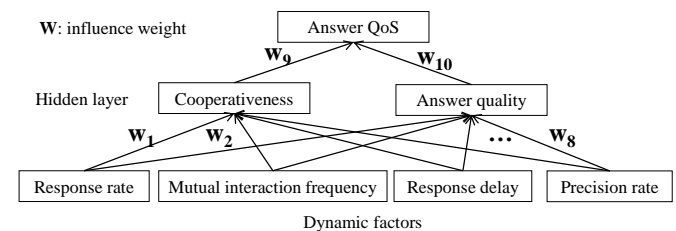


Fig. 2: The neural network model for friend ranking.

#### 3.3.1 Social based Potential Answerer Location

To evaluate the qualification of an asker's friends to answer his question, iASK considers the following factors: answer quality, willingness (cooperativeness) and response delay. In iASK, an asker gives a precision score ranging from 0 to 5 to each received answer [3], which represents the accuracy of this answer. Since a friend may have different degrees of knowledge in different interests, for each interest  $I_j$ , we measure the friend's precision rate to evaluate his answer quality in this interest. To accurately reflect a friend's current qualification to be an answerer, for each of user  $u_a$ 's friends (denoted by  $f_i$ ), iASK periodically calculates the following social and Q&A activities: response rate, mutual interaction frequency, response delay and precision rate.

(1) *Response rate ( $R_{f_i}$ ):*  $R_{f_i} = ACK f_i / Q f_i$  ( $ACK f_i$  is the response from  $f_i$  and  $Q f_i$  is the question sent and forwarded to  $f_i$ ). It is measured by the percentage of questions of  $u_a$  answered or forwarded by  $f_i$ , because forwarding a

question is also considered as a responding behavior. This metric reflects the cooperativeness of a friend.

(2) *Mutual interaction frequency* ( $M_{f_i}$ ):  $M_{f_i} = ACK_{f_i}/T$  ( $T$  is the time of prediction). It is measured by the number of interactions between  $f_i$  and  $u_a$  in a unit time period. This metric reflects the social closeness of the two users.

(3) *Response delay* ( $D_{f_i}$ ):  $D_{f_i} = \sum_{j \in [1, ACK_{f_i}]} D_{f_i}^j / ACK_{f_i}$ . It is measured by the average delay of all interactions between  $f_i$  and  $u_a$  per unit time. This metric reflects the responsiveness of interactions and Q&A activities between the two users.

(4) *Precision rate* ( $P_{f_i}^{I_j}$ ): It is measured by  $P_{f_i}^{I_j} = G_{f_i}^{I_j}/G$ , where  $G$  is the upper bound precision score of an answer in the system, and  $G_{f_i}^{I_j}$  is the average precision score of all answers from friend  $f_i$  under interest  $I_j$ .

The response rate and mutual interaction frequency represent the willingness of friends to answer or forward a question [6]. The response delay represents the timeliness of a friend's response. The precision rate reflects the degree that a friend's answer can precisely answer the user's question.

The satisfaction score of an answer is given by the asker based on different answer QoS factors including the response delay, answer precision, interaction frequency and response rate. If the asker does not receive an answer from an identified potential answerer, he gives 0 precision and satisfaction scores to this user. The 0 answer score helps exclude users who are not appropriate answerers and hence increase the probability to find good answerers. This answer score represents the overall answer QoS to users. iASK aims to identify potential answerers that will receive high answer scores (i.e., high satisfaction) from the asker. For this purpose, iASK depends on a neural network [34], [35]. Due to the dynamism of social networks that iASK leverages, we need a scheme that can dynamically derive the final decision with real time training. As shown in Figure 2, to find out the influence weight of each factor on the QoS of friends' answers, denoted as  $W_{u_a} = \langle w_1, \dots, w_{10} \rangle$ . The training process is the process to determine the  $W_{u_a}$  vector and the non-linear relationship between the four factors and the answer QoS. When a user needs to identify  $K$  friends in his social network to forward a question, he uses the trained neural network to calculate the output QoS value for each friend. Then, he chooses the  $K$  friends having the interest of the question and the highest QoS values to forward the question. Note that  $W_{u_a}$  determined by the training process represents the general influence degree from the factors on the QoS derived from many friends' activities. However, a user may have his own preference priorities on measuring the QoS. For example, users asking simple questions in urgency may prefer short response delay than the precision rate. Also, a user's preference may change over time. Thus, an asker can adaptively adjust the value of  $W_{u_a}$  when evaluating the QoS of each of his friends:

$$\forall i \in [1..8], w_i = \alpha_i w_i \wedge \sum_{i=1}^8 \alpha_i = 1.$$

In this case, the asker needs to forward the question along with his own specified  $W_{u_a}$  to identified top  $K$  friends. Each question receiver uses the received  $W_{u_a}$  in selecting the top  $K$  friends to forward the question in order to meet the QoS

preference of the asker.

The  $W_{u_a}$  vector is updated periodically through training. The training time period represents a tradeoff between the sensitivity of environment variance and computation cost for training. A smaller time period leads to more accurate derived  $W_{u_a}$ , but also generates a higher computation cost due to the frequent updating. When a user receives a question, if he cannot answer it, he further forwards it to his friends. After the question is forwarded by TTL hops, the receiver forwards the question to the central servers.

### 3.3.2 Global based Potential Answerer Location

iASK builds the central servers into a virtual server tree overlay to efficiently identify potential answerers that have the question's interest in the global user base. The entire interest space can be classified into pre-defined categories. For example, Yahoo! Answers has 17 categories such as "Pets", "Travel" and "Sports". Each category can be classified into sub-categories, each of which can further be classified into smaller categories and so on. Based on such classification, an interest tree can be established. Assume that in the interest tree, each node has at most  $d$  children. Then, iASK builds a  $d$ -nary virtual server tree, as shown in Figure 3, to map to the interest tree.

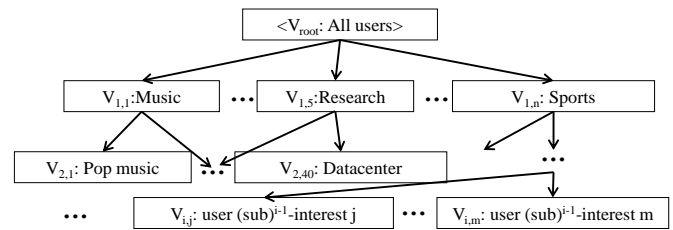


Fig. 3: The virtual server tree in the central servers.

In the tree,  $v_{i,j}$  represents the  $j^{th}$  virtual server on the  $i^{th}$  level of the tree. Each child is responsible for a sub-category of the category in its parent. Each physical server runs a number of virtual servers, and iASK can deploy its virtual server tree to a cloud. This tree is a locality-aware tree, where virtual servers in the same subtree are physically close to each other and also physically close to their parent in order to reduce the communication overhead.

A virtual server responsible for category interest  $I_i$  records all users with interest  $I_i$ , and also is responsible for locating the answerer candidates among these users for questions in interest  $I_i$ . When user  $u_a$  enters his interests, the system translates each interest to identifier  $v_{i,j}$  in the tree accordingly. The virtual server with identifier  $v_{i,j}$  in the tree becomes  $u_a$ 's server holder. The server holder stores the information of  $u_a$ , and the information is forwarded in the bottom-up manner until reaching the tree root and stored in the virtual servers along the path. When  $u_a$  sends a question to the central servers,  $u_a$ 's server holder finds answerer candidates for the question. Specifically, this question is forwarded in the bottom-up manner until it reaches a virtual server responsible for the question's interest. Then, this question is forwarded in the top-down manner until it reaches a virtual server responsible for this question's smallest interest category. This virtual server then identifies the answerer candidates from its responsible users. In this way, the workload of answerer candidate identification is

distributed among the different central servers, thus avoiding single point of failure and workload bottlenecks.

