

# 1 The Impacts of Structural Difference and Temporality of Tweets on 2 Retrieval Effectiveness

3 LIFENG JIA and CLEMENT YU, University of Illinois at Chicago  
4 WEIYI MENG, Binghamton University

5 To explore the information seeking behaviors in microblogosphere, the microblog track at TREC 2011 intro-  
6 duced a real-time ad-hoc retrieval task that aims at ranking relevant tweets in reverse-chronological order.  
7 We study this problem via a two-phase approach: 1) retrieving tweets in an ad-hoc way; 2) utilizing the  
8 temporal information of tweets to enhance the retrieval effectiveness of tweets. Tweets can be categorized  
9 into two types. One type consists of short messages not containing any URL of a Web page. The other type  
10 has at least one URL of a Web page in addition to a short message. These two types of tweets have dif-  
11 ferent structures. In the first phase, to address the structural difference of tweets, we propose a method to  
12 rank tweets using the divide-and-conquer strategy. Specifically, we first rank the two types of tweets sep-  
13 arately. This produces two rankings, one for each type. Then we merge these two rankings of tweets into  
14 one ranking. In the second phase, we first categorize queries into several types by exploring the temporal  
15 distributions of their top-retrieved tweets from the first phase; then we calculate the time-related relevance  
16 scores of tweets according to the classified types of queries; finally we combine the time scores with the IR  
17 scores from the first phase to produce a ranking of tweets. Experimental results achieved by using the TREC  
18 2011 and TREC 2012 queries over the TREC Tweets2011 collection show that: (i) our way of ranking the two  
19 types of tweets separately and then merging them together yields better retrieval effectiveness than rank-  
20 ing them simultaneously; (ii) our way of incorporating temporal information into the retrieval process yields  
21 further improvements, and (iii) our method compares favorably with state-of-the-art methods in retrieval  
22 effectiveness.

23 Categories and Subject Descriptors: H.3.3 [Information Storage and Retrieval]: Information Search and  
24 Retrieval; H.4.m [Information Systems Applications]: Miscellaneous

25 General Terms: Experimentation, Performance

26 Additional Key Words and Phrases: Ad-hoc retrieval of tweets, learning to rank, query temporal  
27 categorization

## 28 ACM Reference Format:

29 Jia, L., Yu, C., and Meng, W. 2013. The impacts of structural difference and temporality of tweets on retrieval  
30 effectiveness. *ACM Trans. Inf. Syst.* 31, 4, Article 21 (November 2013), 38 pages.  
31 DOI: <http://dx.doi.org/10.1145/2500751>

## 32 1. INTRODUCTION

33 Twitter, a worldwide popular microblog service, has a daily volume of over 340 million  
34 tweets,<sup>1</sup> which motivates research interests in studying the information seeking  
35 behaviors within microblogosphere. The microblog track at TREC 2011 introduced  
36 a real-time ad-hoc retrieval task, whereby a user wishes to see the most recent and

<sup>1</sup><http://en.wikipedia.org/wiki/Twitter>

---

Authors' addresses: L. Jia and C. Yu, Computer Science Department, University of Illinois at Chicago, IL;  
email: [ljia2@uic.edu](mailto:ljia2@uic.edu); W. Meng, Computer Science Department, Binghamton University.

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

© 2013 ACM 1046-8188/2013/11-ART21 \$15.00

DOI: <http://dx.doi.org/10.1145/2500751>

37 relevant information to a query within Twitter [Ounis et al. 2011]. To respond to  
 38 a query with a timestamp  $t$ , the retrieved tweets should satisfy the following three  
 39 conditions: (1) relevant to the query, (2) published on or before time  $t$ , and (3) ranked  
 40 in reverse-chronological order of their publishing times.

41 Some studies have been done in information retrieval of tweets. These studies can be  
 42 categorized into two major classes. The techniques in the first class [Choi et al. 2012;  
 43 Duan et al. 2010; Han et al. 2012; Metzler and Cai 2011; Zhang et al. 2012] rank tweets  
 44 by measuring the lexical similarities between tweets and queries. The methods in the  
 45 second class [Amati et al. 2012; Choi and Croft 2012; Dong et al. 2010b; Efron and  
 46 Golovchinsky 2011] rank tweets by exploring temporal information (the publishing  
 47 times of tweets and the timestamps of queries). Some studies [Efron et al. 2012; Liang  
 48 et al. 2012; Massoudi et al. 2011] employ both lexical similarity and temporality in  
 49 ranking tweets. However, there are two important issues that are not well addressed  
 50 by these existing works.

51 The first issue is the impact of the structural difference of tweets on retrieval ef-  
 52 fectiveness. Specifically, there are two types of tweets that have different structures.  
 53 The first type (to be defined as T-tweet in Section 3.2) is just a short text message  
 54 with no more than 140 characters. The second type (to be defined as TU-tweet in  
 55 Section 3.2) contains at least one URL of a Web page in addition to a short text mes-  
 56 sage. All existing studies simultaneously rank both types of tweets. However, we be-  
 57 lieve it is important to utilize the structural difference of tweets in retrieval. Let us  
 58 illustrate the motivation by the following example.

59 *Example 1.* Consider a query  $q = \text{“phone hacking British politicians”}$ , a tweet  
 60  $d_1 = \text{“@jamesrae andy Gray is suing the NOTW... just got fired from Sky for footage$   
 61  $\text{that should never have been seen. I smell Murdoch!”}$ , a second tweet  $d_2 = \text{“Ten-}$   
 62  $\text{sions simmer as ‘frustrated’ Rupert Murdoch flies in to face phone-hacking affair}$   
 63  $\text{http://t.co/b3kOppY via @guardian”}$  and a third tweet  $d_3 = \text{“Windows Phone 7 gets}$   
 64  $\text{USB Tethering Hack http://tinyurl.com/4lafss6”}$ .  $d_1$  is a T-tweet that only has a short  
 65 message.  $d_1$  is relevant to  $q$  but has no query terms.  $d_2$  and  $d_3$  are two TU-tweets.  
 66 Each of them has not only a message but also a URL.  $d_2$  is relevant to  $q$ . It contains  
 67 two query terms “*phone*” and “*hacking*” in its message and all four query terms in the  
 68 web page of the URL in  $d_2$ .  $d_3$  is irrelevant to  $q$ . It contains two query terms “*Phone*”  
 69 and “*hack*” in its message. The Web page of the URL in  $d_3$  has no query terms. The  
 70 content of a TU-tweet is the union of its short message and the contents of the Web  
 71 pages of the URLs in it. It is intuitive that for a TU-tweet, the higher the percentage  
 72 of query terms appearing in it is, the more likely the tweet is relevant. The relevant  $d_2$   
 73 has more query terms than the irrelevant  $d_3$ . However, such an intuition does not ap-  
 74 ply for a T-tweet.  $d_1$  has no query terms but it is relevant to  $q$ . This is because T-tweets  
 75 are so short that some relevant T-tweets may not have any query terms. In addition,  
 76 we find out that (see Section 6.1.2) the sets of the most important features for learning  
 77 to rank the two types of tweets are very different.

78 Motivated by such an observation, we propose to use the divide-and-conquer strat-  
 79 egy to address the structural difference of tweets. Specifically, we learn two rankers  
 80 that are dedicated to ranking T-tweets and TU-tweets separately. This produces two  
 81 tweet type-specific rankers. We then learn a classifier that determines a preference be-  
 82 tween any T-tweet and any TU-tweet with respect to a given query. The details about  
 83 these two tweet type-specific rankers and the classifier are discussed in Sections 3.2  
 84 and 3.3, respectively. Given a query  $q$ , we first obtain a ranking of T-tweets,  $R_1$ , and  
 85 a ranking of TU-tweets,  $R_2$ , by using the two type-specific rankers, respectively. Then  
 86 we apply the classifier to determine the preference between each T-tweet from  $R_1$  and  
 87 each TU-tweet from  $R_2$ . Finally, we merge the tweets from  $R_1$  and  $R_2$  into a single

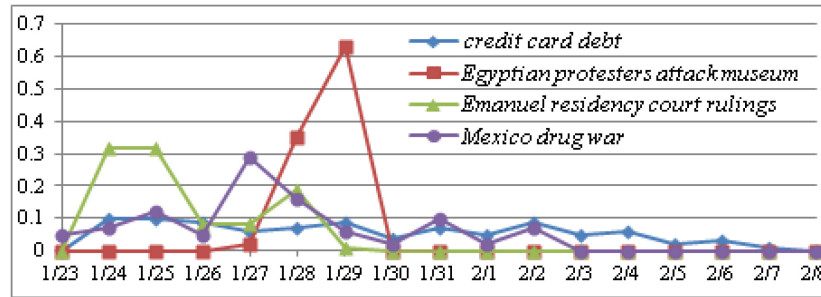


Fig. 1. The distributions of relevant tweets over time.

88 ranking. The merging process considers the preferences from the two rankers and  
 89 the classifier. The discussion of how to merge two rankings of tweets is presented in  
 90 Section 3.4.

91 The second issue is the impact of the temporal sensitivities of queries on the re-  
 92 trieval effectiveness of tweets. Queries can be categorized into time sensitive and time  
 93 insensitive types [Dakka et al. 2012; Jones and Diaz 2007]. For ease of presentation,  
 94 Figure 1 shows the temporal distributions of the relevant tweets with respect to four  
 95 sample TREC queries. These distributions are plotted over the period from 1/23/2011  
 96 to 2/8/2011, when the TREC Tweets2011 collection was sampled from Twitter. The  
 97  $x$ -axis represents time in the unit of day [Efron and Golovchinsky 2011]. The  $y$ -axis  
 98 represents the percentage of relevant tweets published on a particular day. By observ-  
 99 ing these distributions, we claim that there are three types of queries. The first type is  
 100 insensitive to time, while the last two types are time sensitive.

101 — The first type of queries has a relatively flat (uniform) distribution of relevant tweets  
 102 over time, indicating that these queries are insensitive to time. This is exemplified  
 103 by the query “*credit card debt*.”

104 — The second type of queries has a dominant peak in terms of their temporal  
 105 distributions of relevant tweets. The dominant peak contains an extremely large  
 106 portion of relevant tweets concentrating on a single day. This is exemplified by the  
 107 query “*Egyptian protesters attack museum*.” The attack happened during the night  
 108 of 1/28/2011 and a dominant peak in the distribution is formed on 1/29/2011. An  
 109 event related to the topic of such a query is usually the event of a breaking news  
 110 story. The relevant tweets are so concentrated around the peak that the percentage  
 111 of relevant tweets rapidly decreases beyond the peak. In this article, such queries  
 112 are called *dominant peak queries*.

113 — The third type of queries has one or more nondominant peaks. Each peak contains a  
 114 significant portion of relevant tweets on a day but the percentage of relevant tweets  
 115 of a nondominant peak is not as high as that of a dominant peak. A nondominant  
 116 peak of a query is caused by an event that is related to the query. These related  
 117 events trigger people’s intensive discussions about the query topic at different  
 118 times. This is exemplified by two queries: “*Mexico drug war*” and “*Emanuel resi-  
 119 dency court rulings*.” For “*Mexico drug war*,” the nondominant peak on 1/27/2011 is  
 120 caused by a related event, “*Pot-firing catapult found at Arizona-Mexico border*”. For  
 121 “*Emanuel residency court rulings*,” the first two nondominant peaks on 1/24/2011  
 122 and 1/25/2011 correspond to the event: “*Illinois Court Throws Emanuel Off Chicago  
 123 Mayoral Ballot*”; the third peak on 1/28/2011 corresponds to another related event:  
 124 “*Illinois Supreme Court keeps Emanuel on ballot*.” In this article, such queries are  
 125 called *nondominant peak queries*.

126 These three types of queries depend on their temporal distributions of relevant  
 127 tweets. In practice, it is unrealistic to know such distributions for given queries. In  
 128 Efron and Golovchinsky [2011] and Jones and Diaz [2007], the temporal distribution  
 129 of the relevant tweets with respect to a query  $q$  can be approximated by that of the  
 130 top tweets with respect to  $q$ . These top tweets can be retrieved by a ranking model,  
 131 such as BM25 [Robertson et al. 1996]. In this article, we classify queries into different  
 132 types by the temporal distributions of their top tweets. For time-insensitive queries,  
 133 there is no need to employ temporal information; for time sensitive queries, we propose  
 134 two different techniques to calculate the temporal relevance of tweets to dominant  
 135 peak queries and to nondominant peak queries, respectively. The degree of temporal  
 136 relevance is measured by a time-related relevance score (to be given in Section 4.2).  
 137 Our proposed method for categorizing queries and for computing the time-related rel-  
 138 evance scores with respect to the two types of time sensitive queries are presented in  
 139 Section 4.2. In this article, we only study these three types of queries. The studies of  
 140 other types of queries, such as cyclic queries (e.g., “*Halloween*”) are deferred to future  
 141 work.

142 Our work has two novelties: 1) ranking the two types of tweets by a divide-and-  
 143 conquer manner can improve retrieval effectiveness; and 2) our temporal classification  
 144 of queries and two different ways of computing the time-related relevance scores with  
 145 respect to the two different types of time sensitive queries are different from existing  
 146 works. We now summarize the research questions we aim to answer in this article.

147 — Acknowledging that tweets can be classified into the two types by their different  
 148 structures, is the retrieval effectiveness of tweets affected by their structural  
 149 difference?

150 — How to leverage the structural difference of tweets to enhance their retrieval  
 151 effectiveness?

152 — What are the effectiveness and the efficiency of the proposed algorithm?

153 — How can we improve retrieval effectiveness by taking into consideration the  
 154 temporal information (publishing times) of tweets?

155 — How does our method perform compared to various state-of-the-art methods?

156 This article has the following contributions.

157 — We investigate the impact of the structural difference of tweets on retrieval  
 158 effectiveness.

159 — We present a novel algorithm of ranking tweets by using the divide-and-conquer  
 160 strategy. To our knowledge, our work is the first study that leverages the structural  
 161 difference of tweets to enhance retrieval effectiveness.

162 — We present a novel categorization of queries by their sensitivities to time.

163 — We propose different techniques to calculate the degrees of temporal relevance of  
 164 tweets with respect to the different categories of queries.

165 The remainder of this article is organized as follows. We review the related works  
 166 in Section 2. Section 3 introduces our divide-and-conquer method for ranking tweets.  
 167 Section 4 discusses our method for categorizing queries in terms of their temporal  
 168 sensitivities and proposes different techniques to calculate the temporal relevance of  
 169 tweets. Experimental setup and experimental results are provided in Section 5 and  
 170 Section 6, respectively. The article is concluded in Section 7.

## 171 2. RELATED WORK

172 Recently, interests are rising in exploring Twitter for information retrieval of tweets  
 173 by different criteria, such as lexical relevance [Choi et al. 2012; Duan et al. 2010; Han  
 174 et al. 2012; Metzler and Cai 2011; Zhang et al. 2012], temporal relevance [Amati et al.

175 2012; Choi and Croft 2012; Dong et al. 2010b; Efron and Golovchinsky 2011] and  
176 jointly lexical and temporal relevance [Efron et al. 2012; Liang et al. 2012; Massoudi  
177 et al. 2011]. Beyond tweet retrieval, some studies [Amodeo et al. 2011; Dakka et al.  
178 2012; Dong et al. 2010a; Jones and Diaz 2007; Keikha et al. 2011a, 2011b; Li and  
179 Croft 2003] also showed that incorporating the publishing times of documents into  
180 the retrieval process is beneficial for ad-hoc retrieval. Instead of using the publishing  
181 times of documents, some works [Berberich et al. 2010; Dai and Davison 2010; Elsas  
182 and Dumais 2010; Kulkarni et al. 2011] studied how to improve the ranking effec-  
183 tiveness by using the temporal information extracted from the contents of documents.  
184 Moreover, our study is also related to some works [Ailon et al. 2008; Bian et al. 2010;  
185 Dai et al. 2011; Hüllermeier and Fürnkranz 2010] in learning to rank. In the rest of  
186 this section, we review in greater detail the related works.

### 187 **2.1. Lexical Relevance-Based Retrieval**

188 The first thread of related works studied tweet retrieval by measuring their lexical  
189 similarities to queries. Duan et al. [2010] employed RankSVM [Herbrich et al. 2000;  
190 Joachims 2002] to rank tweets by their lexical relevance to queries. Metzler and Cai  
191 [2011] studied the real-time ad-hoc tweet retrieval problem by using RankSVM to  
192 rank tweets with respect to queries and rearranged the top-ranked tweets in reverse-  
193 chronological order. This work achieved the best results reported in TREC 2011. Choi  
194 et al. [2012] showed that the quality of tweets is correlated with their relevance and  
195 applied the quality features in relevance ranking. They assumed that high quality  
196 tweets are more likely to be retweeted than low quality ones and learned a model to  
197 estimate the probability of a tweet being retweeted by exploring its lexical content.  
198 Zhang et al. [2012] proposed a query-specific model to rank tweets by considering the  
199 characteristics unique to a query. Specifically, given a query  $q$ , they treated the top and  
200 the bottom tweets retrieved by a ranking model as positive and negative examples and  
201 then learned a ranking model specific to  $q$ . Efron et al. [2012] expanded each tweet  
202  $d$  with respect to a query  $q$  as follows. The terms of the most similar tweets to  $d$  are  
203 added to  $d$ . The query  $q$  is then compared with the expanded tweets for the similarity  
204 computation, in order to improve retrieval effectiveness. Han et al. [2012] expanded  
205 each tweet  $d$  in a similar manner by the terms from other tweets that are lexically  
206 similar to  $d$ . Our work has two fundamental differences from the works reviewed ear-  
207 lier: 1) we consider the structural difference of the two types of tweets in the retrieval  
208 process while they ranked both types of tweets together; and 2) they only measured the  
209 lexical similarities of tweets to queries while we take into consideration both lexical  
210 similarities and temporal information.

### 211 **2.2. Temporal Relevance-Based Retrieval**

212 The second thread of related works studied the impact of temporal information on re-  
213 trieval effectiveness. Dong et al. [2010a, 2010b] proposed the recency ranking problem  
214 and studied the problem using Twitter data. Amati et al. [2012] assumed that the re-  
215 cent tweets with respect to (the timestamp of) a query  $q$  are more likely to be relevant  
216 than the old tweets. Massoudi et al. [2011] studied a query expansion method where  
217 the expanded query terms are selected from high-quality and recent tweets, instead of  
218 low-quality and old tweets. The quality of tweets can be estimated by some indicators,  
219 such as the number of followers of Twitter users. All the works we have mentioned in  
220 principle prefer recent tweets (or terms from recent tweets) to old ones. However, this  
221 is not always desirable. For example, in Figure 1, for the query “*Mexico drug war,*” a  
222 significant portion of relevant tweets are published on 1/27/2011 and some relevant  
223 tweets are published on 2/2/2011. The tweets on 1/27/2011 are as relevant as those  
224 tweets on 2/2/2011. They should not be assigned lower priorities in retrieval. Our work

225 classifies queries by the temporal distributions of their top tweets and then proposes  
226 different ways of utilizing temporal information of tweets according to the classified  
227 types of queries. Liang et al. [2012] studied the real-time ad-hoc tweet retrieval by a  
228 two-phase approach where 1) an ad-hoc retrieval of tweets is conducted and 2) tweets  
229 are re-ranked to promote the relevant and recent ones. Our two-phase method is differ-  
230 ent from theirs in two aspects. First, they ranked both types of tweets simultaneously  
231 while we leverage the structural difference of tweets. Second, they promoted recent  
232 tweets over old tweets while we classify queries by their time sensitivities before ap-  
233 plying temporal information in different manners according to the classified types of  
234 queries. Choi and Croft [2012] obtained the top tweets (consisting of retweets and non-  
235 retweets) with respect to a query  $q$  from a ranking model. Then they explored the  
236 temporal distribution of the top retweets to measure the importance of each day with  
237 respect to  $q$ . The importance of a day  $t$  to  $q$  is proportional to the number of the top  
238 retweets published on  $t$ . Finally, they arranged non-retweets by considering the im-  
239 portance of each of their publishing days. Our work differs from theirs in that they  
240 use retweets to measure the importance of days while we use top tweets to determine  
241 the importance of days. Moreover, our calculation of the degrees of relevance between  
242 tweets and queries by temporality is quite different from theirs. Efron et al. [2012]  
243 obtained the top tweets with respect to a query  $q$  and then, for each tweet  $d$ , acquired  
244 the most similar (top) tweets to  $d$ . They calculated the temporal similarity between  $q$   
245 and  $d$  based on the temporal distribution of  $q$ 's top tweets and that of  $d$ 's top tweets.  
246 Our work differs from their work in that we classify queries based on the temporal  
247 distributions of their top tweets and then calculate the temporal relevance of tweets to  
248 queries by their classified types.

249 Besides Twitter search, Li and Croft [2003] studied time sensitive queries and  
250 assumed that relevant documents are mostly recent documents. They proposed an  
251 exponential-based age penalty strategy where aged documents are penalized and then  
252 demoted to boost the ranking positions of recent documents. Efron and Golovchinsky  
253 [2011] studied the same problem and proposed a query-specific exponential-based  
254 age penalty method where aged documents are penalized differently with respect to  
255 different queries. Our classification, determination and handling of time sensitive  
256 queries are different from the given works. Moreover, their hypothesis [Efron and  
257 Golovchinsky 2011; Li and Croft 2003] that aged documents should be penalized more  
258 than recent documents is not necessarily true for some time sensitive queries. For ex-  
259 ample, in Figure 1, for the query “*Mexico drug war*”, the relevant tweets on 1/27/2011  
260 should not be penalized relative to those on 2/2/2011. Amodeo et al. [2011] and Keikha  
261 et al. [2011b] presented temporal query expansions by using the terms selected from  
262 the top (blog) documents (with respect to a query  $q$ ) that are published on the days that  
263 are most relevant to  $q$ . The relevance of a day  $t$  to  $q$  is measured by the average similar-  
264 ity of the top documents published on  $t$  to  $q$  in Keikha et al. [2011b] or by the percent-  
265 age of  $q$ 's top documents published on  $t$  [Amodeo et al. 2011]. We do not use temporal  
266 information in query expansion. Keikha et al. [2011a] showed that blog feed retrieval  
267 can benefit from the usage of temporal information. They studied the retrieval of  
268 blog feeds. A blog feed consists of a set of blog documents published on different  
269 days. We study the retrieval of individual tweets. Although both studies use temporal  
270 information, the utilizations of temporality in these two studies are very different.  
271 Dakka et al. [2012] indicated that, for a time sensitive query  $q$ , a document  $d$  can be  
272 represented by two dimensions: the lexical content  $c_d$  and the publishing time  $t_d$ . They  
273 assumed the independence between  $c_d$  and  $t_d$ . Our work differs from theirs in that we  
274 assume the contents of documents (tweets) and their publishing times are not neces-  
275 sarily independent. For example, for the query “*Emanuel residency court rulings*,” the  
276 relevant tweets published on 1/24/2011 and 1/25/2011 discuss the event “*Illinois Court*

277 *Throws Emanuel Off Chicago Mayoral Ballot*” while the relevant tweets published on  
 278 1/28/2011 discuss the event “*Illinois Supreme Court keeps Emanuel on ballot.*” The con-  
 279 tents of tweets with respect to a query can be influenced by other related events which  
 280 happen at different times. Jones and Diaz [2007] categorized queries into time insensi-  
 281 tive ones, temporally ambiguous queries such as “*Iraq War*” (referencing two different  
 282 wars) and temporally unambiguous queries such as “*Turkish earthquake 1999*”. Our  
 283 work categorizes queries by their sensitivities to time, instead of their temporal  
 284 ambiguities.

285 Exploring the temporal information from the contents of documents can improve  
 286 retrieval effectiveness too. Berberich et al. [2010] proposed a language model supple-  
 287 mented with a temporal dimension where the temporal information from a query and  
 288 that from documents are uniformly expressed and matched in retrieval. For exam-  
 289 ple, the query, “*World Cups in 1990s*” should be matched by the documents containing  
 290 “*1998 World Cup,*” because “*1990s*” temporally covers “*1998.*” Elsas and Dumais [2010]  
 291 studied the relationship between the temporal dynamics of document contents and  
 292 the relevance of documents. For example, they showed that the contents of the rele-  
 293 vant documents for navigational queries, such as “*YouTube,*” have great and frequent  
 294 changes over time. Kulkarni et al. [2011] discussed the interaction among the tempo-  
 295 ral changes of query popularity, the temporal changes of document contents and query  
 296 intents. Dai and Davison [2010] utilized the freshness of Web site contents for comput-  
 297 ing Web site authority by examining the frequency of Web site content changes and  
 298 that of Web site hyperlink changes over time. Our work uses the publishing times of  
 299 top documents (tweets) to improve retrieval effectiveness.

### 300 2.3. Learning to Rank

301 Our work is also related to some studies in learning to rank. Bian et al. [2010] provided  
 302 a divide-and-conquer framework for learning to rank documents. Dai et al. [2011] ex-  
 303 tended the same divide-and-conquer framework for learning to rank documents by  
 304 freshness and relevance simultaneously. Our work has a fundamental difference from  
 305 theirs. Both works [Bian et al. 2010; Dai et al. 2011] divided (clustered) queries into  
 306 different clusters where queries within a cluster have a similar set of important learn-  
 307 ing to rank features. However, we divide (partition) documents (tweets) into two sets  
 308 by considering their structural difference. Given some different rankings of a same  
 309 set of documents that yield inconsistencies, Ailon et al. [2008] studied how to obtain a  
 310 ranking of the same set of documents that approximately minimizes the disagreement  
 311 with the given rankings. In our work, we merge two rankings of two different sets  
 312 of tweets, one for T-tweets and the other for TU-tweets. Hüllermeier and Fürnkranz  
 313 [2010] studied the problem where each example (document) is assigned the probabili-  
 314 ties of belonging to different classes. No ranking of examples (documents) is discussed  
 315 in Hüllermeier and Fürnkranz [2010].

## 316 3. A DIVIDE-AND-CONQUER METHOD FOR RANKING TWEETS

317 In this section, we introduce a novel method for ranking tweets. This method explores  
 318 the structural difference of tweets by the divide-and-conquer strategy. It is deployed as  
 319 the first phase to produce a ranking of tweets, taking into consideration their lexical  
 320 similarities to queries only.

### 321 3.1. Method Overview

322 In this method, we differentiate the following two types of tweets: the first type is a  
 323 short plain message without URLs (T-tweet) and the second type is a message con-  
 324 taining at least one URL (TU-tweet). A URL usually leads to a Web page with a sub-  
 325 stantially more content than a short message. To explore such a structural difference,

326 we propose to rank these two types of tweets separately and then merge the two type-  
 327 specific rankings of tweets into a single ranking. The proposed method has two tweet  
 328 type-specific rankers and a classifier. The two type-specific rankers are dedicated to  
 329 ranking T-tweets and TU-tweets. The classifier calculates the preference between any  
 330 T-tweet and any TU-tweet with respect to a query.

331 In this article, we resort to the learning to rank algorithms to produce the two  
 332 rankers. Specifically, RankSVM [Herbrich et al. 2000; Joachims 2002] is employed.  
 333 It can consider not only various lexical similarities between queries and tweets, such  
 334 as BM25 similarity [Robertson et al. 1996], but also some special social network char-  
 335 acteristics that are independent of queries, such as the number of retweets of tweets.  
 336 It leverages different criteria as features to learn the two type-specific rankers. We  
 337 denote as *T-tweet Ranker* the RankSVM model that is dedicated to ranking T-tweets.  
 338 It is learned over the training data consisting of a set of training queries  $Q$  and a set of  
 339 labeled T-tweets with respect to  $Q$ . Let *TU-tweet Ranker* denote the RankSVM model  
 340 that is TU-tweet oriented. It is learned over the training data consisting of the same  
 341 set of training queries  $Q$  but a different set of labeled TU-tweets with respect to  $Q$ .  
 342 A classifier is learned to determine a preference between each T-tweet and each TU-  
 343 tweet. Specifically, it is learned by using the union of the two sets of labeled tweets  
 344 with respect to the same training query set  $Q$ . The classifier indicates for each T-tweet  
 345  $d_1$  and each TU-tweet  $d_2$  whether  $d_1$  is preferred over  $d_2$  or vice versa.

346 The goal of this method is to produce the ranking of tweets for a set of test queries,  
 347  $Q' = \{q'_1, q'_2, \dots, q'_m\}$ . For each test query  $q'_i$ , we apply the *T-tweet Ranker* to obtain a  
 348 ranking of T-tweets  $R_1$ . Then we obtain a ranking of TU-tweets  $R_2$  by the *TU-tweet*  
 349 *Ranker*. For each pair of one T-tweet from  $R_1$  and one TU-tweet from  $R_2$ , the classi-  
 350 fier is employed to determine a preference relationship between them with respect to  
 351  $q'_i$ . There are three sets of preferences: 1) the preference between any two T-tweets  
 352 which is indicated by their relative ranking positions in  $R_1$ ; 2) the preference of any  
 353 two TU-tweets from  $R_2$ ; and 3) the preference between any T-tweet from  $R_1$  and any  
 354 TU-tweet from  $R_2$  indicated by the classifier. Finally, the two rankings,  $R_1$  and  $R_2$ , are  
 355 merged into a ranking by considering all three sets of preferences.

356 Because these three sets of preferences are computed by three different models,  
 357 there may be inconsistent preferences. For example, given two T-tweets  $d_i$  and  $d_j$  and a  
 358 TU-tweet  $d_k$ , the *T-tweet Ranker* may indicate  $d_i > d_j$ , which denotes the preference of  
 359  $d_i$  over  $d_j$ . However, the classifier may indicate  $d_k > d_i$  and  $d_j > d_k$ . In such a circular  
 360 preference situation, no matter how these three tweets are ranked in the merged rank-  
 361 ing, there is at least one inconsistency. Suppose that the degree of the preference of  $d_i$   
 362 over  $d_j$  is 0.5, that of  $d_j$  over  $d_k$  is 0.4, that of  $d_k$  over  $d_i$  is 0.3, and there are no other  
 363 preferences. If we determine that  $d_i$  is ranked above  $d_j$  which is ranked above  $d_k$ , it  
 364 will incur an inconsistency with the degree of 0.3. This is the smallest amount of incon-  
 365 sistency among all possible orderings of these three tweets. In an ideal situation, we  
 366 want to merge the two type-specific rankings into an optimal ranking that agrees best  
 367 with the three sets of preferences. However, such a problem is NP-complete [Cohen  
 368 et al. 1998]. Therefore, we propose a greedy merging algorithm called *GreedyMerging*.  
 369 This algorithm always picks the tweet to be ahead of the remaining tweets, if it incurs  
 370 the least amount of inconsistency relative to any of the remaining tweets. If there is  
 371 no inconsistency among the three sets of preferences, the algorithm will produce the  
 372 optimal merged ranking consistent with all preferences.

### 373 3.2. Tweet Type-Specific Rankers

374 In this section, we present the two rankers: one ranks T-tweets while the other ranks  
 375 TU-tweets. For ease of introduction, we first define T-tweets and TU-tweets.



376 *Definition 3.1 (T-Tweet).* A T-tweet is a tweet whose message body has no URLs.  
 377 The structure of a T-tweet consists of only one field:

378 a) Tweet Message Field: the message body of the tweet.

379 *Definition 3.2 (TU-Tweet).* A TU-tweet is a tweet whose message body has at least  
 380 one URL. A tweet whose message body has URLs only is very rare. The structure of a  
 381 TU-tweet consists of three fields:

382 a) Tweet Message Field: the message body with the exclusion of the embedded URLs.

383 b) URL Title Field: the union of the titles of the Web pages of the embedded URLs.

384 c) URL Body Field: the union of the bodies of the Web pages of the embedded URLs.

385 In a learning problem, the features are essential. Table I presents all the features  
 386 for learning to rank tweets. Some features in Table I are explained in detail in the  
 387 following. For T-tweets, the applicable features are computed based on their tweet  
 388 message fields, whereas for TU-tweets, they are computed based on their three fields  
 389 as well as the union of the three fields. For example, the BM25 similarity between a  
 390 query and a T-tweet  $d$  can be computed based on the tweet message field of  $d$ ; for a  
 391 TU-tweet, four BM25 similarities can be computed, one based on the tweet message  
 392 field, one based on the URL title field, one based on the URL body field and the last  
 393 one based on the union of these three fields. Different degrees of significance can be  
 394 associated with the different fields by the learning model. It has been shown that  
 395 improvement in ranking can be achieved by weighting the fields of documents (for  
 396 example, the titles of documents vs. the bodies of documents) differently [Robertson  
 397 et al. 2004]. In our opinion, the same can apply to the tweets. Thus, we propose the  
 398 features whose calculations are based on the different fields of tweets together with  
 399 queries. During the establishment of the rankers, different weights are learned for  
 400 those different field-based features.

401 Moreover, the features can be categorized into two types: tweet-related (*TR* for short)  
 402 and query-tweet-related (*QTR* for short). The former type is calculated purely based  
 403 on the tweets themselves. For example, for feature  $F_{13}$ , it is a Boolean feature indi-  
 404 cating whether the tweet has at least an embedded URL. Studies [Duan et al. 2010;  
 405 McCreddie et al. 2011; Metzler and Cai 2011] showed that whether a tweet has a URL  
 406 is an effective feature for ranking tweets. Intuitively, the Web pages of the URLs em-  
 407 bedded in tweets often provide more information than tweets' 140 characters. Thus,  
 408 a tweet with embedded URLs has a higher probability of being relevant than a tweet  
 409 without embedded URLs [Duan et al. 2010].

410 Besides the tweet-related features, the query-tweet-related features are also used  
 411 to calculate different lexical similarities between queries and tweets. In addition to  
 412 capturing term similarities, such as BM25 similarities discussed before, our method  
 413 also computes concept similarities as features. A concept is a proper noun (*PN*), a  
 414 dictionary phrase (*DP*), a simple noun phrase (*SNP*), or a complex noun phrase (*CNP*).  
 415 A dictionary phrase is a noun phrase that can be looked up in dictionaries such as  
 416 Wikipedia but is not a proper noun. A simple noun phrase (complex noun phrase)  
 417 consists of two (more than two) nonstop terms but is neither a proper noun nor a  
 418 dictionary phrase. A concept is recognized in a document if all of its nonstop terms  
 419 appear in the document within a text window of certain size, with the smallest window  
 420 size for *PNs*, then a bigger window size for *DPs*, an even bigger window size for *SNPs*,  
 421 and the largest window size for *CNPs*. Please refer to the papers [Liu et al. 2004; Zhang  
 422 et al. 2007] for the details about these concepts. In this article, we adopt the phrase  
 423 recognition tool [Zhang et al. 2007] to identify the four types of concepts from queries  
 424 and tweets. This tool can achieve an accuracy of 92% in recognizing concepts.