For the candidate identification, a virtual server  $v_{i,j}$  needs to rank its responsible users, i.e., the users with interest  $v_{i,j}$ . In order to measure a user's cooperative behavior and his expertise in the category of a question, a virtual server calculates each user's reputation as introduced in the next section, and then selects the users with the highest reputations as potential answerers.

### 3.3.3 Weak Tie Assisted Social based Potential Answerer Location

A user may ask a question outside of his/her interests and hope that the question can be answered from the social community of the question's interest rather than from an expert in the global collective intelligence domain. For example, an asker without a soccer interest at Clemson University may ask a question "how is the soccer team at Clemson University?". Based on our previous algorithm, since the asker is not within the soccer interest community, the question will be answered in the global collective intelligence and may be sent to a soccer expert. However, the expert may not have localized or personalized knowledge about the soccer team at Clemson University though (s)he has a wide knowledge on the soccer teams nationally or globally. Therefore, we need to find a potential answerer in the soccer social community at Clemson University to answer this question through social links, who has a higher probability to answer the question. This social community should be close to the asker's social communities since they belong to the same organization. In a nutshell, for an asker's question which is not in his/her interests but needs localized or personalized answers, the potential answers can be found from the social community that has the question's interest and also is close to the asker's social community.

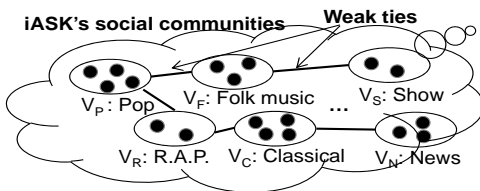
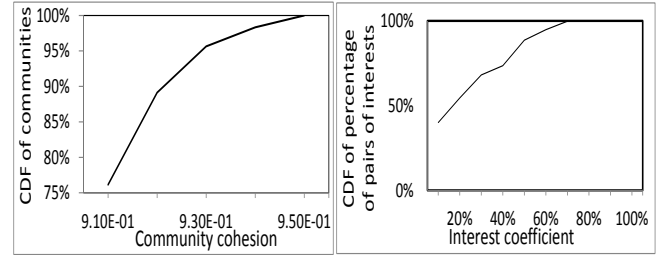


Fig. 4: The architecture of iASK.

In this section, the question we discuss does not belong to the interest of the asker's social community. Since the asker is outside of the social community that has the question's interest, the question needs to be forwarded to a user inside this social community first. To achieve this, as shown in Figure 4, iASK leverages the weak tie, which is a bridge connection between two social communities [36]. Before we present the weak tie assisted social based potential answerer location algorithm, we first conduct trace data analysis to confirm the existence of weak ties.

We used the Yahoo! Answers question/answer trace data from [11] as a showcase to study the existence of weak ties between interest categories in a Q&A system. The trace includes 119,175 users, their profiles and their asked and answered questions. Recall that in iASK, we cluster users according to different interests into different social communities. We form all users that asked or answered at least a question in an interest category into a



(a) CDF of community cohesion (b) CDF of interest coefficient  
 Fig. 5: Design rationality.

social community. Finally, we created 26 communities. For such a social community  $i$ , we use  $A_i$  to denote the total number of question/answer pairs that are within this social community's interest. We use  $A_{i,o}$  to denote the number of questions which are not in this social community's interest but have been asked or answered by the users in this social community. We can calculate the community cohesion of a social community in iASK by  $\frac{A_i}{(A_i + A_{i,o})}$ . The cohesion of a community represents its density. If it is high, most of the asking and answering communications of the community users are within the community, which indicates that the inter-community communications from this community are not frequent, and vice versa. The weak ties are the communications to other communities from this community; that is, the users in this community ask questions belonging to other communities and the questions are answered by the users in other communities.

Figure 5(a) shows the CDF of social communities versus the community cohesion. It shows that all 26 communities have cohesion between 91% and 96%, and 95% of them have community cohesion larger than 93%. The figure indicates that most of the question asking and answering activities are within a community, which shows a strong cohesion of all communities. However, it also shows that there indeed exist weak ties among communities, through which the question can be forwarded to other communities. The cohesion strength is negatively proportional to the strength of the weak ties of a community [11]. Since the community cohesion is high, the communication frequency between communities is rare. Therefore, we need an algorithm to forward a question created in a community that does not have the question's interest to a specific weak tie connecting to the target community that has the question's interest.

It is not effective and efficient to randomly select a friend to forward the question hop by hop to find the weak tie, since most of the links within a social community in a Q&A system is within a social community. To find the weak tie towards the target community, a user can forward his/her question to a friend having more interests that are frequently asked by the users in the ask's community. Then, the probability that the question is forwarded to its target community is increased. To this end, the root of the virtual server tree periodically collects all Q&A activities. Based on the collected activity information, it calculates the interest coefficient between each pair of interests  $I_i$  and  $I_j$ , denoted by  $C_{I_i, I_j}$ , as:

$$C_{I_i, I_j} = \frac{U_i \cap U_j}{U_i \cup U_j}, \quad (1)$$

where  $U_i$  and  $U_j$  denote the set of users in the system asking or answering a question in interest  $I_i$  and interest



$I_j$ , respectively. The virtual server root then publishes the interest coefficient through the virtual server tree to all users. For a friend  $f_j$ , we define a metric called accumulated interest coefficient with the interest  $I_t$  of the question (i.e., the interest of the target community), denoted by  $S_{f_j, I_t}$ . It is calculated by

$$S_{f_j, I_t} = \sum_{I_j \in V_{f_j}} C_{I_t, I_j}, \quad (2)$$

where  $V_{f_j}$  is the set of all interests of friend  $f_j$ . When an asker needs to forward his/her question to a weak tie, it calculates  $S_{f_j, I_t}$  for each of his/her friends based on the coefficient between interest pairs. Then, the asker finds the friend with the interests having the highest accumulated interest coefficient with the interest  $I_t$  of the question. Intuitively, a friend with more interests similar to interest  $I_t$  has a higher probability to have a friend or is much socially closer to users inside the social community of  $I_t$ . As a result, the question has a high probability to be forwarded to the target community.

The effectiveness of this algorithm for finding potential answerers in other social communities depends on whether there are different interest coefficients existing between different pairs of interests (i.e., social communities). To verify this, we regarded different question categories in Q&A trace [11] as interests and calculated the interest coefficients between each pair of interests. Figure 5(b) shows the CDF of interest pairs versus the interest coefficient. It shows that 40% of the interest coefficient are less than 10%. However, there are 5% of interest pairs that have coefficient larger than 50%. It indicates that the interest coefficients are different. That is, some interests are much more closely related to a specific interest than the others. Therefore, forwarding the question to the social community (i.e., interest) for which the users in the asker's community more frequently ask questions should be effective in forwarding the question to its target community.

Below, we present the details of this question forwarding algorithm. When a user needs to forward a question to its friend, it first checks whether (s)he has a friend with the interest  $I_t$  of the question. If such a friend exists, it means that the user has found the weak tie connecting to the target community and (s)he forwards the question to the friend, who will then start a social based potential answerer location as in Section 3.3.1. Otherwise, the user calculates the accumulated interest coefficient  $S_{f_j, I_t}$  between each of his/her friends  $f_j$  and the interest  $I_t$  of the question. It selects the friend with the largest  $S_{f_j, I_t}$  to forward the question to. To avoid flooding, the asker and each forwarder selects top  $N$  friends with the largest  $S_{f_j, I_t}$  to forward the question in order to more efficiently find the weak tie connecting the target community.