Table I. Features for Ranking Tweets

ID	Type	Feature Description ( $q$ = query, $T$ = tweet).	No.
$F_1$	<i>QTR</i>	The percentage of the terms of $q$ contained by the hashtags of $T$ . The hashtags are the keywords or topics of $T$ and they appear in the tweet message field of $T$ by prefixing the symbol “#”.	1
$F_2$	<i>QTR</i>	The percentage of the expansion terms of $q$ contained by the hashtags of $T$ . The expansion terms are obtained by the pseudo relevance feedback method [Liu et al. 2004].	1
$F_3$	<i>QTR</i>	Whether the four fields (the three fields of a TU-tweet and their union) contain $q$ as an <i>SNP</i> or <i>CNP</i> respectively.	4
$F_4$	<i>QTR</i>	The frequency of $q$ in $T$ as an <i>SNP</i> or <i>CNP</i> .	1
$F_5$	<i>QTR</i>	Whether the four fields contain a key term of $q$ , if exist. The key term is the nonverb term in $q$ , satisfying the following two conditions: 1) it has the least document frequency among all query terms; 2) it is not a term in a <i>PN</i> or a <i>DP</i> concept.	4
$F_6$	<i>TR</i>	The length of the tweet message field of $T$ . [Duan et al. 2010; McCreddie et al. 2011; Metzler and Cai 2011]	1
$F_7$	<i>QTR</i>	Whether the four fields contain all <i>PN</i> or <i>DP</i> query concepts.	4
$F_8$	<i>QTR</i>	The sum of the frequencies of all <i>PN</i> or <i>DP</i> query concepts in $T$ .	1
$F_9$	<i>QTR</i>	The percentage of the nonverb terms of $q$ contained in the four fields.	4
$F_{10}$	<i>QTR</i>	The (weighted) percentage of the query concepts contained in the four fields. All query concepts are either equally weighted or weighted by their inverse document frequencies.	8
$F_{11}$	<i>QTR</i>	BM25 and TFIDF similarities between $q$ and the four fields. [Duan et al. 2010; McCreddie et al. 2011]	8
$F_{12}$	<i>TR</i>	Whether $T$ (or the Web pages of embedded URLs) has more than 50% content in English. [McCreddie et al. 2011; Metzler and Cai 2011]	1
$F_{13}$	<i>TR</i>	Whether $T$ has at least one URL in its tweet message field. [Duan et al. 2010; McCreddie et al. 2011; Metzler and Cai 2011]	1
$F_{14}$	<i>TR</i>	The count of the Twitter user of $T$ mentioned by the tweets in the collection. [Duan et al. 2010]	1
$F_{15}$	<i>TR</i>	Whether $T$ is a retweet (or a reply tweet). [Duan et al. 2010; Metzler and Cai 2011]	2
$F_{16}$	<i>QTR</i>	The percentage of the related concepts of $q$ contained in the four fields. The related concepts of $q$ are the top three frequent <i>PN</i> concepts among the top 10 web documents retrieved by Google with respect to $q$ .	4
$F_{17}$	<i>QTR</i>	The percentage of the related nouns of $q$ contained in the four fields. The related nouns are the nouns with the top three document frequencies among the top 10 web documents retrieved by Google with respect to $q$ .	4
$F_{18}$	<i>QTR</i>	Whether the order of query terms appearing in the four fields is the same as that in $q$ .	4

425 A query can be represented by a set of concepts as illustrated by the following  
 426 example.

427 *Example 2.* Given a query of “*Australian Open Djokovic vs. Murray*”, it contains five  
 428 concepts. They are three *PN* concepts, “*Australian Open*,” “*Djokovic*” and “*Murray*,” an  
 429 *SNP* concept, “*Djokovic Murray*” (“*vs.*” is omitted as a stop word) and a *CNP* concept,  
 430 “*Australian Open Djokovic Murray*.”

431 We propose the features (say  $F_{10}$ ) involving query concepts because they capture the  
 432 similarities between queries and tweets better than query terms as illustrated by the  
 433 following example.

434 *Example 3.* Given the query  $q$  = “*Australian Open Djokovic vs. Murray*”, a T-tweet  
 435  $d_1$  = “*and Djokovic it is... Murray becoming more like England football team...failing*  
 436 *where it matters...*” and a T-tweet  $d_2$  = “*Can’t stop watching the Australian Open!*”,  $d_1$   
 437 contains two query terms, “*Djokovic*” and “*Murray*” and  $d_2$  also contains two query  
 438 terms, “*Australian*” and “*Open*”. But  $d_1$  is relevant to  $q$  while  $d_2$  is irrelevant. In terms  
 439 of query concepts,  $d_1$  contains three out of five query concepts, “*Djokovic*”, “*Murray*”  
 440 and “*Djokovic Murray*” but  $d_2$  contains only one query concept, “*Australian Open*”.

441 There are eight features with the ID of  $F_{10}$ . One of the features is the percentage of  
 442 the query concepts contained in the tweet message field. As illustrated by Example 3,  
 443 the more query concepts a tweet contains, the more likely the tweet is relevant to  
 444 the query. The value of this feature for  $d_1$  is  $3/5$  while that for  $d_2$  is  $1/5$ . Another  
 445 member feature is the weighted percentage of the query concepts contained in the  
 446 tweet message field. Since a concept can be weighted by its inverse document frequency  
 447 (*idf* for short), the weighted percentage of the query concepts contained in the tweet  
 448 message field is the ratio of the sum of the *idfs* of the query concepts contained in the  
 449 tweet message field over the sum of the *idfs* of all query concepts. If we consider the  
 450 four fields of TU-tweets (the three fields and their union), eight such features can be  
 451 calculated over the four fields of TU-tweets accordingly.

452 The features with the IDs of  $F_{16}$  and  $F_{17}$  calculate the numbers of the related con-  
 453 cepts and the related nouns of queries in the different fields of tweets. A person who  
 454 writes a tweet specifies an event by a set  $S_1$  of concepts or terms. A person who queries  
 455 the same event may utilize another set  $S_2$  of concepts or terms. The concepts or terms  
 456 in  $S_1$  are related to those in  $S_2$ . Let us illustrate these features with the following  
 457 Example.

458 *Example 4.* Given a query “White House spokesman replaced” and a T-tweet  $d_1 =$   
 459 “Jay Carney named as Barack Obama’s press secretary,”  $d_1$  is relevant to the query,  
 460 although it does not contain any query concepts or terms. “Jay Carney” is a related  
 461 concept to the query, as it is one of the three most frequent *PN* concepts from the top  
 462 10 Web documents retrieved by Google with respect to the query. Therefore, the match  
 463 of “Jay Carney” is an indicator of  $d_1$ ’s relevance to the query.

464 To build the two tweet type-specific rankers, we partition TREC relevance judg-  
 465 ments of tweets into a set of labeled T-tweets and a set of labeled TU-tweets. We  
 466 use the former set of T-tweets as the training data for learning a *T-tweet Ranker*  
 467 and the latter set of TU-tweets for learning a *TU-tweet Ranker*, respectively. For  
 468 building a *T-tweet ranker*, we convert each training example (T-tweet) into a vector  
 469 of the proposed features that are applicable for T-tweets. Then, we feed the vectors  
 470 of features into RankSVM to generate a *T-tweet Ranker*. We repeat the same pro-  
 471 cedure as before by using the training data for TU-tweets to generate a *TU-tweet*  
 472 *Ranker*.

### 473 3.3. Preference Classifier

474 The two tweet type-specific rankers only provide the preference between two tweets of  
 475 the same type. In order to merge the rankings of T-tweets and TU-tweets, a classifier  
 476 is proposed to determine the preference of each T-tweet with respect to each TU-tweet.  
 477 We employ the SVM model [Joachims 1999] to perform such determination. In partic-  
 478 ular, each training example is a triple of  $\langle d_1, d_2, label \rangle$ , where  $d_1$  is a T-tweet,  $d_2$  is  
 479 a TU-tweet and the *label* indicates whether  $d_1$  is preferred over  $d_2$  or vice versa. We  
 480 again use TREC relevance judgments as the training data. Specifically, for a training  
 481 query, a labeled T-tweet  $d_1$  and a labeled TU-tweet  $d_2$  form a training example (pair),  
 482 only if their labels of relevance to that query are different. The different labels of  $d_1$   
 483 and  $d_2$  imply that  $d_1$  is preferred over  $d_2$  or vice versa.

484 To learn such a classifier, we reuse the features in Table I and they are referred to  
 485 as *ranking features*. We also propose a set of new features that captures the differ-  
 486 ence of the corresponding (ranking) features of  $d_1$  and  $d_2$  with respect to a query. Let  
 487 us call this set of new features *dependent features*. Each dependent feature aims at a  
 488 direct comparison of relevance between  $d_1$  and  $d_2$ . It is calculated by a T-tweet (rank-  
 489 ing) feature minus a corresponding TU-tweet (ranking) feature. For example, given

490 the feature group  $F_{11}$ , a T-tweet feature is the BM25 similarity between a query  $q$  and  
 491 the tweet message field of  $d_1$ . But four corresponding TU-tweet features are the BM25  
 492 similarities between  $q$  and the four fields of  $d_2$ , respectively. Thus, four dependent fea-  
 493 tures are obtained by subtracting the four TU-tweet features from the T-tweet feature,  
 494 respectively. A (preference) classifier can be learned by using these features and the  
 495 training examples. In our preliminary experiments, the classifier using both ranking  
 496 features and dependent features performed better than the classifiers that just use  
 497 either ranking features or dependent features.

### 498 3.4. Greedy Merging Algorithm

499 After we build the two tweet type-specific rankers and the preference classifier, we can  
 500 rank tweets with respect to a test query  $q'$ . First, we use these two rankers to rank  
 501 T-tweets and TU-tweets with respect to  $q'$  separately. Then, we employ the preference  
 502 classifier to compute the preference between any two tweets, one from each ranking.  
 503 This constitutes three sets of preferences: one for any two T-tweets, one for any two  
 504 TU-tweets and one for any T-tweet and any TU-tweet. The goal is to merge the two  
 505 rankings into a ranking that agrees with these three sets of preferences as much as  
 506 possible. Cohen et al. [1998] showed that the problem of finding the ordering that  
 507 agrees best with a given set of preferences is NP-complete. Therefore, we propose a  
 508 quadratic greedy merging algorithm. To merge a ranking of T-tweets and a ranking  
 509 of TU-tweets, this algorithm always picks the tweet that has the smallest sum of the  
 510 degrees of the preferences of other tweets (that have not been picked) over it. This  
 511 makes the merged ranking consistent with the three sets of preferences, if there is no  
 512 inconsistency among the three sets of preferences.

513 Let  $T$  and  $TU$  be a ranking of T-tweets and a ranking of TU-tweets, respectively.  
 514 They are defined as follows. We assign (numerical) subscripts to the T-tweets in  $T$  so  
 515 that the T-tweets with smaller subscripts have higher preferences. The same applies  
 516 to  $TU$ . For convenience of presentation, we give the T-tweets in  $T$  the subscripts from 1  
 517 to  $m$  and the TU-tweets in  $TU$  the subscripts from  $m + 1$  to  $m + n$ . But the comparison  
 518 between a subscript of a T-tweet and that of a TU-tweet does not indicate a preference  
 519 between them.

$$520 \quad \begin{aligned} T &= [d_1, \dots, d_m] && \text{s.t. } d_i > d_j, 1 \leq i < j \leq m \\ TU &= [d_{m+1}, \dots, d_{m+n}] && \text{s.t. } d_i > d_j, m + 1 \leq i < j \leq m + n. \end{aligned} \quad (1)$$

522 Let  $f_p : \Omega^T \times \Omega^{TU} \rightarrow R$  be a preference function which maps a pair of a T-tweet  $d_i$   
 523 and a TU-tweet  $d_j$  to a real number.  $\Omega^T$  and  $\Omega^{TU}$  are the T-tweet space and the TU-  
 524 tweet space, respectively. If the real number is positive,  $d_i > d_j$ ; if it is negative, the  
 525 reverse is true; if it is zero, there is no preference between  $d_i$  and  $d_j$ . The magnitude  
 526 of the number indicates the degree of the preference. We assume that the real number  
 527 being zero does not occur, which is true in practice. This function corresponds to the  
 528 preference classifier (see Section 3.3). Let  $D$  be the union of  $T$  and  $TU$ ,  $D = T \cup TU =$   
 529  $[d_1, \dots, d_m, d_{m+1}, \dots, d_{m+n}]$ . Let  $Pref(d_i, d_j)$  denote the preference between a tweet  $d_i$   
 530 and another tweet  $d_j$  in  $D$ .  $Pref(d_i, d_j)$  can be defined as follows.

$$531 \quad Pref(d_i, d_j) = \begin{cases} d_i > d_j & 1 \leq i < j \leq m \\ d_i > d_j & m + 1 \leq i < j \leq m + n \\ d_i > d_j & 1 \leq i \leq m < m + 1 \leq j \leq m + n \text{ and } f_p(d_i, d_j) > 0 \\ d_j > d_i & 1 \leq i \leq m < m + 1 \leq j \leq m + n \text{ and } f_p(d_i, d_j) < 0. \end{cases} \quad (2)$$

532 Let  $RP(i)$  be the ranking position of a tweet  $d_i$  in  $T$  or  $TU$ . Due to the subscript  
533 assignments given to the T-tweets in  $T$  and the TU-tweets in  $TU$ ,  $RP(i)$  is defined as  
534 follows.

$$535 \quad RP(i) = \begin{cases} i & 1 \leq i \leq m \\ i - m & m + 1 \leq i \leq m + n. \end{cases} \quad (3)$$

536 Let  $M = [M_{ij}]_{(m+n) \times (m+n)}$  be the preference matrix for  $D$  as defined here. It is con-  
537 sistent with Equation (2) and has the following interpretation: 1)  $M_{ij} > 0$  indicates  
538  $d_i > d_j$ ; 2)  $M_{ij} < 0$  indicates  $d_j > d_i$ ; 3) the absolute value of  $M_{ij}$  represents the degree  
539 of the preference, which is normalized between 0 and 1. Moreover, we propose three  
540 weighting parameters,  $\lambda_T (> 0)$ ,  $\lambda_{TU} (> 0)$  and  $\lambda_{Pairwise} (> 0)$ , to be set to the degrees  
541 that we trust the three sets of preferences.

$$542 \quad [M_{ij}]_{(m+n) \times (m+n)} = \begin{cases} \lambda_T \cdot \frac{RP(j) - RP(i)}{\max\{RP(i), RP(j)\}} & 1 \leq i, j \leq m \\ \lambda_{TU} \cdot \frac{RP(j) - RP(i)}{\max\{RP(i), RP(j)\}} & m + 1 \leq i, j \leq m + n \\ \lambda_{Pairwise} \cdot \frac{f_p(d_i, d_j)}{\max_{1 \leq s \leq m < m+1 \leq t \leq m+n} \{|f_p(d_s, d_t)|\}} & 1 \leq i \leq m < m + 1 \leq j \leq m + n \\ -M_{ji} & 1 \leq j \leq m < m + 1 \leq i \leq m + n. \end{cases} \quad (4)$$

543 We now explain why  $M$  is defined in such a manner. Specifically, we elaborate the  
544 intuition of each of the four components of  $M$ .

- 545 (1) The first component  $\left(\lambda_T \cdot \frac{RP(j) - RP(i)}{\max\{RP(i), RP(j)\}}\right)$  indicates the preference between any two  
546 T-tweets,  $d_i$  and  $d_j$ . If  $1 \leq i < j \leq m$ , then  $RP(i) < RP(j)$  and therefore  $M_{ij} > 0$ ,  
547 indicating  $d_i > d_j$ ; if  $1 \leq j < i \leq m$ , then  $RP(j) < RP(i)$  and therefore  $M_{ij} < 0$ ,  
548 indicating  $d_j > d_i$ . The degree of the preference is normalized between 0 and 1  
549 by  $\max\{RP(i), RP(j)\}$ . Moreover, it is also easy to verify that  $M_{ij} < M_{i(j+1)}$  if  $1 \leq$   
550  $i \leq m, 1 \leq j \leq m - 1$ . This is reasonable, because as the separation between two  
551 T-tweets increases, so is the degree of the preference. We propose such a heuristic  
552 method to measure the degree of the preference between two T-tweets, because  
553 most learning to rank algorithms, such as RankSVM, produce the ranking scores  
554 that have no meaning in an absolute sense and can only be used for ordering.
- 555 (2) The second component  $\left(\lambda_{TU} \cdot \frac{RP(j) - RP(i)}{\max\{RP(i), RP(j)\}}\right)$  has the same interpretation as the  
556 first component, except that it indicates the preference between any two TU-  
557 tweets,  $d_i$  and  $d_j$ .
- 558 (3) The third component  $\left(\lambda_{Pairwise} \cdot \frac{f_p(d_i, d_j)}{\max_{1 \leq s \leq m < m+1 \leq t \leq m+n} \{|f_p(d_s, d_t)|\}}\right)$  indicates the prefer-  
559 ence between a T-tweet  $d_i$  and a TU-tweet  $d_j$ . If  $f_p(d_i, d_j) > 0$ , then  $M_{ij} > 0$  and  
560  $d_i > d_j$ ; if  $f_p(d_i, d_j) < 0$ , then  $M_{ij} < 0$  and  $d_j > d_i$ . The degree of the preference is  
561 normalized between 0 and 1 by  $\max_{1 \leq s \leq m < m+1 \leq t \leq m+n} \{|f_p(d_s, d_t)|\}$ .
- 562 (4) The fourth component indicates that the preference between a TU-tweet  $d_i$  and a  
563 T-tweet  $d_j$  is the negation of the preference between  $d_j$  and  $d_i$ .

564 Let us illustrate the preference matrix  $M$  with the following example.

565 *Example 5.* Given two T-tweets,  $d_1$  and  $d_2$  and three TU-tweets:  $d_3, d_4$  and  $d_5$ , the  
566 three sets of the preferences of these tweets are shown in Table II.

Table II. Three Sets of Preferences for Example 5

Rankers and Classifier	Tweet Preferences and Their Ranking Positions
<i>T-tweet Ranker</i>	$d_1 > d_2$ ; $RP(d_1) = 1$ and $RP(d_2) = 2$ ;
<i>TU-tweet Ranker</i>	$d_3 > d_4 > d_5$ ; $RP(d_3) = 1$ , $RP(d_4) = 2$ and $RP(d_5) = 3$ ;
<i>Preference Classifier</i>	$f_p(d_1, d_3) = -0.9(d_3 > d_1)$ ; $f_p(d_2, d_3) = -1(d_3 > d_2)$ $f_p(d_1, d_4) = -0.7(d_4 > d_1)$ ; $f_p(d_2, d_4) = -0.8(d_4 > d_2)$ $f_p(d_1, d_5) = 0.6(d_1 > d_5)$ ; $f_p(d_2, d_5) = 0.5(d_2 > d_5)$

567 For simplicity, we assume that the three weighting parameters:  $\lambda_T$ ,  $\lambda_{TU}$  and  $\lambda_{Pairwise}$   
 568 are all equal to 1. The preference matrix for Example 5 is shown here.

$$569 \quad M = \begin{bmatrix} 0 & 0.5 & -0.9 & -0.7 & 0.6 \\ -0.5 & 0 & -1 & -0.8 & 0.5 \\ 0.9 & 1 & 0 & 0.5 & 0.67 \\ 0.7 & 0.8 & -0.5 & 0 & 0.33 \\ -0.6 & -0.5 & -0.67 & -0.33 & 0 \end{bmatrix}.$$

570 To merge the two rankings, we propose a greedy merging algorithm. To explain the  
 571 proposed merging algorithm, we first define *Dispreferness*.

572 *Definition 3.3 (Dispreferness)*. Given the preference matrix  $M$  and a tweet  $d_i$ , the  
 573 *Dispreferness* of the tweet  $d_i$  is calculated by

$$574 \quad \text{Dispreferness}(M, d_i) = \sum_j |\min\{0, M_{ij}\}|. \quad (5)$$

575 Given a tweet  $d_i$ , if it is preferred over a tweet  $d_j$ , then  $M_{ij} > 0$  and  $|\min\{0, M_{ij}\}| = 0$   
 576 will not contribute to  $\text{Dispreferness}(M, d_i)$ . On the other hand, if  $d_j$  is preferred over  
 577  $d_i$ , then  $M_{ij} < 0$  and  $|\min\{0, M_{ij}\}|$  contributes a positive value to  $\text{Dispreferness}(M, d_i)$ .  
 578  $\text{Dispreferness}(M, d_i)$  is the sum of the degrees of the preferences of other tweets over  $d_i$ .  
 579 The greedy merging algorithm, called *GreedyMerging*, merges two rankings of tweets  
 580 by placing the tweet  $d$  with the least  $\text{Dispreferness}(M, d)$  in the first position of the  
 581 merged ranking  $L$ . Placing  $d$  in such a position of  $L$  may incur a certain amount of  
 582 inconsistency and this amount is  $\text{Dispreferness}(M, d)$ . Compared to any other tweet  
 583 placed at the first position, this amount of inconsistency is the least. Then, after re-  
 584 moving  $d$  from the matrix  $M$  and re-computing the  $\text{Dispreferness}$  of other tweets, it  
 585 iteratively places the tweet that has the least  $\text{Dispreferness}$  in the next position in  $L$ .  
 586 The algorithm always picks the tweet that incurs the least amount of inconsistency at  
 587 the time it is picked. Details of the algorithm are shown in Algorithm 1.

588 The following proposition demonstrates that the proposed algorithm is theoretically  
 589 reasonable, because if there is no inconsistency among the three sets of preferences,  
 590 the optimal ranking of tweets will be achieved by *GreedyMerging*.

591 **PROPOSITION 3.4.** *If there is no inconsistency among all the preferences from the T-*  
 592 *tweet Ranker, the TU-tweet Ranker and the pairwise classifier, GreedyMerging produces*  
 593 *the optimal ranking.*

594 **PROOF.** Assuming no inconsistency among all the preferences, there must be a lin-  
 595 ear order of tweets in terms of their preferences:  $d_{i_1} > d_{i_2} > \dots > d_{i_n}$ . This linear order  
 596 is an optimal ranking of tweets because any pair of tweets is ordered by their prefer-  
 597 ences. The first tweet  $d_{i_1}$  has zero  $\text{Dispreferness}$  because no tweet has preference over  
 598 it. Moreover, no other tweet, say  $d$ , has zero  $\text{Dispreferness}$ , since  $d_{i_1}$  is preferred over  
 599  $d$ , causing  $\text{Dispreferness}(M, d) > 0$ . *GreedyMerging* inserts  $d_{i_1}$  into the first position of  
 600 the merged ranking  $L$ . After  $d_{i_1}$  is chosen and the matrix is updated by deleting the

**ALGORITHM 1:** The *GreedyMerging* Algorithm

**Input:** A ranking of T-tweets:  $T$ ; a ranking of TU-tweets:  $TU$ ; the preferences of pairs of a T-tweet and a TU-tweet,  $f_p$ ; Three weighting parameters:  $\lambda_T$ ,  $\lambda_{TU}$  and  $\lambda_{Pairwise}$ ;

**Output:** A merged ranking of tweets  $L$ ;

1. Union two rankings of tweets  $D = T \cup TU$ ;
2. Create the preference matrix  $M_{|D| \times |D|}$  for  $D$ , based on  $T$ ,  $TU$ ,  $f_p$ ,  $\lambda_T$ ,  $\lambda_{TU}$  and  $\lambda_{Pairwise}$ ;
3. **while**( $D \neq \emptyset$ )
  4. Find the tweet  $d$  with the least *Dispreferness*( $M, d$ );
  5.  $d = \arg \min_{d \in D} \{Dispreferness(M, d)\}$ ;
  6. Insert  $d$  into the merged ranking  $L$ ;
  7. Update  $D$  and  $M$ :
  8.  $D = D - \{d\}$ ;
  9.  $M_{|D-1| \times |D-1|} = M_{|D| \times |D|} - [d]$ ; // deleting the row and column representing  $d$ ;
10. **end**

601 row and the column representing  $d_{i_1}$ , the second tweet  $d_{i_2}$  has no tweet preferred over  
 602 it among the remaining tweets and only its *Dispreferness* is zero. *GreedyMerging* in-  
 603 serts  $d_{i_2}$  into the second position of  $L$ . The same argument is applied repeatedly until  
 604 all tweets are inserted into  $L$ .  $\square$

605 After a ranking of T-tweets and a ranking of TU-tweets are merged by *GreedyMerg-*  
 606 *ing*, we obtain a ranking  $L$  of both types of tweets but their IR scores are absent. We  
 607 need to assign some (pseudo) IR scores to the tweets in  $L$  so that the time-related rele-  
 608 vance scores of tweets (to be given in Section 4.2) can be combined with the IR scores  
 609 to yield the similarity scores for the final ranking of tweets (see Section 4.3). The rank-  
 610 ing of the tweets in descending order of their pseudo IR scores should be identical  
 611 to  $L$ . We adopt the conversion proposed in Lee [1997]. Given a ranking of  $n$  tweets,  
 612  $L = [d_1, \dots, d_n]$ , where the subscript  $i$  of tweet  $d_i$  is its ranking position, we assign  $d_i$   
 613 an IR score  $IR(d_i)$  as follows.

$$614 \quad IR(d_i) = 1 - \frac{i-1}{n}. \quad (6)$$

#### 615 4. TEMPORAL USAGE IN RETRIEVAL

616 In the first phase, tweets are ranked by only considering their lexical similarities to  
 617 queries. In this section, we discuss how to use the temporal information (publishing  
 618 times) of tweets to improve retrieval effectiveness.

##### 619 4.1. Time Representation

620 In this section, we describe the temporal representation of tweets with respect to  
 621 queries. Each query  $q$  has a timestamp  $t$  and only the tweets published on or before  
 622  $t$  are considered to be relevant. Given a tweet  $d$  with a publishing time  $t_d$ , we adopt  
 623 the time representation  $f(t_d, t)$  proposed in Efron and Golovchinsky [2011] with the  
 624 interpretation that  $f(t_d, t) = 0$  means the tweet  $d$  is published on the same day as  $t$   
 625 and  $f(t_d, t) = n$  ( $n > 0$ ) indicates the tweet  $d$  is published  $n$  days before  $t$ .

##### 626 4.2. Query Type Determination

627 In this section, we first propose a method to classify queries by the temporal distri-  
 628 butions of their top tweets and then present different ways to measure the temporal  
 629 relevance of tweets to classified queries.

630 There are three types of queries as discussed in Section 1. We utilize the top tweets  
 631 from the first phase to classify a query into one of these three types. Specifically, for  
 632 a query  $q$  with a timestamp  $t$ , let  $D = \{d_1, \dots, d_K\}$  be the top  $K$  tweets retrieved by

633 the divide-and-conquer method in the first phase. Let  $T = \{t_1, \dots, t_K\}$  be the set of  
 634 publishing times associated with those top  $K$  tweets, where each publishing time  $t_i$   
 635 presents either the same time  $t$  or a time before  $t$ . Let  $T_D = \{t'_i | t'_i = f(t_m, t), t_m \in T\}$  be  
 636 the set of the unique time representations of their publishing times. Let  $I(t_j, t, t'_i)$  be an  
 637 indicator function.

$$638 \quad I(t_j, t, t'_i) = \begin{cases} 1 & f(t_j, t) = t'_i \\ 0 & \text{otherwise.} \end{cases} \quad (7)$$

639 The type of  $q$  can be classified as follows.

640 —  $q$  is a time insensitive query if the largest proportion of the top  $K$  tweets pub-  
 641 lished on a single day is less than or equal to a certain threshold  $p (\leq 0.5)$ , that is,  
 642 Equation (8) holds.

$$643 \quad \max_{t'_i \in T_D} \left\{ \frac{1}{K} \sum_{t_j \in T} I(t_j, t, t'_i) \right\} \leq p \leq 0.5. \quad (8)$$

644 —  $q$  is a dominant peak query if the largest proportion of the top  $K$  tweets pub-  
 645 lished on a certain single day (say  $t'$ ) is greater than a threshold  $s (> p)$ , that is,  
 646 Equation (9) holds. Its dominant peak is on  $t'$ .

$$647 \quad \max_{t'_i \in T_D} \left\{ \frac{1}{K} \sum_{t_j \in T} I(t_j, t, t'_i) \right\} > s > p. \quad (9)$$

648 —  $q$  is a nondominant peak query if the largest proportion of the top  $K$  tweets pub-  
 649 lished on a single day is less than or equal to  $s$  but greater than  $p$ , that is, Equation  
 650 (10) holds. It can have a set of nondominant peaks and the proportion of the top  $K$   
 651 tweets at each peak is less than or equal to  $s$  but greater than  $p$ .

$$652 \quad s \geq \max_{t'_i \in T_D} \left\{ \frac{1}{K} \sum_{t_j \in T} I(t_j, t, t'_i) \right\} > p. \quad (10)$$

653 The parameters  $K$ ,  $p$  and  $s$  are estimated empirically. After a query  $q$  is classified  
 654 into one of the three types, the tweets from the first phase are assigned time-related  
 655 relevance scores (*TRSs* for short) to  $q$  as follows.

656 — If  $q$  is a time-insensitive query, all the tweets retrieved from the first phase are not  
 657 assigned any *TRSs*. This implies that time has no impact on ranking the tweets  
 658 with respect to  $q$ .

659 — If  $q$  is a dominant peak query, that is, the temporal distribution of its top  $K$  tweets  
 660 has a dominant peak on  $t'_i$  (the  $t'_i$  days before  $t$ ), a tweet  $d$  (published on  $t_d$ ) is  
 661 assigned a *TRS* as follows.

$$662 \quad TRS(t_d, t) = \frac{1}{2\delta} \exp \left\{ -\frac{|f(t_d, t) - t'_i|}{\delta} \right\}. \quad (11)$$

663 This function is in the form of the Laplace distribution [Laplace 1774]. When the  
 664 tweet occurs at the peak, its *TRS* is normalized by  $\max_{t_d} \{TRS(t_d, t)\}$  to be 1. The  
 665 farther the tweet  $d$  is temporally away from the peak, the smaller the *TRS* of  $d$   
 666 is. In other words, tweets temporally closer to the peak are given higher *TRSs*.  
 667 We tested different exponential functions and found that the Laplace-like function



668 performed best. It has a single peak on  $t'_i$  and its variance  $2\delta^2$  can be estimated by  
 669 the maximum likelihood method.

$$670 \quad \hat{\delta} = \frac{1}{|T_D|} \sum_{t'_i \in T_D} \left| \frac{1}{K} \sum_{t_j \in T} I(t_j, t, t'_i) - \hat{\mu} \right| \text{ s.t. } \hat{\mu} = \frac{1}{|T_D|} \sum_{t'_i \in T_D} \left( \frac{1}{K} \sum_{t_j \in T} I(t_j, t, t'_i) \right). \quad (12)$$

671 — If  $q$  is a nondominant peak query, that is, the temporal distribution of its top  
 672  $K$  tweets has a set of nondominant peaks at a set of time representations  $P =$   
 673  $\{t'_1, \dots, t'_p\}$ , a tweet (published on  $t_d$ ) is assigned a *TRS* as follows.

$$674 \quad TRS(t_d, t) = \begin{cases} \frac{\sum_{t_j \in T} I(t_j, t, t'_n)}{\max_{t'_m \in P} \left\{ \sum_{t_j \in T} I(t_j, t, t'_m) \right\}} \cdot \frac{\sum_{d' \in D_{t'_n}} BM25(d, d')}{|D_{t'_n}|} & f(t_d, t) \notin P \\ \frac{\sum_{t_j \in T} I(t_j, t, f(t_d, t))}{\max_{t'_m \in P} \left\{ \sum_{t_j \in T} I(t_j, t, t'_m) \right\}} & f(t_d, t) \in P \end{cases} \quad (13)$$

$$\text{s.t. } D_{t'_m} = \{d' | f(t_{d'}, t) = t'_m, t'_m \in P\}, t'_n = \arg \max_{t'_m \in P} \left\{ \frac{\sum_{d' \in D_{t'_m}} BM25(d, d')}{|D_{t'_m}|} \right\}$$

675 Let us explain the intuition of Equation (13) as follows.

676 (1) Suppose that the distribution of  $q$ 's top  $K$  tweets has multiple nondominant  
 677 peaks.

678 (a) For a tweet (published on  $t_d$ ) belonging to the highest peak at time  
 679  $f(t_d, t)$ , its *TRS* is assigned to be 1, that is,  $\max_{t'_m \in P} \left\{ \sum_{t_j \in T} I(t_j, t, t'_m) \right\} =$

$$680 \quad \sum_{t_j \in T} I(t_j, t, f(t_d, t)) \Rightarrow \frac{\sum_{t_j \in T} I(t_j, t, f(t_d, t))}{\max_{t'_m \in P} \left\{ \sum_{t_j \in T} I(t_j, t, t'_m) \right\}} = 1$$

681 (b) For a tweet  $d$  (published on  $t_d$ ) belonging to a nonhighest peak at time  
 682  $f(t_d, t)$ , its *TRS* is the ratio of the number of the top  $K$  tweets at that peak  
 683 to that at the highest peak, that is,  $TRS(t_d, t) = \frac{\sum_{t_j \in T} I(t_j, t, f(t_d, t))}{\max_{t'_m \in P} \left\{ \sum_{t_j \in T} I(t_j, t, t'_m) \right\}}.$

684 (c) For a tweet  $d$  (published on  $t_d$ ) not belonging to any peak, we first deter-  
 685 mine which peak contains the tweets that are most similar to  $d$ . We use  
 686 *BM25* to measure the average similarity of  $d$  to the tweets at a peak.<sup>2</sup>  
 687 Then we pick the peak with the highest average similarity to  $d$ , say the

688 peak at time  $t'_n$ . Let  $S_2 \left( = \frac{\sum_{d' \in D_{t'_n}} BM25(d, d')}{|D_{t'_n}|} \right)$  denote that highest average  
 689 similarity. Each tweet in that picked peak is assigned the same *TRS*. Let

690  $S_1 \left( = \frac{\sum_{t_j \in T} I(t_j, t, t'_n)}{\max_{t'_m \in P} \left\{ \sum_{t_j \in T} I(t_j, t, t'_m) \right\}} \right)$  denote that *TRS* of a tweet in that picked  
 691 peak. Finally we assign  $d$  a *TRS* that is the product of  $S_1$  and  $S_2$ . In other  
 692 words, the tweets in different peaks describe different events related to  $q$ .

693 We first determine which related event  $d$  is likely to describe. The likeli-  
 694 hoods of  $d$  describing different events are measured by the average similar-  
 695 ities of  $d$  to those tweets at different peaks. We then assign  $d$  a *TRS* that  
 696 is equal to the highest average similarity multiplied by the *TRS* of a tweet  
 697 describing the same related event as  $d$  does.

<sup>2</sup>We utilize the tweet message field without exploring the Web pages of URLs if present.

- 698 (2) Suppose that the distribution of  $q$ 's top  $K$  tweets has a single nondominant  
 699 peak, the same approach is used.  
 700 (a) For a tweet belonging to the unique peak, its  $TRS$  is assigned to be 1.  
 701 (b) For a tweet  $d$  that does not belong to that peak, the average similarity of  $d$   
 702 to the tweets in that peak is computed. It is multiplied by the  $TRS$  of any  
 703 tweet at the peak (having a value of 1 due to the single peak) to yield the  
 704  $TRS$  of  $d$ .

#### 705 4.3. Aggregation of IR Scores and Time-Related Relevance Scores

706 The first phase calculates the IR scores of tweets with respect to a query  $q$ . The second  
 707 phase of the method calculates the time-related relevance scores of tweets by using  
 708 temporal information. Given a tweet  $d$ , let  $IR(d)$  and  $TRS(d)$  be the IR score of  $d$   
 709 and the time-related relevance score of  $d$ , respectively. An aggregation score  $AGS(d)$   
 710 can be calculated in the manner of F-measure [Rijsbergen 1979] (see Equation (14)).  
 711 The tweets are arranged in descending order of the aggregation scores. Although the  
 712 F-measure is usually used as an evaluation measure, it can be employed to balance  
 713  $IR(d)$  and  $TRS(d)$ . The parameter  $\beta$  aims at balancing the contributions of  $IR(d)$  and  
 714  $TRS(d)$  to the aggregation score. The appropriate value of  $\beta$  is estimated in the experi-  
 715 ments. Experimental results demonstrate that such an aggregation outperforms other  
 716 aggregations, such as CombSUM and CombMNZ [Shaw et al. 1994].

$$717 \quad AGS(d) = (1 + \beta^2) \frac{IR(d) \cdot TRS(d)}{\beta^2 \cdot IR(d) + TRS(d)} \quad (14)$$

### 718 5. EXPERIMENT SETUP

#### 719 5.1. TREC Tweets2011 Collection

720 TREC 2011 released a tweet collection called Tweets2011 for the real-time ad-hoc re-  
 721 trieval task of the microblog track. The collection consists of about 16 million tweets  
 722 sampled from Twitter over 17 days (from 1/23/2011 to 2/8/2011). Instead of directly giv-  
 723 ing those tweets, TREC 2011 provided two tools for participating groups to crawl the  
 724 collection. One tool employing a Twitter API provides an information-rich collection of  
 725 tweets in the JSON format. The other one just crawls the HTML pages of tweets. The  
 726 efficiency of the first tool is very low, crawling about 150 tweets per hour due to the  
 727 limitation of the Twitter API. The second tool only crawls the HTML pages of tweets  
 728 and it is far more efficient than the first tool. However, some social information, such as  
 729 Twitter user profile, is absent in the HTML collection of tweets. We utilize the second  
 730 tool in this article. Since Twitter users might delete their tweets at any time, change  
 731 their usernames or change the public sharing properties of their tweets, it is possible  
 732 that some tweets are successfully crawled by some groups while become unavailable  
 733 when other groups are crawling. The statistics of our crawled tweet collection is shown  
 734 in Table III. In the TREC Tweets2011 collection crawled by us, 16.7% of tweets are  
 735 TU-tweets. We crawled the Web pages whose URLs are linked by the TU-tweets in the  
 736 collection, which results in another collection of about 2.3 million Web pages.<sup>3</sup>

#### 737 5.2. TREC 2011 and 2012 Queries and TREC Relevance Judgments

738 TREC 2011 released 50 queries and TREC 2012 released 60 queries. TREC required  
 739 both sets of queries to be retrieved over the TREC Tweets2011 collection. Each query  
 740 represents an information need at a specific time. An example query is shown in  
 741 Figure 2. The num tag encloses the ID of the query. The query tag encloses the query.