### 3.3.4 Interest Coefficient based Uncategorized Question Forwarding

The above algorithm assumes that the interest of the question can be identified by the asker. An asker sometimes may not be able to identify the interest category of his/her question if it is out of his/her knowledge scope. We call such a question *uncategorized question*. For example, an asker without the knowledge of computer science may not be able to distinguish a question belonging to "Cloud Com-

puting" or "High Performance Computing". In this case, the question needs to be first forwarded to a user that can successfully categorize a question. Intuitively, for questions whose interests are farther away from the asker's interests, the asker has less capability to categorize it and vice versa. Therefore, it is reasonable to assume that the interest of the uncategorized question is different from the asker's interests. An asker then should aim to forward his/her uncategorized question to a friend with more interests far from the asker's interests until a forwarder can distinguish the question's interest.

To find a friend with more interests far away from the asker's interest, we use the interest coefficient to measure interest difference between a friend's interests and the asker's interests. To do this, we first determine the interest closeness between asker  $i$  and friend  $j$  by calculating the overall interest coefficient between each interest of friend  $j$  and the set of interests of asker  $i$ . We use  $S_i$  and  $S_j$  to indicate the asker  $i$ 's interest set and friend  $j$ 's interest set, respectively. For each interest  $I_j$  in  $S_j$ , we measure its interest coefficient to each interest  $I_i$  in  $S_i$ . We use the maximum interest coefficient as  $I_j$ 's interest similarity to  $S_i$  as  $D_{I_j, S_i} = \max\{C_{I_i, I_j}\}_{I_i \in S_i}$ . We then use the average interest similarity of all interests within  $S_j$  to calculate the interest difference between the friend  $j$  and asker  $i$  as

$$1 - \frac{\sum_{I_j \in S_j} D_{I_j, S_i}}{|S_j|}. \quad (3)$$

During each forwarding, a friend with the largest interest difference is selected to forward the uncategorized question. Similarly, to avoid flooding, the asker and each forwarder selects top  $N$  friends with the largest interest difference to forward the question in order to more efficiently find a user who can identify the interest of the question.

When a question receiver can categorize the question, if the forwarder is within the target interest's community, the social based potential answerer location algorithm introduced in Section 3.3.1 is used to find the potential answerer. Otherwise, the weak tie assisted social based potential answerer location algorithm introduced in Section 3.3.3 is used to find the target community.

### 3.4 A Fine-Grained Reputation System

A virtual server calculates each user  $u_j$ 's rank score by two different reputations: the global reputation denoted as  $R_{u_j}^g$ , and an expertise reputation in an interest  $I_i$  denoted as  $R_{u_j}^{I_i}$ . The root server, which holds all users, is responsible for calculating  $R_{u_j}^g$  for every  $u_j$  in the system. Recall that iASK has a contact-fan network. As users like to be fans of others who are more knowledgeable than them [11], a more trustable and knowledgeable answerer usually has more fans. Then, the root server considers the global reputations of a user's fans to estimate the user's global reputation:  $\sum_{u_i \in f(u_j)} R_{u_i}^g / |f(u_j)|$ , where  $f(u_j)$  is the set of  $u_j$ 's fans, and  $R_{u_i}^g$  is fan  $u_i$ 's reputation. As in Yahoo! Answers, users select the best answer for each question in iASK. We use  $B_{u_j}$  to denote the percentage of  $u_j$ 's best answers in his answers, which reflects  $u_j$ 's expertise. Then,  $R_{u_j}^g$  is calculated as the harmonic mean of user  $u_j$ 's expertise ( $B_{u_j}$ ) and the reputations of his fans.

$$R_{u_j}^g = \frac{1}{\frac{1}{2} * \left( \frac{1}{B_{u_j}} + \frac{1}{\sum_{u_i \in f(u_j)} R_{u_i}^g / |f(u_j)|} \right)}, \quad (4)$$

The virtual server for interest  $I_i$  calculates  $R_{u_j}^{I_i}$ :

$$R_{u_j}^{I_i} = N_{u_j}^{I_i} / N^{I_i}, \quad (5)$$

where  $N_{u_j}^{I_i}$  is the number of best answers under interest  $I_i$  provided by  $u_j$  and  $N^{I_i}$  is the total number of best answers in interest  $I_i$ .  $R_{u_j}^{I_i}$  reflects  $u_j$ 's expertise in interest  $I_i$ . The virtual server requests the global reputations of its responsible users from the root server and calculates the harmonic mean of  $R_{u_j}^g$  and  $R_{u_j}^{I_i}$  as the final reputation of each user  $u_j$ :

$$R_{u_j} = \frac{1}{\frac{1}{2} * (\frac{1}{R_{u_j}^g} + \frac{1}{R_{u_j}^{I_i}})}. \quad (6)$$

It identifies the top  $K$  users with the highest  $R_{u_j}$  values as the answerer candidates and sends the question to them. If there is no best answer after a timeout, the virtual server posts the question on the forum, where each user can see and answer the question.

### 3.5 Reputation based Reward Strategy

We propose an adaptive reward strategy based on user reputations to further improve the answer quality and response cooperativeness and timeliness. In iASK, an asker needs to provide virtual currency for a question, and users receive rewards for answering questions. Since the friendship is altruistic and trustable [30], our reputation based reward strategy is mainly to motivate users to answer questions from non-friends. A user can set the virtual currency threshold  $v_t$  for answering a question, so questions with reward less than  $v_t$  will not be sent to this user. In this way, the high-reputed experts will not be disturbed by a vast number of questions by setting a high  $v_t$  even though they are always chosen as the potential answerers. When an asker asks a question, he chooses a reward  $C_o$  in the range of [2...10] as virtual currency that he will pay. The reward represents the response quality and timeliness the asker expects, since the higher the reward is, the higher reputed answerers the question will be sent to. Each answerer receives the default smallest reward as 2, and an answerer  $u_j$  for the first best answer receives high reward, which is calculated as:

$$C_{q_j} = \text{Max}\{C_o, (C_o + R_{u_j}^g)/2\}. \quad (7)$$

Since a higher global reputation  $R_{u_j}^g$  leads to higher final reward to an answerer [37], [38], this strategy motivates users to gain reputation as high as possible in order to increase their rewards. Thus, users are incentivized to provide high quality (trustable and accurate) answers timely. This is very important to the sustainable development of Q&A systems considering that 64% are one-time users (i.e., users with only one question) [11].

## 4 PERFORMANCE EVALUATION

We conducted trace-driven experiments to evaluate the performance of iASK. We used the Yahoo! Answers question/answer trace data from [11] and Facebook user friendship trace from [39]. The Yahoo! Answers trace has 119,175 users and their profiles, including the number of contacts and fans, and the asked and answered questions. The Facebook trace has a list of all user-to-user links for 60,101 unique users from the Facebook New Orleans networks. We constructed a virtual server tree overlay with three layers according to the categories and subcategories in the trace.

To construct the social network in the simulation, we randomly selected 100,000 users from the Yahoo! Answers trace. For each user, we regard his/her most frequent select subcategories as his/her interests, which include at least 80% of his/her total questions. The distribution of the number of friends of all users follows the Facebook trace. According to the trace, a user  $u_i$  has  $c_i$  contacts and  $f_i$  fans. To construct the contact-fan network,  $u_i$  randomly selected  $c_i$  other users as his/her contacts. The number of fans of each contact of  $u_i$  should be no larger than  $f_i$ . For the answers of each question, we randomly assigned the best answers to users with the question's subcategory according to the distribution of the number of best answers of each user in this subcategory in the trace. We then randomly assigned the other answers to users that have not been assigned any answers of this question.