<sup>3</sup>Some URLs given by the TU-tweets are not available during our crawling.

Table III. The Statistics of Our Crawled TREC Tweets2011 Collection

HTTP Response Code	Tweet Count	Description
200 (OK)	14437978	Successfully downloaded tweets.
302 (Found)	1612080	Downloaded retweets via redirects.
403 (Forbidden)	339147	The tweets without public sharing properties.
404 (Not Found)	707403	The tweets no longer available.

```

<top>
<num> Number: MB075 </num>
<query> Aguilera super bowl fail </query>
<querytime> Tue Feb 08 21:56:22 +0000 2011 </querytime>
<querytweetime> 35094611483426816 </querytweetime>
</top>

```

Fig. 2. An example of TREC query.

742 The querytime tag gives the timestamp of the query in the form of ISO standard. Each  
743 tweet is assigned a unique tweet ID. The descending ordering of the IDs of tweets can  
744 be interpreted as the reverse-chronological order of their publishing times. The query-  
745 tweetime tag represents the timestamp of the query. In response to a query with a  
746 timestamp  $t$ , only the tweets whose IDs are not greater than  $t$  need to be considered.

747 TREC also provided the relevance judgments of tweets with respect to those two  
748 sets of queries. TREC assessors read tweets, then followed the URLs inside them and  
749 finally labeled them in a three point scale: “highly relevant,” “relevant,” and “irrele-  
750 vant.” For the TREC 2011 queries, 49 (out of 50) queries have at least one relevant  
751 or highly relevant tweet and 33 (out of 50) queries have at least one highly relevant  
752 tweet. For the TREC 2012 queries, 59 (out of 60) queries have at least one relevant or  
753 highly relevant tweet and 56 (out of 60) queries have at least one highly relevant tweet.  
754 “Highly relevant” tweets are preferred over “relevant” tweets that are preferred over  
755 “irrelevant” tweets. For the set of TREC 2011 queries, we use the set of TREC 2012  
756 queries as the training query set and their corresponding TREC relevance judgments  
757 as the training data and vice versa.

### 758 5.3. Relevance Criteria

759 There are two relevant criteria: 1) both relevant and highly relevant tweets are con-  
760 sidered relevant; 2) only the highly relevant tweets are considered relevant. In our  
761 experiments, we denote these two relevant criteria as the *relevant criterion* and the  
762 *highly relevant criterion*, respectively. Our results are evaluated by these two criteria.

### 763 5.4. Evaluation Measures

764 In this article, we employ the precision at top 30 tweets (P30 for short), the mean  
765 average precision (MAP for short) and the normalized discounted cumulative gain at  
766 top 30 tweets (NDCG@30 for short) as the evaluation measures. To evaluate the re-  
767 trieval effectiveness of our method that does not involve ranking tweets in reverse-  
768 chronological order, we use MAP as the primary measure and P30 and NDCG@30 as  
769 the secondary measures. However, we use P30 as the primary measure and MAP as  
770 the secondary measure to evaluate the performance of our method in ranking tweets in  
771 reverse-chronological order, as TREC 2011 stipulated that P30 is the official measure  
772 for the reverse-chronological rankings of tweets [Ounis et al. 2011]. In this article, we  
773 only consider statistical significance at  $p < 0.05$  according to one-sided paired t-test.

## 774 6. EXPERIMENTAL RESULTS

775 In this section, we evaluate our method by using both TREC 2011 and TREC 2012  
 776 queries over the TREC Tweets2011 collection. Two sets of experiments are conducted  
 777 to evaluate our two-phase method. One set evaluates the retrieval performance of the  
 778 divide-and-conquer method in the first phase; the other set evaluates that of utilizing  
 779 temporal information of tweets in the second phase. We also compare the performance  
 780 of our two-phase method with various state-of-the-art methods. In particular, we con-  
 781 duct the experiments to reveal the answers to the following research questions.

- 782 — Is it beneficial to apply the divide-and-conquer strategy on ranking tweets? In other  
 783 words, would there be any benefit to rank the two types of tweets separately, com-  
 784 pared with the method of ranking them simultaneously? Experiments are conducted  
 785 to verify the motivation of leveraging the structural difference of tweets.
- 786 — What are the important features for learning to rank tweets? We study the degrees  
 787 of importance of the proposed features for ranking T-tweets, TU-tweets and both  
 788 types of tweets together.
- 789 — What are the effectiveness and the efficiency of the proposed divide-and-conquer  
 790 algorithm for ranking tweets?
- 791 — How many queries do benefit from the divide-and-conquer algorithm and how many  
 792 queries do not? In particular, we conduct a result analysis of the proposed algorithm  
 793 and discuss the reasons why our algorithm helps or hurts some typical queries.
- 794 — Is it necessary to have two different types of time sensitive queries (dominant peak  
 795 queries vs. nondominant peak queries)? Experiments are conducted to validate the  
 796 benefit of our proposed categories of temporal queries.
- 797 — How to estimate the parameters  $K$ ,  $p$ ,  $s$  and  $\beta$  that are used by our temporal classi-  
 798 fication of queries?
- 799 — Does the utilization of temporal information provide further improvement over the  
 800 algorithm using the divide-and-conquer strategy?
- 801 — How many queries do benefit from the usage of temporal information and how many  
 802 queries do not? We analyze the performance of our method query by query and  
 803 discuss the reasons why our method improves or deteriorates the performance of  
 804 some queries.
- 805 — How is the performance of our two-phase method that combines the usage of tem-  
 806 poral information with the divide-and-conquer approach, compared with various  
 807 state-of-the-art methods?

### 808 6.1. Relevance Ranking Analysis

809 In this section, we first demonstrate the necessity of considering the structural differ-  
 810 ence of tweets. Second, we study the degrees of importance of the proposed features for  
 811 ranking tweets. Third, we study the effectiveness of the divide-and-conquer method by  
 812 comparing it with various baselines. Fourth, we discuss the efficiency of the proposed  
 813 method. Finally, we conduct a result analysis and discuss why some queries are helped  
 814 or hurt by our method.

815 *6.1.1. The Motivation of Considering Structural Difference of Tweets.* To validate the motiva-  
 816 tion of using the divide-and-conquer strategy to address the structural difference of  
 817 tweets, we analyze a uniform ranker (denoted by *Uniform Ranker*) and the two tweet  
 818 type-specific rankers (denoted by *T-tweet Ranker* and *TU-tweet ranker* respectively).  
 819 The *Uniform Ranker* is constructed by using RankSVM. It is learned over the training  
 820 data consisting of a set of training queries and both types of labeled tweets. It ranks  
 821 both types of tweets simultaneously. We first apply the *Uniform Ranker* to produce  
 822 a ranking of tweets. This ranking  $R$  consists of both types of tweets and is then

Table IV. *Uniform Ranker* vs. Tweet Type-Specific Rankers

	TREC 2011					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
<i>Uniform Ranker</i> (for T-tweets)	0.0613	0.1497	0.1142	0.0231	<b>0.0202</b>	0.1030
<i>T-tweet Ranker</i>	<b>0.0768</b> †	<b>0.1639</b>	<b>0.1327</b> †	<b>0.0297</b>	0.0152	<b>0.1151</b>
<i>Uniform Ranker</i> (for TU-tweets)	0.4440	0.5013	0.4762	0.3966	<b>0.2364</b>	0.4831
<i>TU-tweet Ranker</i>	<b>0.4715</b> †	<b>0.5102</b>	<b>0.4952</b> †	<b>0.4042</b>	0.2242	<b>0.4923</b>
	TREC 2012					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
<i>Uniform Ranker</i> (for T-tweets)	0.0474	0.1266	0.0670	0.0182	0.0411	<b>0.0706</b>
<i>T-tweet Ranker</i>	<b>0.0510</b>	<b>0.1373</b>	<b>0.0678</b>	<b>0.0184</b>	<b>0.0446</b>	0.0696
<i>Uniform Ranker</i> (for TU-tweets)	0.2882	0.4226	0.2798	0.2447	0.2435	0.2722
<i>TU-tweet Ranker</i>	<b>0.2926</b> †	<b>0.4367</b> †	<b>0.2949</b>	<b>0.2489</b> †	<b>0.2548</b>	<b>0.2829</b>

Note: † indicates statistically significant improvements over the corresponding baselines

823 partitioned into two rankings,  $R_1$  for T-tweets and  $R_2$  for TU-tweets. The relative  
824 order of the tweets in each  $R_i$  ( $i = 1, 2$ ) is the same as that in  $R$ . Two tweet type-specific  
825 rankers are constructed by using RankSVM too. The *T-tweet Ranker* is learned by only  
826 using the portion of T-tweets in the training data and the *TU-tweet ranker* is learned  
827 by only using the portion of TU-tweets. They are used to rank the two types of tweets  
828 separately. Finally, we compare the performance of these two rankings  $R_1$  and  $R_2$  with  
829 those of the two corresponding rankings from the two tweet type-specific rankers.  
830 The performance is evaluated by using both relevant criteria. The comparison of their  
831 performance is shown in Table IV.

832 We make three observations based on the information shown in Table IV. First,  
833 the *Uniform Ranker* achieves decent performance in ranking TU-tweets but it per-  
834 forms poorly in ranking T-tweets with respect to both sets of TREC 2011-2012 queries.  
835 Second, the two tweet type-specific rankers consistently outperform the *Uniform*  
836 *Ranker* in terms of MAP, P30 and NDCG@30 by the relevant criterion over both sets  
837 of queries. Third, for the highly relevant criterion, the two type-specific rankers show  
838 somewhat stronger performance than the *Uniform Ranker*. Specifically, for the TREC  
839 2011 queries, the two rankers consistently outperform the *Uniform Ranker* in MAP  
840 and NDCG@30 but get marginal deteriorations in P30. For the TREC 2012 queries,  
841 the *TU-tweet Ranker* consistently outperforms the *Uniform Ranker* in all three mea-  
842 sures. The *T-tweet Ranker* outperforms the *Uniform Ranker* in terms of MAP and P30  
843 but gets a negligible deterioration in NDCG@30. These three observations validate the  
844 motivation and the necessity of treating the two types of tweets separately.

845 **6.1.2. Feature Analysis.** It is worth investigating the degrees of importance of the pro-  
846 posed features for learning to rank tweets. We sort the proposed features in descend-  
847 ing order of their degrees of importance that are calculated by RankSVM [Bian et al.  
848 2010]. Specifically, we study the degrees of importance of the features applicable for  
849 the *T-tweet Ranker*, the *TU-tweet Ranker* and the *Uniform Ranker*. Table V shows the  
850 top 10 important features for each of these three rankers.

851 From Table V, several observations can be made. First, *QTR* features (the features  
852 whose calculations depend on tweets and queries) are more important than *TR* fea-  
853 tures (the features whose calculations depend on tweets only) in ranking T-tweets,  
854 TU-tweets or ranking them simultaneously, because *QTR* features dominate the top  
855 10 features for these three rankers. Second, the top 10 features for the *T-tweet Ranker*  
856 are very different from those for the *TU-tweet Ranker*. In particular, only 3 of the top

Table V. Top 10 Features for *T-tweet Ranker*, *TU-tweet Ranker*, and *Uniform Ranker*

Top 10 Features for <i>T-tweet Ranker</i>				Shared by Rankers Below	
Rank	ID	Type	Feature Description	TU-Tweet	Uniform
1	$F_{17}$	<i>QTR</i>	The percentage of related nouns of $Q$ contained in the tweet message field.	✓	✓
2	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the tweet message field	✓	✓
3	$F_{10}$	<i>QTR</i>	The weighted percentage of query concepts contained in the tweet message field	✓	✓
4	$F_{12}$	<i>TR</i>	Whether the tweet message field has more than 50% content in English.		
5	$F_{16}$	<i>QTR</i>	The percentage of related concepts of $Q$ contained in the tweet message field.		
6	$F_{18}$	<i>QTR</i>	Whether the order of query terms in the tweet message field is the same as that of in the query		
7	$F_3$	<i>QTR</i>	Whether the tweet message field contains the whole query as a <i>SNP</i> or <i>CNP</i> .		
8	$F_1$	<i>QTR</i>	The percentage of query terms contained by the hashtags in the tweet.		
9	$F_5$	<i>QTR</i>	Whether the tweet message field contains the key query term.		
10	$F_2$	<i>QTR</i>	The percentage of expansion terms contained by the hashtags in the tweet.		
Top 10 Features for <i>TU-tweet Ranker</i>				Shared by Rankers Below	
Rank	ID	Type	Feature Description	T-Tweet	Uniform
1	$F_{10}$	<i>QTR</i>	The weighted percentage of query concepts contained in the union of all three fields.		✓
2	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the URL title field.		✓
3	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the URL body field.		✓
4	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the tweet message field.	✓	✓
5	$F_3$	<i>QTR</i>	Whether the URL title field contains the whole query as a <i>SNP</i> or <i>CNP</i> .		✓
6	$F_{17}$	<i>QTR</i>	The percentage of related nouns of $Q$ contained in the tweet message field.	✓	✓
7	$F_{10}$	<i>QTR</i>	The weighted percentage of query concepts contained in the URL body field.		✓
8	$F_{10}$	<i>QTR</i>	The weighted percentage of query concepts contained in the tweet message field.		✓
9	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the union of all three fields.	✓	✓
10	$F_3$	<i>QTR</i>	Whether the URL body field contains the whole query as a <i>SNP</i> or <i>CNP</i> .		✓
Top 10 Features for <i>Uniform Ranker</i>				Shared by Rankers Below	
Rank	ID	Type	Feature Description	T-Tweet	TU-tweet
1	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the tweet message field.	✓	✓
2	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the URL title field.		✓
3	$F_{10}$	<i>QTR</i>	The weighted percentage of query concepts contained in the tweet message field.	✓	✓
4	$F_{17}$	<i>QTR</i>	The percentage of related nouns of $Q$ contained in the tweet message field.	✓	✓

Table V. Continued

Top 10 Features for <i>Uniform Ranker</i>				Shared by Rankers Below	
Rank	ID	Type	Feature Description	T-Tweet	TU-tweet
5	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the union of all three fields.		✓
6	$F_{10}$	<i>QTR</i>	The percentage of query concepts contained in the URL body field.		✓
7	$F_3$	<i>QTR</i>	Whether the URL title field contains the whole query as a <i>SNP</i> or <i>CNP</i> .		✓
8	$F_{10}$	<i>QTR</i>	The weighted percentage of query concepts contained in the union of all three fields.		✓
9	$F_{10}$	<i>QTR</i>	The weighted percentage of query concepts contained in the URL body field.		✓
10	$F_3$	<i>QTR</i>	Whether the URL body field contains the whole query as a <i>SNP</i> or <i>CNP</i> .		✓

857 10 features for the *T-tweet Ranker* appear among the top 10 important features for  
858 *TU-tweet Ranker* and they are not among the top 3 features for the *TU-tweet Ranker*.  
859 This observation shows that the *T-tweet Ranker* and the *TU-tweet Ranker* emphasize  
860 different features and thus again verifies the motivation and the necessity of ranking  
861 these two types of tweets separately. Third, the top 10 important features for the  
862 *T-tweet Ranker* are quite different from those for the *Uniform Ranker* while the top  
863 10 important features for the *TU-tweet Ranker* are very similar to those for the *Uni-*  
864 *form Ranker*. In particular, only 3 of the top 10 features for the *T-tweet Ranker* appear  
865 among those for the *Uniform Ranker* while all the top 10 features for the *TU-tweet*  
866 *Ranker* are the same as those for the *Uniform Ranker* but with a different order. This  
867 observation explains why the *Uniform Ranker* achieves decent performance in ranking  
868 TU-tweets but suffers poor performance in ranking T-tweets.

869 **6.1.3. The Impact of the Divide-and-Conquer Method.** To study the impact of our divide-  
870 and-conquer method, four systems are configured. The first system is BM25 similarity  
871 [Robertson et al. 1996]. We empirically learn the two parameters  $b$  and  $k$  for BM25.  
872 In particular, the parameter  $b$  is learned from 0.5 to 1 with an interval of 0.05 and the  
873 parameter  $k$  is learned from 1.2 to 2.0 with an interval of 0.1. The combination of these  
874 two parameters that optimizes the performance of the TREC 2011 queries is applied  
875 to the TREC 2012 queries and vice versa. The second system is the *Uniform Ranker*  
876 (see Section 6.1.1). These two methods act as the baselines. The third system is the  
877 proposed divide-and-conquer method equipped with a simple merging (called *Simple-*  
878 *Merging*) algorithm. It can act as an alternative to the *GreedyMerging* algorithm to  
879 merge the rankings of T-tweets and TU-tweets. The *SimpleMerging* algorithm works  
880 as follows. Given a ranking of T-tweets, a ranking of TU-tweets and the preferences  
881 of T-tweets relative to TU-tweets, *SimpleMerging* compares the preference between  
882 the first T-tweet and the first TU-tweet. If the first T-tweet is preferred over the  
883 first TU-tweet, *SimpleMerging* puts the first T-tweet into the merged ranking and  
884 then compares the preference between the second T-tweet and the first TU-tweet.  
885 Otherwise, *SimpleMerging* puts the first TU-tweet into the merged ranking and then  
886 compares the first T-tweet with the second TU-tweet. Repeat the given comparison un-  
887 til all tweets are merged into the final ranking. *SimpleMerging* guarantees to preserve  
888 the relative ranking positions of the T-tweets and those of the TU-tweets. Its time  
889 complexity is linear. The fourth system is the divide-and-conquer method equipped  
890 with the *GreedyMerging* algorithm and its time complexity is quadratic. The three  
891 parameters,  $\lambda_T$ ,  $\lambda_{TU}$  and  $\lambda_{Pairwise}$ , of *GreedyMerging* are estimated as follows. We stip-  
892 ulate that the sum of the three parameter values be 1 and each parameter can only be

Table VI. The Comparison of the Divide-and-Conquer Method of *SimpleMerging* or *GreedyMerging* with *Uniform Ranker* and *BM25*

	TREC2011					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
<i>BM25</i>	0.3693	0.3966	0.3747	0.2488	0.1576	0.3474
<i>Uniform Ranker</i>	0.4778 ↑	0.4905 ↑	0.4880 ↑	0.3788 ↑	0.2000 ↑	0.4793 ↑
<i>SimpleMerging</i>	0.4953 ↑	0.5109 ↑	0.4914 ↑	0.3912 ↑	0.2152 ↑	0.4882 ↑
<i>GreedyMerging</i>	<b>0.5006</b> ↑ ‡	<b>0.5143</b> ↑	<b>0.4939</b> ↑	<b>0.4090</b> ↑ ‡	<b>0.2283</b> ↑ †	<b>0.4933</b> ↑
	TREC2012					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
<i>BM25</i>	0.2603	0.3791	0.2207	0.1910	0.2167	0.2319
<i>Uniform Ranker</i>	0.3077 ↑	0.4175 ↑	0.2705 ↑	0.2345 ↑	0.2375 ↑	0.2633 ↑
<i>SimpleMerging</i>	0.3206 ↑	0.4130 ↑	0.2832 ↑	0.2409 ↑ †	0.2357 ↑	0.2710 ↑
<i>GreedyMerging</i>	<b>0.3259</b> ↑	<b>0.4367</b> ↑	<b>0.2966</b> ↑ †‡	<b>0.2590</b> ↑ †‡	<b>0.2583</b> ↑	<b>0.2852</b> ↑ †‡

Note: ↑, †, and ‡ indicate statistically significant improvements over *BM25*, *Uniform Ranker* and *SimpleMerging*, respectively.

893 assigned one of 10 possible values: 0.1, ..., 1.0. The combination of these three parameters  
 894 that optimizes the performance of the TREC 2011 queries is applied to the TREC  
 895 2012 queries and vice versa. The performances of these systems are shown in Table VI.

896 Several observations can be made from the information in Table VI. First, all three  
 897 learning to rank models, the *Uniform Ranker*, the divide-and-conquer method with the  
 898 *SimpleMerging* algorithm (the *SimpleMerging* algorithm for short) and the divide-and-  
 899 conquer method with the *GreedyMerging* algorithm (the *GreedyMerging* algorithm for short)  
 900 consistently and significantly outperform the *BM25* baseline in all measures by  
 901 both criteria with respect to the two sets of queries. This indicates that using learning  
 902 to rank techniques benefits the retrieval effectiveness of tweets. Second, for the set of  
 903 TREC 2011 queries, the *SimpleMerging* algorithm consistently outperforms the *Uni-*  
 904 *form Ranker* in all the measures by both criteria; for the set of TREC 2012 queries,  
 905 the *SimpleMerging* algorithm consistently outperforms the *Uniform Ranker* in MAP  
 906 and NDCG@30 but gets negligible deteriorations in P30 by both criteria. For all the  
 907 measures with respect to the two sets of TREC queries, the *GreedyMerging* algorithm  
 908 consistently outperforms the *Uniform Ranker* baseline by both relevant criteria. This  
 909 observation validates that the retrieval effectiveness of tweets benefits from the em-  
 910 ployment of the divide-and-conquer strategy for handling the structural difference  
 911 of tweets. Third, the *GreedyMerging* algorithm consistently outperforms the *Simple-*  
 912 *Merging* algorithm in all the measures by both relevant criteria with respect to the  
 913 two sets of queries. This indicates the performance of the *GreedyMerging* algorithm is  
 914 superior to that of the *SimpleMerging* algorithm.

915 **6.1.4. The Efficiency of GreedyMerging Algorithm.** To merge a ranking of  $m$  T-tweets and  
 916 a ranking of  $n$  TU-tweets, the time complexity of the *GreedyMerging* algorithm is  
 917  $O((m+n)^2)$ . It consists of the construction of the preference matrix  $M$  and the merging  
 918 process based on  $M$ . Compared with the *SimpleMerging* algorithm, the *GreedyMerg-*  
 919 *ing* algorithm is not very efficient when  $m$  and  $n$  are large. However, its quadratic time  
 920 complexity should not be problematic when only merging the top  $m'$  T-tweets and the  
 921 top  $n'$  TU-tweets, where  $m' \ll m$  and  $n' \ll n$ . Merging the top  $m'$  T-tweets and the top  
 922  $n'$  TU-tweets makes the construction of the preference matrix efficient, since we only  
 923 construct the submatrix based on these top tweets. It also makes the merging process



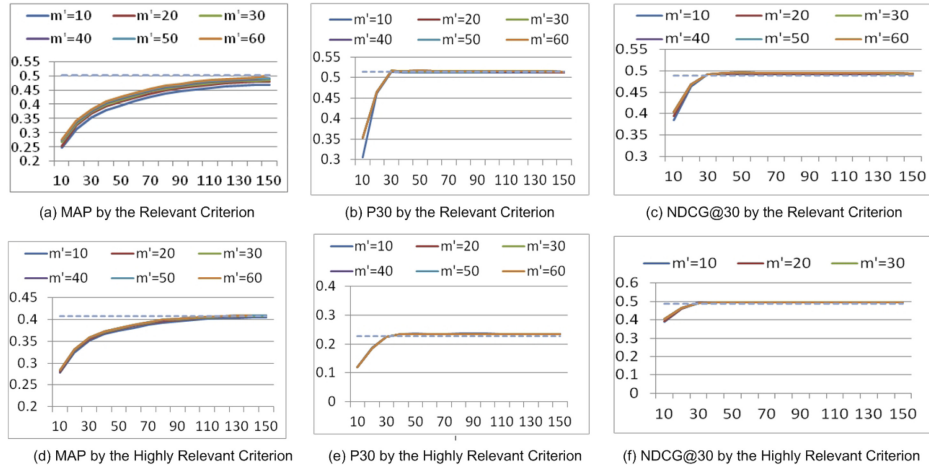


Fig. 3. Performance of *GreedyMerging* with the Varying Values of  $m'$  and  $n'$  with respect to the TREC 2011 queries.

924 efficient, because the *GreedyMerging* algorithm merges tweets by their *Dispreferness*  
 925 that now is calculated on this small submatrix.

926 We study the effectiveness of the *GreedyMerging* algorithm when  $m'$  and  $n'$  are as-  
 927 signed small values. Figures 3 and 4 show the MAP, P30 and NDCG@30 performance  
 928 of the *GreedyMerging* algorithm with varying small values of  $m'$  and  $n'$  for both sets  
 929 of TREC 2011 and TREC 2012 queries, respectively. For the TREC 2011 queries, the  
 930 value of  $m'$  varies from 10 to 60 and that of  $n'$  varies from 10 to 150. For the TREC  
 931 2012 queries, the value of  $m'$  varies from 10 to 60 and that of  $n'$  varies from 10 to 300.  
 932 In all the component figures of Figures 3 and 4, the  $x$  axes represent the varying values  
 933 of  $n'$  and the  $y$  axes represent the MAP, P30 and NDCG@30 performance by either the  
 934 relevant criterion or the highly relevant criterion. The different curves represent the  
 935 varying values of  $m'$ . The dash lines represent the corresponding performance of the  
 936 *GreedyMerging* algorithm by merging all  $m$  T-tweets with all  $n$  TU-tweets. For ease of  
 937 presentation, let us denote as *FullGreedyMerging* the *GreedyMerging* algorithm that  
 938 merges all  $m$  T-tweets and all  $n$  TU-tweets.

939 Figure 3 shows the performance of the *GreedyMerging* algorithm by both relevant  
 940 criteria with respect to the set of TREC 2011 queries. According to Figure 3(a), which  
 941 shows the MAP performance by the relevant criterion, the *GreedyMerging* algorithm  
 942 achieves a comparable MAP score of 0.4987 when merging only the top 60 ( $m' = 60$ ) T-  
 943 tweets and the top 150 ( $n' = 150$ ) TU-tweets, relative to a MAP score of 0.5006 achieved  
 944 by *FullGreedyMerging*. A similar observation can be made based on Figure 3(d) where  
 945 the MAP performance is evaluated by the highly relevant criterion. The *GreedyMerg-*  
 946 *ing* algorithm achieves a comparable MAP score of 0.4017 when merging only the top  
 947 40 ( $m' = 40$ ) T-tweets and the top 90 ( $n' = 90$ ) TU-tweets, compared with a MAP score  
 948 of 0.4090 achieved by *FullGreedyMerging*. If the users are interested in the top tweets,  
 949 we can achieve comparable performance in terms of P30 and NDCG@30, when merg-  
 950 ing very few T-tweets and TU-tweets. According to Figure 3(b) (Figure 3(e)) where the  
 951 P30 performance is evaluated by the (highly) relevant criterion, we can achieve a P30  
 952 score of 0.5122 (0.2263) when just merging the top 10 ( $m' = 10$ ) T-tweets and the top  
 953 30 ( $n' = 30$ ) TU-tweets, compared with the P30 score of 0.5143 (0.2283) achieved by  
 954 *FullGreedyMerging*. Similar observations can be made based on the NDCG@30 per-  
 955 formance shown by Figure 3(c) and Figure 3(f). This indicates that we can make the

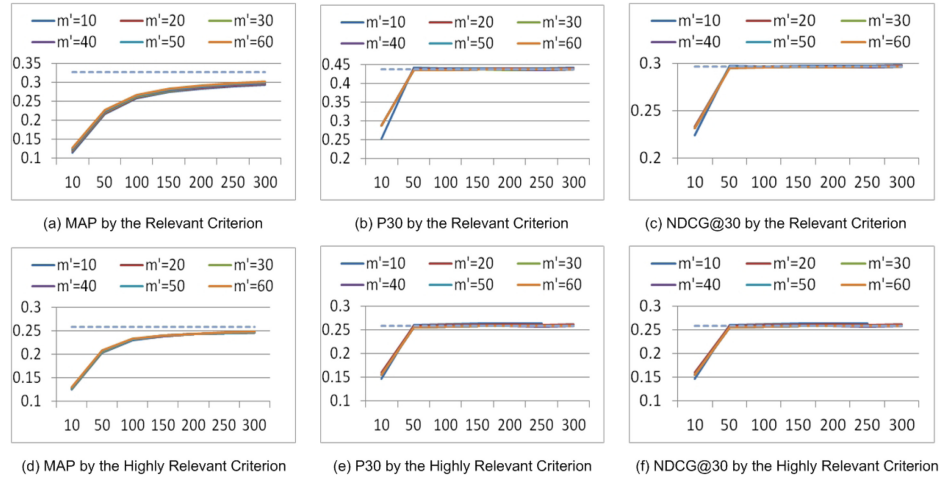


Fig. 4. Performance of *GreedyMerging* with the Varying Values of  $m'$  and  $n'$  with respect to the TREC 2012 queries.

956 *GreedyMerging* algorithm much more efficient without significantly hurting its effec-  
 957 tiveness for the TREC 2011 queries.

958 Figure 4 shows the performance of the *GreedyMerging* algorithm by both rele-  
 959 vant criteria with respect to the set of TREC 2012 queries. As shown in Figure 4(a)  
 960 (Figure 4(d)) where the MAP performance is evaluated by the (highly) relevant crite-  
 961 rion, the *GreedyMerging* algorithm achieves a reasonable MAP score of 0.3013 (0.2476)  
 962 by merging only the top 60 ( $m' = 60$ ) T-tweets and the top 300 ( $n' = 300$ ) TU-tweets,  
 963 relative to a MAP score of 0.3256 (0.2590) achieved by *FullGreedyMerging*. We note  
 964 that the TREC 2012 queries are harder than the TREC 2011 queries to achieve good  
 965 performance, which explains why we only achieve reasonable MAP performance for  
 966 the TREC 2012 queries by merging more top tweets than the TREC 2011 queries. If  
 967 only the top tweets are interested by users, we can achieve comparable performance  
 968 in P30 and NDCG@30 by merging very few top T-tweets and top TU-tweets. In partic-  
 969 ular, according to Figure 4(b) (Figure 4(e)) where the P30 performance is evaluated by  
 970 the (highly) relevant criterion, the *GreedyMerging* algorithm achieves a comparable  
 971 P30 score of 0.4340 (0.2571) by merging only the top 10 ( $m' = 10$ ) T-tweets and the  
 972 top 30 ( $n' = 30$ ) TU-tweets, relative to the P30 score of 0.4367 (0.2583) achieved by  
 973 *FullGreedyMerging*. Similar observations can be made based on the NDCG@30 perfor-  
 974 mance shown in Figure 4(c) and Figure 4(f). All these observations indicate that the  
 975 *GreedyMerging* algorithm can be much more efficient by achieving reasonable MAP  
 976 performance and comparable P30 and NDCG@30 performance for the TREC 2012  
 977 queries.

978 **6.1.5. Result Analysis.** In this section, we conduct an analysis for both sets of TREC  
 979 queries. Specifically, we compare the MAP performance of the *Uniform Ranker* with  
 980 that of the divide-and-conquer method using the *GreedyMerging* algorithm (see Ta-  
 981 ble VI). This comparison shows whether our way of handling the structural difference  
 982 of tweets can improve retrieval effectiveness. We analyze the average precision (*AP* for  
 983 short) changes query by query. Figure 5 shows the *AP* changes by both relevant  
 984 criteria with respect to the two sets of queries. For example, Figure 5(a) shows the *AP*  
 985 changes for the TREC 2011 queries by the relevant criterion. The changes are displayed  
 986 from the most improved query (on the left) to the most deteriorated query (on the right).  
 987 This displaying style continues from Figure 5(b) to 5(d). According to Figure 5, our

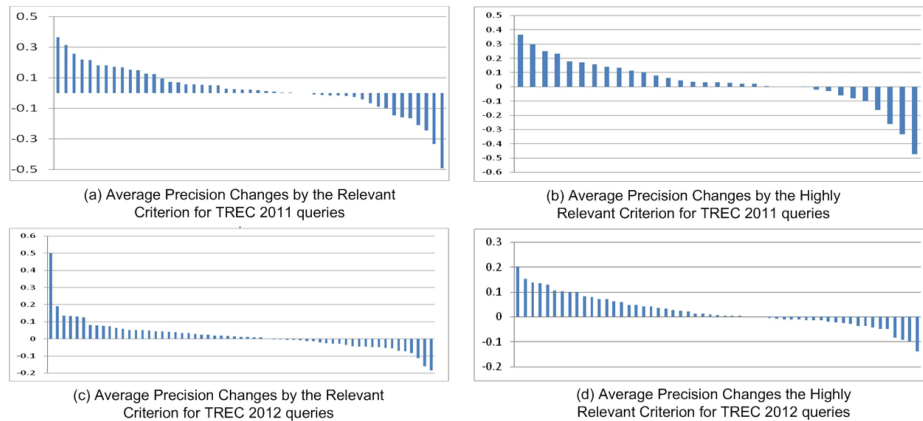


Fig. 5. AP Changes of the TREC 2011-2012 queries (*Uniform Ranker* vs. Divide-and-Conquer Method).

988 proposed method based on the divide-and-conquer strategy can improve the majority  
 989 of the queries by both relevant criteria for the two sets of queries. This validates the  
 990 effectiveness of our method.

991 We perform a deeper analysis of our results to find out how many queries are sig-  
 992 nificantly improved or hurt in their APs ( $\Delta AP \geq 0.1$ ) by our method and discuss the  
 993 corresponding reasons. For the TREC 2011 queries, Figure 5(a) shows that 13 queries  
 994 are significantly improved in their APs while 7 queries are significantly hurt accord-  
 995 ing to the relevant criterion. Figure 5(b) shows that 10 queries are significantly im-  
 996 proved while 4 queries are significantly hurt according to the highly relevant crite-  
 997 rion. For the TREC 2012 queries, Figure 5(c) shows that 6 queries are significantly im-  
 998 proved in their APs while 3 queries are significantly hurt according to the relevant  
 999 criterion. Figure 5(d) shows that 9 queries are significantly improved while only 1 query is sig-  
 1000 nificantly hurt according to the highly relevant criterion.

1001 One reason why our method improves some queries in their APs is that the *T-tweet*  
 1002 *Ranker* (see Section 6.1.1) outperforms the *Uniform Ranker* in ranking T-tweets for  
 1003 them. Let us illustrate this reason with an example.

1004 *Example 6.* The query  $q = \text{“Assange Nobel peace nomination”}$  has four concepts: two  
 1005 *PN* concepts, “Assange” and “Nobel peace,” and two *CNP* concepts, “Nobel peace nomi-  
 1006 nation” and “Assange Nobel peace nomination.” Given a T-tweet  $d_1 = \text{“Nobel war prize$   
 1007  $\text{for wikileaks... only if the nukes are fired... \#cablegate \#wikileaks \#assange \#anony-$   
 1008  $\text{mous”}$  and another T-tweet  $d_2 = \text{“\#unlikelyheadlines GEORGE BUSH WINS NOBEL$   
 1009  $\text{PEACE PRIZE! Ha,}”$ ,  $d_1$  is relevant to  $q$  while  $d_2$  is irrelevant to  $q$ . The *T-tweet Ranker*  
 1010 ranks  $d_1$  on top of  $d_2$ , because its most important feature is “the percentage of related  
 1011 nouns of the query contained in the tweet message field” (see Table V).  $d_1$  contains one  
 1012 related noun, “wikileaks”, while  $d_2$  does not contain any related nouns. The merged  
 1013 ranking preserves the ranking of  $d_1$  ahead of  $d_2$ . However, the *Uniform Ranker* ranks  
 1014  $d_2$  above  $d_1$ . For the most important feature of the *Uniform Ranker*, “the percentage  
 1015 of query concepts contained in the tweet message field” (see Table V),  $d_1$  contains a  
 1016 *PN* concept, “Assange” in its message field and  $d_2$  contains another *PN* concept “Nobel  
 1017 peace” in its message field too.  $d_1$  and  $d_2$  are tied. The second most important feature  
 1018 of the *Uniform Ranker*, “the percentage of query concepts contained in the URL title  
 1019 field” (see Table V), is not applicable for T-tweets. For the third most important feature  
 1020 of the *Uniform Ranker*, “the weighted percentage of query concepts contained in the  
 1021 tweet message field” (see Table V),  $d_2$  beats  $d_1$ , because the weight of “Nobel peace”

1022 is higher than that of “Assange.” We use the inverse document frequency of a concept  
 1023 as its weight. There are more tweets containing “Assange” than the tweets containing  
 1024 “Nobel peace” in our collection. Therefore, the *Uniform Ranker* ranks  $d_2$  above  $d_1$ .