The score for a best answer was set to 10, and the score for a non-best answer was set to a random value from [0, 5]. The average precision score of all answers of the same interest of a user indicates the precision rate, which represents the users intelligence on that interest and used as an input for this users QoS. When forwarding questions to friends or global potential answers, the number of selected answerers  $K$  was set to 10. The distribution of response time to a question follows the trace in [10]. When receiving a question, a user decides whether to respond to or drop it based on his/her response rate which is the response rate between two friends follows a bonded Pareto distribution with a lower bound, an upper bound and a shape as 0.2, 0.8 and 2, respectively. In responding, if the user has the answer, the question will be answered; otherwise, the question is forwarded to the potential answerers based on the iASK algorithms. Therefore, based on how many answers received and the total answers, we can derive the recall rate, and similarly based on how many best answers are received, we can get the precision rate. The timeout for a question routing inside the social network was set to 800 minutes. In the weak tie assisted social based potential answerer location algorithm and the interest coefficient based uncategorized question forwarding algorithm, we set the TTL for the question forwarding through social network to 4, and the asker first forwards the question to 10 friends nearby. If the weak tie has been found or the user that can categorize the question has been found, a social based potential answerer location or the weak tie assisted social based potential answerer location algorithm is conducted. Otherwise, the global based potential answerer location is conducted to find the experts in the global collective intelligence domain.

Recall that iASK allows users to set different weights to factors (Figure 2) in QoS calculation for answerer selection. Since current Q&A systems do not have such a function, the weights of all factors of all friends were set to 0.5 initially. Before each experiment, we let each user ask 100 questions to initialize the weights of the factors. We use  $BA$  to denote the set of best answers of asked questions in the simulation, and use  $RA$  to denote the set of retrieved answers in the system from the trace. The following metrics are used to evaluate iASK's performance:

- (1) Response rate. It is the number of successful interactions (including forwarding and answering) divided by the total number of all interactions.



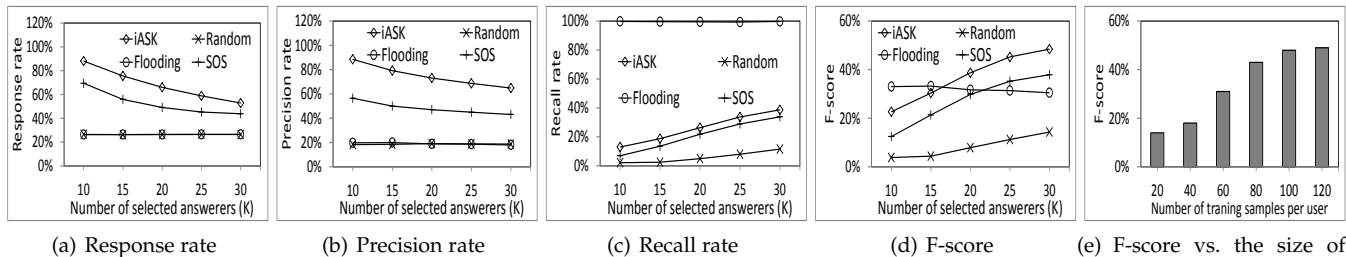


Fig. 6: Effectiveness of different Q&A systems in the social community.

(2) Precision rate. It is defined as  $|RA \cap BA|/|RA|$  to represent the quality of received answers.

(3) Recall rate. It is defined as  $|Unique(BA \cap RA)|/|Unique(BA)|$  to denote the completeness of received answers, where  $Unique(s)$  retrieves the set of all unique elements contained in  $s$ .

(4) F-score [40]. It is defined as  $F = \frac{2 * precision * recall}{precision + recall}$  to measure the overall effectiveness of searching good potential answerers considering both the precision and recall rates.

(5) Response delay. It is time period between asking a question and receiving the first best answer for it.

We compared iASK's friend selection algorithm in the social community intelligence domain with three other algorithms: i) Random [41], which randomly selects  $K$  friends, ii) Flooding [42], which floods questions to all friends, and iii) SOS [10], which select  $K$  friends with highest score calculated by equal weights of social closeness and interest similarity. The Random method can simulate current web-based Q&A websites, in which a question is randomly visited by different users. The Flooding method can simulate previously proposed social-based Q&A systems, in which a question is flooded to all nodes in the social network. To compare the performance of the entire Q&A system, we compared the iASK system incorporating both global collective intelligence and social community intelligence with three other systems: i) Global(Tree) which selects potential answerers using iASK's virtual server tree, ii) Global(Flat) which selects potential answerers based on one-level categories without subcategories, and iii) SOS [10] without a forum to post unsolved questions. Global(Flat) can represent the previously proposed centralized social-based Q&A systems. To show the performance of the weak tie assisted social based potential answerer location algorithm and the interest coefficient based uncategorized question forwarding algorithm, we compared iASK with a modified algorithm, denoted by iASK-Random, in which friends are randomly selected in order to find the weak tie or a user who can categorize the question. We also compared them with Global(Tree) and Global(Flat), respectively, since the question cannot be categorized by the asker in the second scenario. To show the performance of our reputation-based reward strategy, we compared our method denoted as *Varying Reward* with: i) *No Reward*, which does not have an award strategy; ii) *Determined Reward*, which has a fixed award for answering a question.

#### 4.1 Performance in Social Community Intelligence

In this experiment, we measure the performance of iASK's friend selection algorithm in the social community intelligence domain. The number of selected potential answerers at each hop is increased from 10 to 30 with step size of

5. Each user in the system in turn asked one question. In order to measure the sole performance of the friend selection algorithm, askers generated questions within their interests.

Figure 6(a) shows that the response rate follows  $iASK > SOS > Random \approx Flooding$ . In iASK, users choose friends with higher QoS, including the response rate, to answer or forward questions. Therefore, iASK generates higher response rate than others. SOS considers the interest similarity and social closeness in friend selection. Since SOS does not consider the response rate directly, it leads to lower response rate than iASK. Random and Flooding do not consider the response rates of friends, leading to similar lower response rates than SOS. We also see that the response rate of iASK and SOS decreases as the number of selected answerers increases because friends with lower response rates are more likely to be selected. This result implies that iASK's social based answerer identification method is the best to find cooperative friends.

Figure 6(b) shows the precision rate of each method, which follows  $iASK > SOS > Random \approx Flooding$ . Random and Flooding do not consider the answer precision rate of friends, so they have the lowest precision rate in all methods. SOS chooses friends with similar interests as the question, who are likely to give best answers, leading to higher precision rate than Random and Flooding. However, unlike iASK, SOS does not always choose friends with high precision rate due to the large number of friends with this interest. Consequently, iASK has the highest precision rate in all methods. Also, due to the same reason as in Figure 6(a), the precision rates of both SOS and iASK decrease as  $K$  increases. Figures 6(a) and 6(b) together indicate that iASK outperforms other methods regarding both response rate and answer quality.

Figure 6(c) shows the recall rates of all methods, which follows  $Flooding > iASK > SOS > Random$ . Flooding sends a question to all friends, thus it produces the highest recall rate close to 100%. However, Flooding generates many more messages for question forwarding than other methods. Since both iASK and SOS consider interests, they supply many more high-quality answers than Random, leading to a higher recall rate. iASK has a higher recall rate than SOS due to its higher response rate and precision rate as shown in Figures 6(a) and 6(b), respectively. This figure indicates that iASK can resolve more questions with best answers than other non-flooding methods.

Figure 6(d) shows the F-scores of all methods. We see that the F-score follows  $iASK > SOS > Random$ . This is because iASK has the largest precision and recall rates and Random has the lowest rates as shown in Figures 6(b) and 6(c). Flooding has the largest F-score when  $K \leq 15$ , because it generates the largest recall rate (100%) while all

other methods' recall rates are small. However, Flooding generates the largest overhead by flooding the question to all users. When  $K \geq 20$ , iASK and SOS outperform Flooding due to their increasing recall rates. These results indicate that iASK achieves a better overall accuracy than all other methods by considering both precision and recall rates when  $K$  is large.

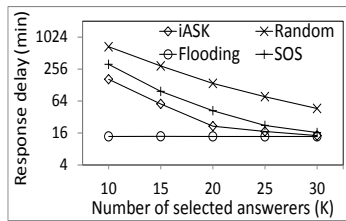


Fig. 7: Response delay of different Q&A systems in the social community.