1025 Another reason why our method improves some queries in their APs is that the  
 1026 *TU-tweet Ranker* (see Section 6.1.1) outperforms the *Uniform Ranker* in ranking TU-  
 1027 tweets for them. Let us illustrate this reason with an example.

1028 *Example 7.* The query  $q =$  “Supreme Court cases” has two concepts, a *PN* concept  
 1029 “Supreme Court” and a *CNP* concept “Supreme Court cases”. Given a TU-tweet  
 1030  $d_1 =$  “@enrogers Only FOX news... [http://www.foxnews.com/opinion/2010/01/22/  
 1031 ken-klukowski-supreme-court-amendment-mccain-feingold/](http://www.foxnews.com/opinion/2010/01/22/ken-klukowski-supreme-court-amendment-mccain-feingold/)” and another TU-tweet  
 1032  $d_2 =$  “Letter to Julia Gillard by Peter H Kemp - Solicitor of the Supreme Court of NSW  
 1033 <http://wlcentral.org/node/1175#assange#wikileaks>,”  $d_1$  is relevant while  $d_2$  is irrel-  
 1034 evant. The *Uniform Ranker* ranks  $d_2$  on top of  $d_1$ , because its most important feature is  
 1035 “the percentage of query concepts contained in the tweet message field” (see Table V).  
 1036  $d_1$  has no query concept in its message field while  $d_2$  has a query concept “Supreme  
 1037 Court” in its message field. The *TU-tweet Ranker* ranks  $d_1$  above  $d_2$ , because its most  
 1038 important feature is “the weighted percentage of query concepts contained in the union  
 1039 of all three fields” (see Table V). The Web page linked by the URL in  $d_1$  contains both  
 1040 query concepts, “Supreme Court” and “Supreme Court cases” while the Web page linked  
 1041 by the URL in  $d_2$  only contains “Supreme Court”. The merged ranking preserves the  
 1042 ranking of  $d_1$  ahead of  $d_2$ .

1043 The *T-tweet Ranker* and the *TU-tweet Ranker* are superior to the *Uniform Ranker* in  
 1044 ranking T-tweets and TU-tweets. Our merging algorithm preserves most of the pref-  
 1045 erences indicated by those two tweet type-specific rankers. So our divide-and-conquer  
 1046 method improves the majority of the queries.

1047 The reason why some queries suffer significant drops in their APs is that the *TU-*  
 1048 *tweet Ranker* falsely ranks some irrelevant TU-tweets over some relevant TU-tweets  
 1049 with respect to them. This happens for a small minority of queries, because the *TU-*  
 1050 *tweet Ranker* is not perfect. Let us illustrate this reason with a query “Michelle Obama  
 1051 fashion” whose performance is hurt most by both relevant criteria among the TREC  
 1052 2012 queries.

1053 *Example 8.* The query  $q =$  “Michelle Obama fashion” has two concepts, a *PN* con-  
 1054 cept, “Michelle Obama” and a *CNP* concept, “Michelle Obama fashion”. Given a TU-  
 1055 tweet  $d_1 =$  “Michelle Obama & Jill Biden Coordinate With Pearls On Monday (PHO-  
 1056 TOS, POLL) <http://huff.to/h4PmSg>” and a TU-tweet  $d_2 =$  “Fashionista: Fashion News  
 1057 Roundup: Franca Sozzani Trashes Fashion Bloggers, Cathy Horyn Throws Down Over  
 1058 Michell... <http://bit.ly/fTKBEs>,”  $d_1$  is relevant but  $d_2$  is irrelevant. The Web page  
 1059 linked by the URL in  $d_1$  presents the following excerpt: “And the pair did a little co-  
 1060 ordinating of their own – blazers and pearls. Michelle opted for a gray suit, with a  
 1061 necklace secured with a safety pin (super punk rock!), while Jill mixed cream with  
 1062 metallics and long necklace strands.”. This excerpt implicitly talks about the fashion  
 1063 aspect of “Michelle Obama”, although the query term “fashion” does not occur at all.  
 1064 However, consider an excerpt from the Web page linked by the URL in  $d_2$ , “Fashion  
 1065 News Roundup: Franca Sozzani Trashes Fashion Bloggers, Cathy Horyn Throws Down  
 1066 Over Michelle Obamas McQueen, and Naomi Campbells a No-Show in Court”. This  
 1067 excerpt contains all the query terms that form a *CNP* concept but is irrelevant to the  
 1068 query. The *TU-tweet Ranker* ranks  $d_2$  on top of  $d_1$ , because its most importance feature  
 1069 is “the weighted percentage of query concepts contained in the union of all three fields”  
 1070 (see Table V).  $d_2$  contains all the query concepts in the union of its three fields while  
 1071  $d_1$  only has a query concept, “Michelle Obama” in the union of its three fields. The

1072 ranking of  $d_2$  ahead of  $d_1$  is preserved in the merged ranking. However, the *Uniform*  
 1073 *Ranker* ranks  $d_1$  above  $d_2$ , because its most important feature is “the percentage of  
 1074 query concepts contained in the tweet message field” (see Table V).  $d_1$  has a query con-  
 1075 cept “*Michelle Obama*” in its message field while  $d_2$  does not have any query concepts  
 1076 in its message field.

1077 Again our merging algorithm can preserve most of the preferences of TU-tweets  
 1078 indicated by the *TU-tweet Ranker*. Given a query  $q$ , if its *AP* performance of TU-tweets  
 1079 achieved by the *TU-tweet Ranker* is significantly deteriorated, compared with that of  
 1080 the *Uniform Ranker*,  $q$  suffers a significant drop in the *AP* performance of its merged  
 1081 ranking.

## 1082 6.2. Improving Relevance Ranking via Temporal Information

1083 In this section, we present a set of experiments to evaluate our method that han-  
 1084 dles temporality. Specifically, we first validate our proposed three temporal cate-  
 1085 gories of queries. Then we evaluate our proposed F-measure aggregation method (see  
 1086 Section 4.3) by comparing it with two baseline aggregation methods, combSUM and  
 1087 combMNZ [Shaw et al. 1994]. Third, we show that the incorporation of the temporal  
 1088 information of tweets can further improve the retrieval effectiveness of the divide-and-  
 1089 conquer method in the first phase. Finally, we present a result analysis and discuss the  
 1090 reasons why our way of using temporality helps or hurts some queries.

1091 *6.2.1. The Validation of Temporal Query Categorizations.* In this section, we validate our  
 1092 three temporal categories of queries. We conduct the experiments in three scenarios  
 1093 where queries are classified into either time insensitive ones or time sensitive ones.  
 1094 In the first scenario, a classified time sensitive query is always treated as a dominant  
 1095 peak query, no matter how its top tweets are temporally distributed. The time-related  
 1096 relevance scores of the tweets with respect to it are calculated by the Laplace-like func-  
 1097 tion given by Equation (11). In the second scenario, a classified time sensitive query  
 1098 is always treated as a nondominant peak query, irrespective of the temporal distribu-  
 1099 tion of its top tweets. The time-related relevance scores of the tweets with respect to  
 1100 it are thus computed by Equation (13). In the third scenario, a time sensitive query is  
 1101 classified to be either a dominant peak query or a nondominant peak query. The time-  
 1102 related relevance scores of the tweets with respect to that query are calculated by  
 1103 Equation (11) or Equation (13), depending on its type. This is what we propose in this  
 1104 article. By comparing the results from the third scenario with those from the first two  
 1105 scenarios, we can conclude whether our temporal query categorization is necessary.

1106 Several parameters are proposed to temporally categorize queries, so we first discuss  
 1107 how to estimate them, which is followed by a description of how to configure the three  
 1108 scenarios. There are four parameters to be estimated. They are  $K$ ,  $p$ ,  $s$  and  $\beta$  (see  
 1109 Equations (8) to (10) and (14)). Given a query  $q$ , we first empirically use the top  $K$   
 1110 tweets of  $q$  to approximate the temporal distribution of the relevant tweets to  $q$ ; then  
 1111 we categorize  $q$  into one of three classes, after comparing the proportion of its top  
 1112  $K$  tweets at each day by  $p$  and  $s$ ; we calculate the time-related relevance scores of  
 1113 the tweets according to the classified type of  $q$ ; finally, we aggregate the IR scores of  
 1114 the tweets with their time-related relevance scores by  $\beta$ . We perform a grid search  
 1115 for estimating them. Specifically,  $K$  is estimated within the range from 10 to 60 at  
 1116 intervals of 10; the parameter  $p$  is estimated within the range from 0 to 0.5 at intervals  
 1117 of 0.1 to ensure  $p \leq 0.5$ ; the parameter  $s$  is estimated within the range from  $p + 0.1$   
 1118 to 1 at intervals of 0.1 to ensure  $s > p$ ; the parameter  $\beta$  is estimated within the range  
 1119 from 0 to 1 using the same interval length as  $p$  and  $s$ . For the TREC 2011 queries, we  
 1120 employ the TREC 2012 queries and their relevance judgments as the training data to

1121 estimate these four parameters and vice versa. We present the parameter estimation  
1122 method here.

- 1123 (1) Given a combination of a  $K$  value, a  $p$  value, and an  $s$  value within their ranges,  
1124 the training query set  $Q$  can be categorized into three subsets of queries:  $Q_A$  (time  
1125 insensitive queries),  $Q_B$  (time sensitive dominant peak queries) and  $Q_C$  (time sen-  
1126 sitive nondominant peak queries). The class of a training query  $q(\in Q)$  is deter-  
1127 mined by exploring the temporal distribution of its top  $K$  tweets (from the first  
1128 phase) jointly with the  $p$  value and the  $s$  value.
- 1129 (2) For the subset of the training queries,  $Q_A$ , parameter  $\beta$  is not estimated. Let  $MAP_A$   
1130 denote the MAP performance of the ranking of the tweets by their IR scores with  
1131 respect to  $Q_A$ .
- 1132 (3) For the subset of the training queries,  $Q_B$ , we iteratively test all possible values of  
1133 the parameter  $\beta$  for  $Q_B$ . Let  $\beta_B$  denote this parameter  $\beta$  for  $Q_B$ .
- 1134 (a) For each possible value of  $\beta_B$ , we first calculate the time-related relevance  
1135 scores ( $TRS$ s) of tweets with respect to  $Q_B$  by Equation (11), then aggregate the  
1136 IR scores of the tweets (with respect to  $Q_B$ ) with their  $TRS$ s by Equation (14)  
1137 by using the  $\beta_B$  value; finally we obtain a ranking of the tweets in descend-  
1138 ing order of their aggregated scores. The performance of this ranking can be  
1139 measured by a MAP score.
- 1140 (b) Find the  $\beta_B$  value (from Step 3.a) that corresponds to the highest MAP score.  
1141 Let  $MAP_B$  denote this highest MAP score for  $Q_B$ .
- 1142 (4) For the subset of the training queries  $Q_C$ , we iteratively test all possible values  
1143 of  $\beta$  for  $Q_C$ . Let  $\beta_C$  denote this parameter  $\beta$  for  $Q_C$ . Apply a similar method to  
1144 Step 3 on  $Q_C$  except that the  $TRS$ s of the tweets with respect to  $Q_C$  are computed  
1145 by Equation (13). Find the  $\beta_C$  value that corresponds to the highest MAP score for  
1146  $Q_C$  (denoted by  $MAP_C$ ).
- 1147 (5) Union the  $K$  value, the  $p$  value and the  $s$  value from Step 1, the  $\beta_B$  value from  
1148 Step 3, and the  $\beta_C$  value from Step 4 into a combination of five parameters. This  
1149 combination corresponds to a MAP performance for all training queries  $Q(= Q_A \cup$   
1150  $Q_B \cup Q_C)$  that can be calculated as follows. Let  $MAP_Q$  denote this MAP score.

$$1151 \quad MAP_Q = \frac{MAP_A \cdot |Q_A| + MAP_B \cdot |Q_B| + MAP_C \cdot |Q_C|}{|Q_A| + |Q_B| + |Q_C|} \quad (15)$$

- 1152 (6) Iteratively repeat Step 1–Step 5 with another combination of a  $K$  value, a  $p$  value  
1153 and an  $s$  value until all their possible combinations are iterated. Find the combina-  
1154 tion of  $K$ ,  $p$ ,  $s$ ,  $\beta_B$  and  $\beta_C$  that corresponds to the highest  $MAP_Q$ . This combination  
1155 is the set of the estimated parameter values.

1156 We provide some explanations for this method. Given a possible combination of a  
1157 value of  $K$ , a value of  $p$  and a value of  $s$ , we find out the value of the parameter  $\beta$  that  
1158 maximizes the MAP performance of all the training queries. Since the calculations of  
1159 the time-related relevance scores for dominant peak queries and nondominant peak  
1160 queries are defined differently, the parameter  $\beta$  for them should be estimated differ-  
1161 ently. Therefore, we technically have five parameters to estimate:  $K$ ,  $p$ ,  $s$ ,  $\beta_B$  and  $\beta_C$ .  
1162 After the parameters are estimated, we can apply them to test queries. Specifically,  
1163 given the top  $K$  tweets of a test query  $q'$ , we categorize  $q'$  into one of three classes by  
1164 exploring the temporal distribution of the top  $K$  tweets via the estimated  $p$  value and  
1165 the estimated  $s$  value. If  $q'$  is categorized to be a time insensitive query, there is no  
1166 estimated parameter  $\beta$  for  $q'$ ; if  $q'$  is categorized to be a dominant peak query, the esti-  
1167 mated parameter  $\beta_B$  value is used to aggregate the IR scores of the tweets for  $q'$  with  
1168 their  $TRS$ s; if  $q'$  is categorized to be a nondominant peak query, the estimated param-  
1169 eter  $\beta_C$  value is used to aggregate the IR scores of the tweets for  $q'$  with their  $TRS$ s.

Table VII. The Comparison of Three Systems

	TREC 2011					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
System I ( $p = s$ )	0.5062	<b>0.5224</b>	0.4983	0.4116	<b>0.2323</b>	0.4981
System II ( $s = 1$ )	0.5142	<b>0.5224</b>	0.5061	0.4141	0.2273	0.4962
System III	<b>0.5270</b> <sup>†‡</sup>	0.5218	<b>0.5076</b>	<b>0.4357</b>	0.2283	<b>0.5125</b>
	TREC 2012					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
System I ( $p = s$ )	0.3236	0.4362	0.2939	0.2564	0.2554	0.2817
System II ( $s = 1$ )	0.3360	0.4492	0.2933	0.2644	0.2589	0.2858
System III	<b>0.3415</b> <sup>†</sup>	<b>0.4695</b> <sup>†‡</sup>	<b>0.3018</b>	<b>0.2719</b>	<b>0.2738</b> <sup>†‡</sup>	<b>0.2911</b>

Note: † and ‡ indicate statistically significant improvements over System I and System II respectively.

1170 Our proposed system that uses the given estimation method corresponds to the third  
 1171 scenario. Let System III denote the system in the third scenario. System III is com-  
 1172 pared against two systems, each having only one type of time sensitive queries.

1173 System I is obtained by stipulating  $p = s < 1$ , ignoring the restrictions of  $p \leq 0.5$   
 1174 and  $s > p$ . It assumes that if a query does not satisfy Equation (8), it is time sensitive.  
 1175 Because Equation (10) cannot hold when  $p = s$ , all time sensitive queries are assumed  
 1176 to be the dominant peak queries, regardless of the distributions of their top tweets.  
 1177 If a query has multiple peaks, then the highest peak serves as the dominant peak.  
 1178 System I has two parameters  $p(= s)$  and  $\beta$  that can be estimated in a similar way as  
 1179 discussed before. In particular, by assuming  $p = s < 1$ , all training queries can be  
 1180 partitioned into a set of time insensitive queries and a set of dominant peak queries.  
 1181 The combination of a  $p$  value and a  $\beta$  value that yields the largest MAP score for the  
 1182 training queries is utilized to categorize a test query  $q'$ . If  $q'$  is a time sensitive query,  
 1183 then the Laplace-like function is used to calculate the time-related relevance scores  
 1184 for the tweets for  $q'$ , as this is the only type of time sensitive queries for this system.  
 1185 System I corresponds to the proposed system in the first scenario described earlier.

1186 System II is configured by stipulating  $s = 1$ . Because Equation (9) cannot hold when  
 1187  $s = 1$ , all time sensitive queries are assumed to be the nondominant peak queries.  
 1188 Equation (13) is applied to calculate the time-related relevance scores. The two pa-  
 1189 rameters  $p$  and  $\beta$  are estimated using a similar method to that used by System I. For  
 1190 each test query  $q'$ , the combination of a  $p$  value and a  $\beta$  value which yields the largest  
 1191 MAP score for the training queries is applied to  $q'$ . System II corresponds to the pro-  
 1192 posed system used in the second scenario. Table VII presents their performances.

1193 As shown in Table VII, for the set of TREC 2011 queries, compared with System I,  
 1194 System III suffers slight deteriorations in P30 by both relevant criteria. However, Sys-  
 1195 tem III consistently outperforms System I in MAP and in NDCG@30 by both relevant  
 1196 criteria. We also see that System III consistently outperforms System II in almost all  
 1197 measures by both relevant criteria except a negligible deterioration in P30 by the rele-  
 1198 vant criterion. For the set of TREC 2012 queries, System III consistently outperforms  
 1199 System I and System II in all measures by both relevant criteria. These improvements  
 1200 validate our temporal query categorizations.