Recall each user in iASK needs a training process during the social based potential answerer location as introduced in Section 3.3.1. We varied the size of training samples, which are the questions used for training, to measure its effect on the overall searching accuracy of iASK. In this experiment, we set the number of selected answerers to 30. Figure 6(e) shows the F-score of iASK versus the number of training samples per user. We see that the F-score increases as the size of training sample increases. That is because with more training samples, iASK can more accurately calculate the weights in the neural network in Figure 2, which leads to more accurate location of potential answerers. Then, both the recall and precision rates increase, leading to the increase of F-score. The result indicates that a larger number of training samples leads to higher effectiveness of searching potential answerers. Figure 6(e) also shows that F-score stays relatively stable when the number of training samples is larger or equal to 100. It indicates that 100 training samples per user are appropriate to have a good answerer searching effectiveness in iASK.

Figure 7 shows the average response delay for all questions. It follows Flooding < iASK < SOS < Random due to the same reason as in Figure 6(c). This figure indicates that iASK leads to shorter response delay for askers than other non-flooding methods. However, Flooding generates a low precision rate and also high overhead for dispatch messages to all friends in every hop.

#### 4.2 Performance in Global Collective Intelligence

In this experiment, we measure the performance of the iASK system performance without interest coefficient based question forwarding with different user scales. The number of users in the system was increased from 20,000 to 100,000 with step size of 20,000. Different sets of users were randomly chosen from the selected 100,000 users. We assume that each user has equal probability to ask factual and non-factual questions. For non-factual questions, social friends supply better answers than the global users [10]. Thus, if a user is more than two hop social distance away from the asker, the probability to assign a best answer to this user is decreased by one half. The actual response rate of a global user in a virtual server is the smallest actual response rate to all of his/her friends, since friendship is more altruistic and trustable [30].

Figure 8(a) shows that the response rate follows iASK > SOS > Global(Tree) > Global(Flat). iASK has a larger response rate than SOS due to the same reason as in Figure 6(a). Both iASK and SOS depend on the social friends to answer questions, who are more willing to answer questions than strangers as the global users. Thus, they both have higher response rates than the two Global systems. Global(Tree) has a fine-grained user and interest clustering compared to Global(Flat). Since some global users with the highest reputations may have interests in several subcategories rather than all subcategories in a category, these users generate a low response rate when being asked questions in other subcategories. Thus, Global(Tree) is more effective to find global experts than Global(Flat). This figure indicates that iASK is the most effective system to find cooperative answerers by leveraging both social community intelligence and global collective intelligence, and the fine-grained virtual server tree overlay is effective in locating cooperative global experts.

Figure 8(b) shows the precision rate of each system, which follows iASK > Global(Tree) > Global(Flat) > SOS. iASK has the highest precision rate by choosing answerers with high QoS that considers precision rate. Without using the social networks, two Global systems choose global users that may have low precision rate for non-factual answers. Due to the same reason as in Figure 8(a), Global(Tree) generates a better precision rate than Global(Flat). SOS does not directly consider precision rates to locate the experts; thus, it generates the lowest precision rate. This figure indicates that iASK supplies the highest quality answers.

Figure 8(c) shows the recall rate of each system, which follows the same distribution as in Figure 8(b) due to the same reasons. The experimental result confirms that neither a social-based Q&A system nor a web-based global Q&A system can supply a good question recall rate.

Figure 8(d) shows the F-score of each system versus the number of total users. From the figure, we see that the F-score follows iASK > Global(Tree) > Global(Flat) > SOS. That is because both the recall and precision rates of all systems follow the same order in Figures 8(b) and 8(c). The figure indicates that iASK achieves a better overall accuracy than all other methods in searching good potential answerers.

Figure 8(e) shows the average response delay for all systems. It follows Global(Flat)  $\approx$  Global(Tree) < iASK < SOS. iASK and SOS generate longer response delay due to their question routing time over the social network. SOS generates longer response delay than iASK due to the same reasons as in Figure 8(c). Both Figures 8(c) and 8(e) indicate that iASK generates shorter response delay than social-based Q&A systems, and a better recall rate than all others by incorporating both the social community intelligence and global collective intelligence.

#### 4.3 Performance of Interest Coefficient based Forwarder Selection

In this section, we measured these two enhancement algorithms in the same scenario as in Section 4.2 unless otherwise specified. We first measured the performance of the weak tie assisted social based potential answerer location algorithm compared to other algorithms. All questions outside of the asker's interests can be successfully

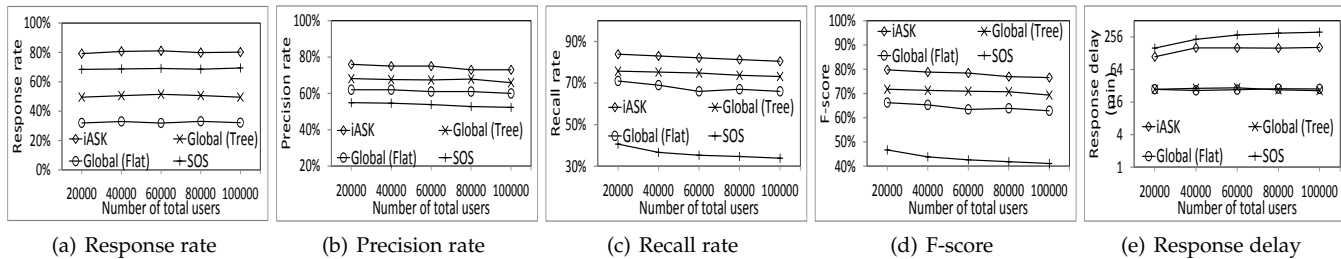


Fig. 8: Effectiveness and efficiency of different Q&A systems.

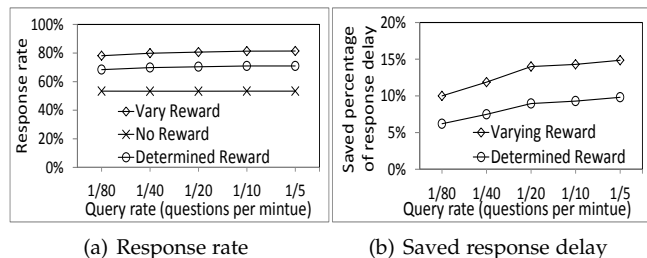


Fig. 9: Effectiveness of the reward strategy.

categorized by the asker. Figure 10(a) shows the precision rate of all algorithms versus the number of users in the system. It shows that the precision rate follows  $iASK > iASK\text{-Random} > Global(Tree)$ .  $iASK$  and  $iASK\text{-Random}$  have a better precision rate than  $Global(tree)$  due to the same reasons as in Figure 8(b). That is,  $iASK$  has the highest precision rate by choosing answerers with high QoS metric that considers precision rate.  $iASK$  uses interest coefficient to quickly find the weak tie, through which the answerers in nearby social communities can be quickly found to answer the non-factual questions with high quality. Without using the social networks, the Global systems choose global users that may have low precision rate for non-factual answers.  $iASK\text{-Random}$  randomly selects a friend to forward the question. However, because of the high cohesion of each interest community as shown in Figure 5(a),  $iASK\text{-Random}$  cannot always successfully find the weak tie towards the target community with the same interest as the question. Therefore,  $iASK$  generates a higher precision rate than  $iASK\text{-Random}$ . The figure indicates that the weak tie assisted social based potential answerer location algorithm is effective in supplying the highest quality answers to non-factual questions (whose interests are not in the askers' social communities) through social networks.

Figure 10(b) shows the recall rate of all algorithms versus the number of users in the system. It shows that both  $iASK$  and  $iASK\text{-Random}$  have a larger recall rate than  $Global(Tree)$  due to the same reasons as in Figure 8(c). It also shows that  $iASK$  generates a larger precision rate than  $iASK\text{-Random}$  due to the same reasons as in Figure 10(a). It indicates that the weak tie assisted social based potential answerer location algorithm is effective in identifying the weak tie to find the potential answerers through the social links. Figure 10(c) shows the F-Score of all algorithms, which indicates the overall performance of both precision rate and recall rate. It shows that the F-Score follows  $iASK > iASK\text{-Random} > Global(Tree)$  due to the same reasons as in Figures 10(a) and 10(b). It indicates that  $iASK$  achieves a better overall accuracy than all other methods in searching good potential answerers through social networks for non-factual questions.

We then measured the performance of interest coefficient based uncategorized question forwarding algorithm. In order to measure the performance of answering uncategorized questions, all questions outside of the asker's interests are uncategorized questions, and the questions are randomly selected from the questions of an interest having a low interest coefficient with the askers' interests. Figure 11(a) shows the precision rate of all algorithms versus the number of users in the system. It shows that  $iASK$  and  $iASK\text{-Random}$  have a larger precision rate than  $Global(Flat)$  due to the same reasons as in Figure 8(b).  $iASK$  uses the interest coefficient based uncategorized question forwarding algorithm to find the social community, which depends on the difference of interest closeness (as shown in Figure 5(b)) to find the friend having more different interests compared to the asker's interests. Due to the large user space, a random walk has a lower probability to reach the target community of the question. Therefore,  $iASK$  generates a larger precision rate than  $iASK\text{-Random}$ . It indicates that the interest coefficient based uncategorized question forwarding algorithm has the largest precision rate by successfully finding the social communities of the uncategorized questions to supply better answer quality.

Figure 11(b) shows the recall rate of all algorithms versus the number of users in the system. It shows that both  $iASK$  and  $iASK\text{-Random}$  have a larger recall rate than  $Global(Flat)$  due to the same reasons as in Figure 8(c). It also shows that  $iASK$  generates a larger precision rate than  $iASK\text{-Random}$  due to the same reasons as in Figure 11(a). It indicates that the interest coefficient based uncategorized question forwarding algorithm is effective in finding the social community to find the potential answerers through the social links. Figure 11(c) shows the F-Score of all algorithms, which indicates the overall performance of both precision and recall rates. The figure shows that the F-Score follows  $iASK > iASK\text{-Random} > Global(Flat)$  due to the same reasons as in both Figures 11(a) and 11(b). It indicates that  $iASK$  achieves a better overall accuracy than all other methods in searching good potential answerers for uncategorized questions.

#### 4.4 Performance of the Reward Strategy

We then measure the performance of our reputation-based reward strategy. In this experiment, the users only ask questions beyond their interests, and they randomly select a reward between [2,10] for each question. The virtual currency threshold for answering a question of each user is equal to his/her global reputation. We increased the question query rate (question/minute) for each user from 1/80 to 1/5, by doubling the rate at each step. Each experiment lasted for



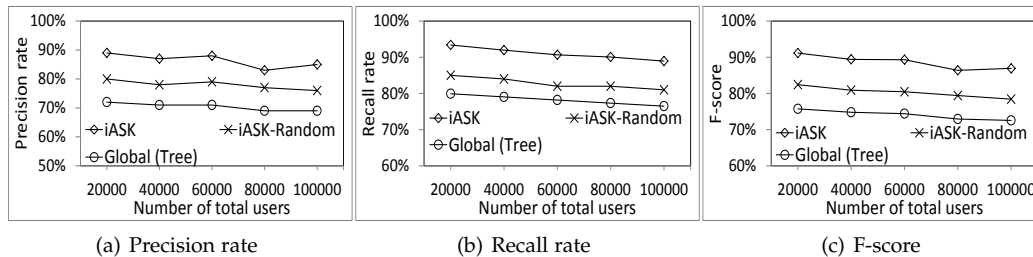


Fig. 10: Effectiveness of the weak tie assisted social based potential answerer location.

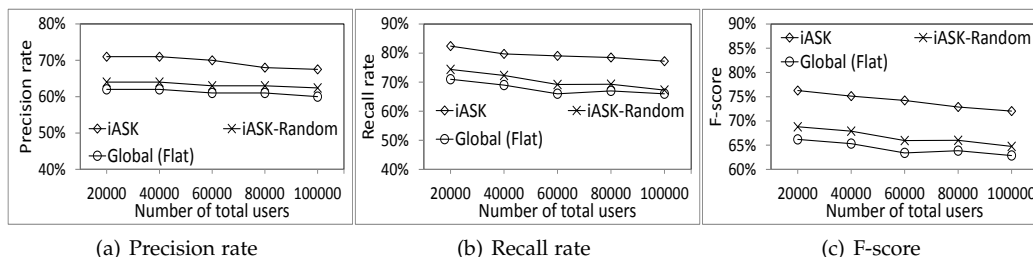


Fig. 11: Effectiveness of interest coefficient based uncategorized question forwarding.

1,000 minutes, and each user immediately responded to the question whenever his/her currency is not enough for his/her question asking.

Figure 9(a) shows the response rate of all methods as the question query rate increases, which follows *Varying reward* > *Determined reward* > *No reward*. Because in reward systems, users are motivated to quickly respond to a question whenever they are shortage of their currency, leading to a higher response rates than the fee-free system. In varying reward system, low reputed users need to answer more questions in order to ask a question with high reputation requirement. Thus, the *Varying reward* system has higher response rate than *Determined reward*. Figure 9(b) shows the additional saved percentage of the response delay of both reward systems compared to *No reward*. The saved percentage follows *Varying reward* > *Determined reward* as in Figure 9(a) due to the same reasons. Figures 9(a) and 9(b) indicate that iASK's reputation-based reward strategy is effective in motivating users to more cooperatively and quickly answer questions.

## 5 IASK IN THE REAL-WORLD TESTING

### 5.1 Real Implementation

We implemented iASK client in Java based on the Applet framework, and built a neural network for friend ranking. We also implemented the virtual server tree overlay in Java running on Tomcat 7.0 with MySQL database. Each virtual server was implemented as an independent thread. In order to avoid overloading a physical server, we ran each ten threads on a server in Palmetto [43], which has 771 8-core servers. The client can run in any browser supporting Java runtime environment 1.7. When asking or forwarding questions, each client selects  $K$  potential answerers to send a question independently according to iASK's algorithms. The screen shots for iASK are shown in Figure 13.

Figure 13(a) shows the main menu of iASK. Users can manage their profiles, ask and answer questions to help each other, manage personal friendship and contact-fan network, and rate the answers in order to update the weights of different factors for their QoS preference. Figures 13(b)

and 13(c) show the interfaces for asking and answering questions, respectively. Users choose interest categories of their questions. In this example, the user wants to ask a question in the "Research" category, which has subcategories including "Social network", "Cloud computing" and "Data mining". Each question will be forwarded to two users with the highest scores. Each potential answerer can answer, forward and drop each question. The TTL was set to 3. If a question cannot be resolved within TTL hops, it will be sent to the central servers. Based on the virtual server tree, all users with the interest of the question are located, and then two global potential answerers with the highest reputation values are selected to forward the question. We will present the experimental results from this real-world prototype for daily use in Section 5.

### 5.2 Performance Evaluation

We organized a testing with 42 students at our university. They built the social network according to their actual friendship between each other. In our experiment, the users are selected from different departments and they have different interests in the predefined category lists. We encourage each user to ask questions within or outside his/her interests. An asker needs to rate each answer with 0-10 stars, where 0 is totally unsatisfied, 5 is correct and 10 is very satisfied. In order to estimate the factors and weights, we first let users to ask five questions, rate all answers and follow others as fans. Then, we let users to ask another five questions for the measurement. We compared iASK with other four systems: i) iASK-R, which randomly selects two answerers; ii) iASK-L, which chooses the answerers with the lowest scores; iii) Global, which always sends questions to global experts and simulates Yahoo! Answers [3]; iv) Google, in which the asker gives the score for the first three answers from the Google search engine.

Figure 12(a) and Figure 12(b) show the rating scores of answers of factual questions and non-factual questions, respectively. The factual questions are like "What are the service models in Cloud computing?". These questions can be easily answered by an expert in this interest. The non-factual questions are like "How to learn data mining in

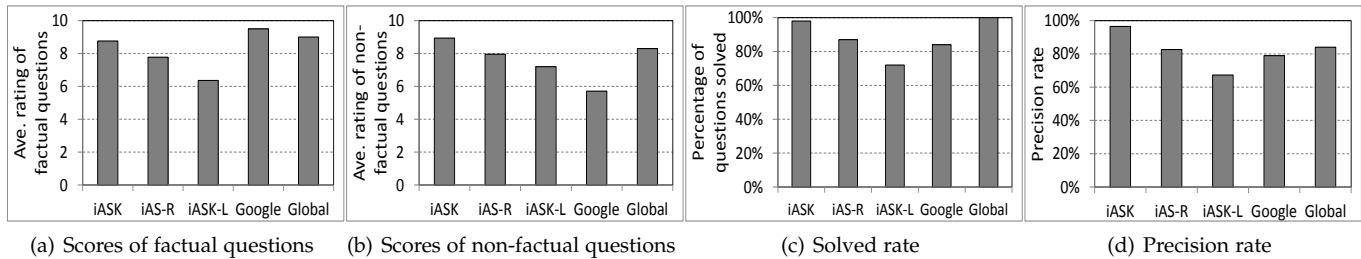


Fig. 12: Effectiveness of Q&A systems in the real-world testing.

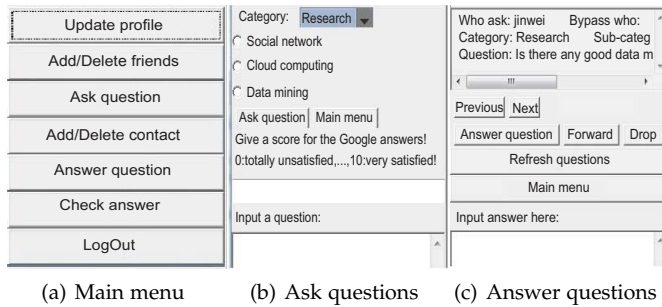


Fig. 13: Client software execution in a web browser.

our university?”. As shown in Figure 12(a), the scores of answers follows Google>Global>iASK>iASK-R>iASK-L. Google has the highest answer quality, because an expert among all users has limited knowledge compared to Google for factual questions. Global has a larger average score than all iASK methods because the expert is chosen from all users, who may have better knowledge in this interest. iASK chooses friends with better QoS scores, so it has a better performance than iASK-R. iASK-L has the worst performance because it always chooses friends with the lowest scores. The figure indicates the effectiveness of iASK’s social based answerer identification method to locate the expert, and the lower rating score of iASK than Google should be improved under a larger user scale with more friends to choose from.

Figure 12(b) shows the rating scores of answers of non-factual questions of all methods, which follows iASK>Global≈iASK-R>iASK-L>Google. Google has the lowest score without considering the askers’ preferences. iASK has better performance than iASK-R and iASK-L due to the same reasons as in Figure 12(a). iASK has the highest score because it always chooses answerers with high QoS values evaluated by its neural network friend ranking method that considers many factors. This figure indicates that iASK can supply the quality of best answers for non-factual questions.

Figure 12(c) shows the question solved rate of different methods, which is measured by the percentage of questions, each of which has at least one answer with rating no less than 5. It follows Global>iASK>iASK-R>Google>iASK-L. Global always chooses users with high reputations, and due to the small size of the users, the selected answerer may know the asker’s preferences. Thus, it generates 2% higher solved rate than iASK. iASK chooses friends with the best QoS scores, so it has a better performance than iASK-R. iASK-L has worse performance than iASK-R because it always chooses friends with the lowest scores. This figure indicates that iASK is effective in solving questions.

Figure 12(d) shows the precision rate of all different methods. The precision rate is measured by the percentage of answers, which have scores no less than 5. It shows

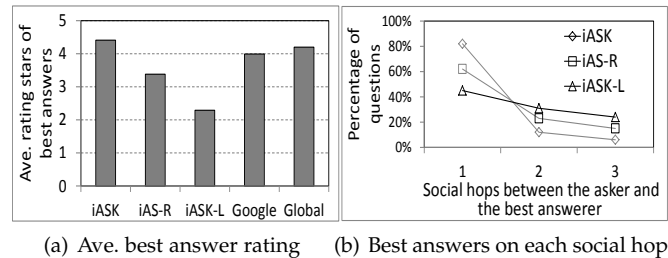


Fig. 14: The quality of best answers.

that the precision rate of all methods except Global follows iASK>iASK-R>Google>iASK-L due to the same reasons as in Figure 12(c). However, since Global cannot always supply correct answers for non-factual questions, it has a lower precision rate than iASK. Figures 12(c) and 12(d) together show that iASK solves more questions with better answer quality than other systems.

In our test, if a question does not have a best answer, the rating of its best answer was set to 0. We then measure the average star ratings of best answers as shown in Figure 14(a). It shows that the star ratings of all methods follows iASK>Global>Google>iASK-R>iASK-L due to the same reasons as in Figure 12(d), except that Google and Global have a better performance than iASK-R. That is due to the lower solved rate of iASK-R than Google and Global as shown in Figure 12(c). Figure 14(b) shows the percentage of best answers distribution over each social distance hop between the best answerer and asker. It shows that there are more best answers given by direct friends in iASK than in other two methods, due to the same reason as in Figure 12(d). Both Figures 14(a) and 14(b) indicate the effectiveness of iASK to select cooperative answerers in the social community intelligence.

## 6 CONCLUSION

In this paper, we propose iASK, a unified distributed Q&A system incorporating both social community intelligence and global collective intelligence. To find good answerer candidates in a user’s social network, iASK uses a neural network to consider multiple factors in evaluating the answer QoS of the user’s friends. If a question cannot be answered in a user’s social community, the answerer candidates will be located from the global user base. iASK builds central servers into a virtual server tree overlay to efficiently locate answerer candidates in the interest of the question. iASK has a fine-grained reputation system to locate cooperative global experts, and depends on a reputation-based reward strategy that adaptively rewards question answerers based on their reputations, in order to provide cooperative incentives in answering questions. iASK also has the

weak tie assisted social based potential answerer location algorithm and the interest coefficient based uncategorized question forwarding algorithm to further improve its performance. Our comprehensive trace-driven experiments and daily usage results from an iASK's prototype show that iASK outperforms other systems in enhancing answering QoS and efficiency. In our future work, the fault tolerance and the robust enhancement after the failure happen will be studied. We will test iASK on a larger user base in the real world and add more features to rank users in order to more precisely and efficiently locate the experts.

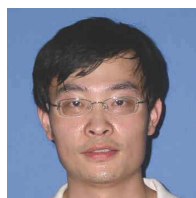
## ACKNOWLEDGEMENTS

This research was supported in part by U.S. NSF grants NSF-1404981, IIS-1354123, CNS-1254006, CNS-1249603, and Microsoft Research Faculty Fellowship 8300751. An early version of this work was presented in the Proc. of P2P [44].

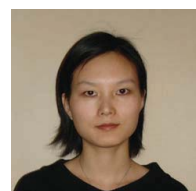
## REFERENCES

- [1] Ask, <http://www.ask.com>, [Accessed in May 2015].
- [2] Answers, <http://www.answers.com>, [Accessed in May 2015].
- [3] Yahoo! Answers, <http://answers.yahoo.com>, [Accessed in May 2015].
- [4] stackoverflow, <http://stackoverflow.com/>, [Accessed in May 2015].
- [5] Quora, <http://www.quora.com>, [Accessed in May 2015].
- [6] M. R. Morris, J. Teevan, and K. Panovich. What Do People Ask Their Social Networks, and Why?: A Survey Study of Status Message Q&A Behavior. In *Proc. of CHI*, 2010.
- [7] F. Harper, D. Raban, S. Rafaei, and J. Konstan. Predictors of Answer Quality in Online Q&A Sites. In *Proc. of SIGCHI*, 2008.
- [8] R. W. White, M. Richardson, and Y. Liu. Effects of Community Size and Contact Rate in Synchronous Social Q&A. In *Proc. of CHI*, 2010.
- [9] D. Horowitz and S.D. Kamvar. The Anatomy of a Large-Scale Social Search Engine. In *Proc. of WWW*, 2010.
- [10] Z. Li, H. Shen, G. Liu, and J. Li. SOS: A Distributed Mobile Q&A System Based on Social Networks. In *Proc. of ICDCS*, 2012.
- [11] Z. Li and H. Shen. Collective Intelligence in the Online Social Network of Yahoo!Answers and Its Implications. In *Proc. of CIKM*, 2012.
- [12] L. Zhang, X. Li, Y. Liu, Q. Huang, and S. Tang. Mechanism Design for Finding Experts Using Locally Constructed Social Referral Web. In *Proc. of INFOCOM*, 2012.
- [13] E. Bakshy, I. Rosenn, C. Marlow, and L. A. Adamic. The Role of Social Networks in Information Diffusion. *CoRR*, 2012.
- [14] M. Granovetter. *The Strength Of Weak Ties*. *American Journal of Sociology* 78, 1360-80, 1973.
- [15] J. Zhang, M. S. Ackerman, and L. Adamic. Expertise Networks in Online Communities: Structure and Algorithms. In *Proc. of WWW*, 2007.
- [16] J. Bian, Y. Liu, D. Zhou, E. Agichtein, and H. Zha. Learning to Recognize Reliable Users and Content in Social Media with Coupled Mutual Reinforcement. In *Proc. of WWW*, 2009.
- [17] M. Bouguessa, B. Dumoulin, and S. Wang. Identifying Authoritative Actors in Question-Answering Forums: the Case of Yahoo! Answers. In *Proc. of KDD*, 2008.
- [18] H. Zhang, T. N. Dinh, and M. T. Thai. Maximizing the Spread of Positive Influence in Online Social Network. In *Proc. of ICDCS*, 2013.
- [19] E. Amitay, D. Carmel, N. Har'El, and S. A. Ofek-Koiman. Golbandi: Social search and discovery using a unified approach. In *Proc. of HT*, 2009.
- [20] D. Carmel, N. Zwerdling, I. Guy, S. Ofek-Koifman, N. Har'el, I. Ronen, E. Uziel, S. Yogeve, and S. Chernov. Personalized social search based on user's social network. In *Proc. of CIKM*, 2009.
- [21] S. Kolay and A. Dasdan. The value of socially tagged URLs for a search engine. In *Proc. of WWW*, 2009.
- [22] S. Bao, G. Xue, X. Wu, Y. Yu, B. Fei, and Z. Su. Optimizing web search using social annotations. In *Proc. of WWW*, 2007.
- [23] B. M. Evans and E. H. Chi. Towards a Model of Understanding Social Search. In *Proc. of CSCW*, 2008.

- [24] J. Carretero, F. Isaila, A.-M. Kermarrec, F. Taian, and J. M. Tirado. Geology: Modular Georecommendation in Gossip-Based Social Networks. In *Proc. of ICDCS*, 2012.
- [25] G. Wang, K. Gill, M. Mohanlal, H. Zheng, and B. Y. Zhao. Wisdom in the social crowd: an analysis of quora. In *Proc. of WWW*, 2013.
- [26] X. Cheng and J. Liu. NetTube: Exploring Social Networks for Peer-to-Peer Short Video Sharing. In *Proc. of INFOCOM*, 2009.
- [27] Z. Wang, L. Sun, X. Chen, W. Zhu, J. Liu, M. Chen, and S. Yang. Propagation-based Social-Aware Replication for Social Video Contents. In *Proc. of ACM Multimedia*, 2012.
- [28] H. Shen, Z. Li, H. Wang, and J. Li. Leveraging Social Network Concepts for Efficient Peer-to-Peer Live Streaming Systems. In *Proc. of ACM Multimedia*, 2012.
- [29] X. Zhang and G. Cao. Efficient Data Forwarding in Mobile Social Networks with Diverse Connectivity Characteristics. In *Proc. of ICDCS*, 2014.
- [30] E. Pennisi. How did Cooperative Behavior Evolve? *Science*, 2005.
- [31] A. Mtibaa, M. May, C. Diot, and M. Ammar. Peoplerrank: Social Opportunistic Forwarding. In *Proc. of Infocom*, 2010.
- [32] G. Liu, H. Shen, and H. Chandler. Selective Data Replication for Online Social Networks with Distributed Datacenters. In *Proc. of ICNP*, 2013.
- [33] I. Stoica, R. Morris, D. Liben-Nowell, D. Karger, M. Kaashoek, F. Dabek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup protocol for internet applications. *IEEE/ACM Trans. Netw.*, 2003.
- [34] S. Haykin. *Neural Networks: A Comprehensive Foundataion. Second Edition*. Prentice-Hall Publisher, 1999.
- [35] P. GunWoo, Y. SoungWoung, L. Soojin, and L. SangHoon. Credible user identification using social network analysis in a q&a site. 2011.
- [36] M.S. Granovetter. The Strength of the Weak Tie: Revisited. *Sociological Theory*, 1093.
- [37] L. Mamykina, B. Manoim, M. Mittal, G. Hripcsak, and B. Hartmann. Design lessons from the fastest q&a site in the west. In *Proc. of SIGCHI*, 2011.
- [38] M. R. Morris, J. Teevan, and K. Panovich. What do people ask their social networks, and why?: a survey study of status message q&a behavior. In *Proc. of SIGCHI*, 2010.
- [39] B. Viswanath, A. Mislove, M. Cha, and K. P. Gummadi. On the evolution of user interaction in facebook. In *Proc. of WOSN*, 2009.
- [40] W. G. Stock. *Information Retrieval*. Butterworth, 2007.
- [41] Q. Lv, P. Cao, E. Cohen, K. Li, and S. Shenker. Search and replication in unstructured peer-to-peer networks. In *Proc. of ISC*, 2002.
- [42] R. Gray, N. B. Ellison, J. Vitak, and C. Lampe. Who wants to know?: question-asking and answering practices among facebook users. In *Proc. of CSCW*, 2013.
- [43] Palmetto Cluster. <http://citi.clemson.edu/palmetto/>, [Accessed in May 2015].
- [44] G. Liu and H. Shen. iASK: A Distributed Q&A System Incorporating Social Community and Global Collective Intelligence. In *Proc. of the IEEE International Conference on Peer-to-Peer Computing (P2P)*, 2015.



**Guoxin Liu** received the BS degree in BeiHang University 2006, and the MS degree in Institute of Software, Chinese Academy of Sciences 2009. He is currently a Ph.D. student in the Department of Electrical and Computer Engineering of Clemson University. His research interests include distributed networks, with an emphasis on Peer-to-Peer, datacenter and online social networks.



**Haiying Shen** received the BS degree in Computer Science and Engineering from Tongji University, China in 2000, and the MS and Ph.D. degrees in Computer Engineering from Wayne State University in 2004 and 2006, respectively. She is currently an Associate Professor in the department of computer science at University of Virginia. Her research interests include distributed computer systems and computer networks with an emphasis on P2P and content delivery networks, mobile computing, wireless sensor networks, and grid and cloud computing. She was the Program Co-Chair for a number of international conferences and member of the Program Committees of many leading conferences. She is a Microsoft Faculty Fellow of 2010, a senior member of the IEEE and a member of the ACM.