1201 *6.2.2. The Evaluation of Aggregation Method.* In this section, we evaluate our proposed  
 1202 aggregation method (i.e., System III in Table VII). We compare its performance with  
 1203 two baselines, CombSUM and CombMNZ [Shaw et al. 1994]. In particular, given  
 1204 a tweet  $d$  with an IR score  $IR(d)$  and a time-related relevance score  $TRS(d)$ , the

Table VIII. The Comparison of Various Aggregations

	TREC 2011					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
CombSUM	0.5028	<b>0.5245</b>	0.4951	0.4025	<b>0.2394</b>	0.4917
CombMNZ	0.4909	0.5156	0.4851	0.3931	0.2343	0.4882
Our Aggregation	<b>0.5270</b>	0.5218	<b>0.5076</b>	<b>0.4357</b>	0.2283	<b>0.5125</b>
	TREC 2012					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
CombSUM	0.3358	0.4616	0.2915	0.2495	0.2631	0.2849
CombMNZ	0.3412	0.4638	<b>0.3020</b>	0.2660	0.2690	0.2897
Our Aggregation	<b>0.3415</b>	<b>0.4695</b>	0.3018	<b>0.2719</b>	<b>0.2738</b>	<b>0.2911</b>

1205 CombSUM method calculates an aggregated score for  $d$ ,  $CombSum(d) = IR(d) +$   
 1206  $TRS(d)$ ; the CombMNZ method calculates an aggregated score for  $d$ ,  $CombMNZ(d) =$   
 1207  $CombSum(d) \cdot m_d$ , where  $m_d$  is the number of nonzero scores for  $d$ . Specifically, if  $d$   
 1208 has a nonzero  $IR(d)$  score and a nonzero  $TRS(d)$  score, then  $m_d = 2$ . If a query  $q$   
 1209 is time insensitive, no  $TRS$ s are assigned to the tweets with respect to  $q$ . Table VIII  
 1210 shows the comparisons of the three aggregation methods. For the TREC 2012 queries,  
 1211 our aggregation method consistently outperforms CombSUM in all measures by both  
 1212 relevant criteria. Our method also outperforms CombMNZ in almost all measures by  
 1213 both relevant criteria except a negligible deterioration in NDCG@30 by the relevant  
 1214 criterion. For the TREC 2011 queries, compared with the two baselines, our method  
 1215 suffers marginal deteriorations in P30 by both relevant criteria. But it outperforms the  
 1216 two baselines in all other measures by both relevant criteria. Overall, our aggregation  
 1217 method shows the strongest performance among all three aggregation methods. How-  
 1218 ever, the improvements over CombSUM and CombMNZ by our aggregation method  
 1219 are not statistically significant.

1220 **6.2.3. The Impact of Temporal Information on Retrieval Effectiveness.** We now study the im-  
 1221 pact of incorporating temporal information on retrieval effectiveness. In this experi-  
 1222 ment, we use two baselines. The first baseline is our divide-and-conquer method using  
 1223 the *GreedyMerging* algorithm (i.e., its performance in Table VI), because we want to  
 1224 see whether the inclusion of temporal information can further improve the perfor-  
 1225 mance of this baseline or not. Let BASELINEI denote the first baseline. The second  
 1226 baseline is the algorithm proposed by [Efron and Golovchinsky 2011]. Given a query  
 1227  $q$  with a timestamp  $t$ , this method ranks the tweets published before or on  $t$  by using  
 1228 their temporal information. Specifically, it prefers the recent tweets close to  $t$  to the old  
 1229 tweets and calculates a score  $P(d|q)$  for a tweet  $d$  (publishing at  $t_d$ ) by Equation (16).

$$1230 \quad P(d|q) \propto P(q|d) \cdot r \cdot e^{-r \cdot f(t_d, t)}, \quad (16)$$

1231 where  $r$  is the rate parameter of the exponential distribution.  $P(q|d)$  is an IR score  
 1232 provided by a retrieval model.  $f(t_d, t)$  is the same time representation we adopt in this  
 1233 article. Efron and Golovchinsky [2011] proposed to do the maximum posterior estima-  
 1234 tion for the parameter  $r$  for each  $q$  as follows. Let  $D_q = \{d_1, \dots, d_k\}$  be the top  $k$  tweets  
 1235 for  $q$  by a ranking model. Let  $T_{D_q} = \{t_1, \dots, t_k\}$  be the set of the time representations  
 1236 of the publishing times associated with  $D_q$ . Then  $r_q^{MAP} = \frac{\rho+k-1}{\sigma + \sum_{i=1}^k t_i}$ . This estimation in-  
 1237 volves three parameters,  $k$ ,  $\rho$  and  $\sigma$ . In order to compare our method with this method  
 1238 (denoted by BASELINEII), we use the IR scores of the tweets from the first phase as  
 1239  $P(q|d)$  for BASELINEII. Moreover, we also do the maximum posterior estimation of



Table IX. The Impacts of Temporal Information

	TREC 2011					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
BASELINEI (Divide-and-Conquer Method)	0.5006	0.5143	0.4939	0.4090	<b>0.2283</b>	0.4933
BASELINEII [Efron and Golovchinsky 2011]	0.3958↓	0.3891↓	0.3909↓	0.3391↓	0.1505↓	0.4117↓
Our Method (System III)	<b>0.5270</b> †‡	<b>0.5218</b> ‡	<b>0.5076</b> ‡	<b>0.4357</b> †‡	<b>0.2283</b> ‡	<b>0.5125</b> ‡
	TREC 2012					
	Relevant			Highly Relevant		
	MAP	P30	NDCG@30	MAP	P30	NDCG@30
BASELINEI (Divide-and-Conquer Method)	0.3259	0.4367	0.2966	0.2590	0.2583	0.2852
BASELINEII [Efron and Golovchinsky 2011]	0.2305↓	0.3283↓	0.2082↓	0.1698↓	0.1923↓	0.2011↓
Our Method (System III)	<b>0.3415</b> †‡	<b>0.4695</b> †‡	<b>0.3018</b> ‡	<b>0.2719</b> ‡	<b>0.2738</b> †‡	<b>0.2911</b> ‡

Note: † and ↓ indicate statistically significant improvements and deteriorations over BASELINEI; ‡ indicates statistically significant improvements over BASELINEII.

1240  $r$  for each test query. Efron and Golovchinsky [2011] showed that their suggested pa-  
 1241 rameter values are effective for the proposed maximum posterior estimations across  
 1242 two collections, including a Twitter collection. In our experiments, for each parameter,  
 1243 we try different values for that parameter including their suggested value. For exam-  
 1244 ple, for the parameter  $k$ , we try 10, 20, and 30, where 20 is their suggested value. The  
 1245 method using their suggested values for the three parameters  $k$ ,  $\rho$  and  $\sigma$  achieves the  
 1246 best performance and thus we just report the best performance in this article. Table IX  
 1247 presents the comparison of our method with two different baselines.

1248 As shown in Table IX, compared with BASELINEI, BASELINEII deteriorates in  
 1249 all measures by both relevant criteria for both sets of TREC queries. We are not sur-  
 1250 prised by this results, because BASELINEII always prefers recent tweets to old tweets.  
 1251 Such a recency-preferred strategy does not apply for our queries, as we discussed in  
 1252 Section 1, which is the reason why we propose our three temporal categorizations of  
 1253 queries. For comparing our method with BASELINEI with respect to the set of TREC  
 1254 2011 queries, our method ties with BASELINEI in P30 by the highly relevant criterion  
 1255 but outperforms BASELINEI in all other measures by both relevant criteria. For the  
 1256 set of TREC 2012 queries, it consistently achieves improvements over BASELINEI  
 1257 in all measures using both relevant criteria. This demonstrates that our proposed  
 1258 method using temporality can effectively further improve the retrieval effectiveness  
 1259 of our divide-and-conquer method in the first phase. Our method also consistently and  
 1260 statistically significantly outperforms BASELINEII in all measures by both relevant  
 1261 criteria for the two sets of queries, which demonstrates the effectiveness of our pro-  
 1262 posed temporality usage.

1263 **6.2.4. Result Analysis.** In this section, we conduct an analysis for our utilization of the  
 1264 temporal information of tweets. In particular, we do a query-by-query analysis by com-  
 1265 paring the MAP performance of BASELINEI with that of our method (see Table IX).  
 1266 Figure 6 shows the average precision ( $AP$  for short) changes for the TREC 2011 and  
 1267 2012 queries by both relevant criteria. It displays the changes from the most improved  
 1268 query to the most deteriorated query. According to Figure 6, our usage of temporality  
 1269 improves the average precisions for the majority of the TREC 2011 and 2012 queries.  
 1270 This demonstrates the effectiveness of our proposed method.

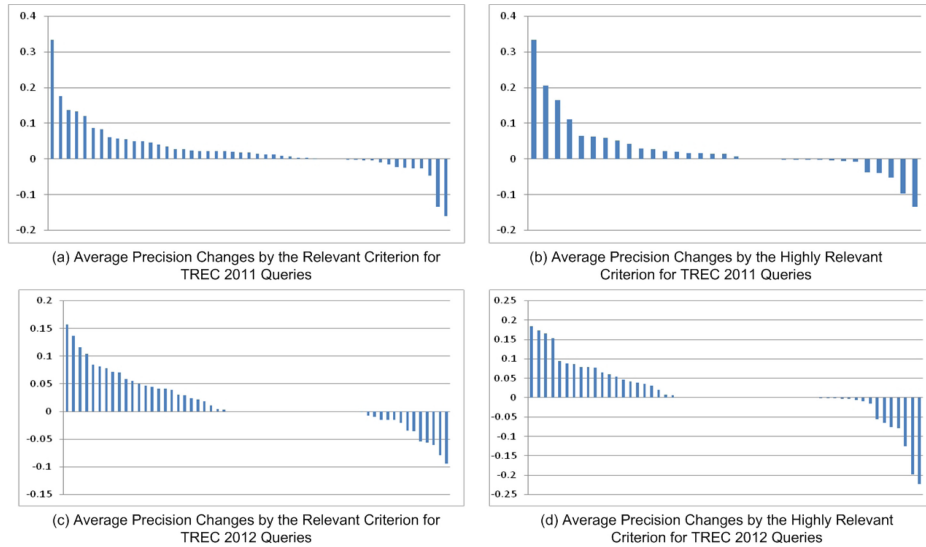


Fig. 6. The average precision changes for temporality.

1271 We provide a deeper analysis to see how many queries are significantly improved  
 1272 or hurt ( $\Delta AP > 0.05$ ) in their  $AP$ s by our temporality usage and present the reasons  
 1273 why their  $AP$ s are improved or deteriorated. For the TREC 2011 queries, 10 (8) queries  
 1274 are significantly improved while 2 (3) queries are significantly hurt according to the  
 1275 (highly) relevant criterion. For the TREC 2012 queries, 11 (13) queries are significantly  
 1276 improved while 5 (7) queries are significantly hurt according to the (highly) relevant  
 1277 criterion. Let us illustrate the reasons with two examples. For example, the query  $q =$   
 1278 “*Aguilera super bowl fail*” is a dominant peak query. 73% of its relevant tweets were  
 1279 published on 2/7/2011. The first phase of our method achieves a very good precision at  
 1280 the top tweets for  $q$ , as 26 out of the top 30 tweets are relevant to  $q$ . Then our method  
 1281 correctly classifies  $q$  as a dominant peak query and predicts its dominant peak is on  
 1282 2/7/2011. This query is significantly improved in its  $AP$ . On the contrary, if the first  
 1283 phase of our method fails to achieve a decent precision by its top tweets with respect to  
 1284 a query  $q$ , then our classification of  $q$  is inaccurate, which may lead to a deterioration  
 1285 in the  $AP$  of  $q$  after we apply our temporal method on  $q$ . For example, for the query  
 1286 “*Michelle Obama fashion*,” we are not surprised to see a significant performance drop  
 1287 for this query again, because the first phase of our method achieves a poor precision  
 1288 at its top tweets, as 10 out of the top 30 tweets are relevant to  $q$ . Our classification for  
 1289 this query and our prediction of the peaks of this query are inaccurate. Overall, the  
 1290 performance of our usage of temporal information depends on an accurate classifica-  
 1291 tion of each query, which in turn depends on how well its top tweets using the first  
 1292 phase of our method approximate its relevant tweets.

### 1293 6.3. Comparison with Related Works

1294 In this section, we compare the performance of our method with those of some related  
 1295 works. TREC 2011 required the retrieved tweets to be ordered in reverse-chronological  
 1296 order [Ounis et al. 2011]. In this experiment, we evaluate the performance of our  
 1297 two-phase method in ranking tweets in reverse-chronological order. Since our method  
 1298 mainly aims at ranking tweets in terms of relevance, we adopt a simple strategy to  
 1299 produce the reverse-chronological ranking of tweets. In particular, we take the top 30

Table X. Comparison of Our Method vs. State-of-the-Art Methods with Respect to the TREC 2011 Queries

	Relevant		Highly Relevant	
	Reverse Chronological Order			
	P30	MAP	P30	MAP
[Liang et al. 2012]	0.4177	0.2365	0.1979	0.2722
[Choi et al. 2012]	0.5068	<b>0.3068</b>	-	-
Our Method	<b>0.5218</b>	0.3018	<b>0.2283</b>	<b>0.3189</b>
	Descending Order of Relevance			
	P30	MAP	P30	MAP
	-	0.3950	-	-
[Amati et al. 2012]	-	0.3950	-	-
Our Method	0.5218	<b>0.5270</b>	0.2282	0.4357
	Descending Order of Relevance Top 100 Tweets			
	P30	MAP	P30	MAP
	-	0.2350	-	-
[Efron et al. 2012]	-	0.2350	-	-
Our Method (Top 100)	0.5218	<b>0.4892</b>	0.2282	0.4262

1300 tweets and rearrange them in reverse order of time. This strategy is the most popular  
1301 strategy adopted by the participants in TREC 2011 [Ounis et al. 2011]. The primary  
1302 evaluation measure is P30 for the reverse-chronological ranking of tweets. Metzler  
1303 and Cai [2011] achieved the best P30 score in TREC 2011 but their results were ob-  
1304 tained in the absence of TREC relevance judgments as training data. So we omit the  
1305 comparison of our results with theirs, because we use TREC relevance judgments as  
1306 training data. We compare our results with other published results with respect to  
1307 the set of TREC 2011 queries. Liang et al. [2012] achieved improvements over the re-  
1308 sults of Metzler and Cai [2011] only by the highly relevant criterion. Moreover, [Choi  
1309 et al. 2012] only reported their performance by the relevant criterion and their results  
1310 outperform the TREC 2011 best results. Some studies [Amati et al. 2012; Efron et al.  
1311 2012] reported their MAP performance by ranking tweets in descending order of rel-  
1312 evance to the TREC 2011 queries, without addressing the requirement of the reverse  
1313 chronological order. We compare our results with these published results in Table X. As  
1314 shown in Table X, our method consistently and significantly outperforms the results  
1315 from Liang et al. [2012] in terms of P30 and MAP by both relevant criteria. Accord-  
1316 ing to the primary evaluation measure P30, our results outperform theirs by 24.9%  
1317 using the relevant criterion and by 15.4% using the highly relevant criterion. Both  
1318 works explore the Web pages whose URLs are embedded in tweets. According to the  
1319 primary measure P30, our results outperform the results from Choi et al. [2012] by the  
1320 relevant criterion and obtains a competitive performance in MAP. For ranking tweets  
1321 in descending order of relevance, our results also significantly outperform the results  
1322 from Amati et al. [2012] and Efron et al. [2012]. Their results were obtained without  
1323 exploring the Web pages linked by tweets while our results use the information from  
1324 those Web pages. Efron et al. [2012] reported their results by only evaluating top 100  
1325 tweets with respect to a given query.

1326 For the set of TREC 2012 queries, we compare our results with the best known re-  
1327 sults reported by the TREC 2012 overview paper [Soboroff et al. 2012]. Unlike TREC  
1328 2011, TREC 2012 required tweets to be ranked in descending order of relevance, in-  
1329 stead of reverse chronological order. Moreover, TREC 2012 only evaluated up to top  
1330 1000 tweets by the highly relevant criterion. The “*hitURLrun3*” run from [Han et al.  
1331 2012] achieved the best P30 and MAP scores [Soboroff et al. 2012]. For the TREC  
1332 2012 participants, the relevance judgments with respect to the TREC 2011 queries are

Table XI. Comparison of Our Method vs. Best Results with Respect to the TREC 2012 Queries

	TREC 2012	
	Highly Relevant	
	MAP	P30
<i>hitURLrun3</i> [Han et al. 2012]	0.2640	0.2701
Our Method	<b>0.2719</b>	<b>0.2738</b>

1333 available as training data, so we compare our corresponding results with the reported  
 1334 best results. Since TREC 2012 required tweets to be ranked in descending order of  
 1335 relevance, MAP is more important than P30. Table XI shows that our results compare  
 1336 favorably with the best results in both measures. Both methods use the Web pages  
 1337 whose links are provided by tweets.

## 1338 7. CONCLUSION AND FUTURE WORK

1339 In this article, we studied the problem of real-time ad-hoc retrieval of tweets intro-  
 1340 duced by TREC 2011. We proposed a two-phase approach to retrieve tweets. Motivated  
 1341 by the observation that tweets have different structures where one type of tweets con-  
 1342 tains just short plain messages (called T-tweets) and the other type of tweets contains  
 1343 short messages with at least one embedded URL (called TU-tweets), we proposed a  
 1344 divide-and-conquer based method for the first phase. Specifically, the method consists  
 1345 of two tweet type-specific rankers and a classifier. We first used the two rankers to ob-  
 1346 tain a ranking of T-tweets and a ranking of TU-tweets. Then we utilized the classifier to  
 1347 determine a preference for every two tweets, one from each type. Finally, we proposed  
 1348 a greedy algorithm to merge the two type-specific rankings into a single ranking for  
 1349 both types of tweets. The merging process takes into consideration all the preferences  
 1350 from the two rankers and the classifier. Experiments showed that our proposed method  
 1351 yields better retrieval effectiveness than the ranker that ranks the two types of tweets  
 1352 simultaneously. We also showed how our method can be made efficient by performing  
 1353 a merging of only the top tweets. In the second phase, we proposed to classify temporal  
 1354 queries by the temporal distributions of their top tweets and calculate the time-related  
 1355 relevance scores of the tweets with respect to different classes of queries accordingly.  
 1356 A ranking of tweets is produced by combining their IR scores from the first phase with  
 1357 their time-related relevance scores. Experimental results demonstrated that the uti-  
 1358 lization of the temporal information can further improve the retrieval effectiveness  
 1359 of the first phase. Our method is also compared favorably with some state-of-the-art  
 1360 methods.

1361 For future work, we plan to investigate whether we can further improve the per-  
 1362 formance of the divide-and-conquer method by the social aspects of tweets. Such in-  
 1363 formation can be found in the JSON version of the TREC Tweets2011 collection. We  
 1364 also plan to study other categorizations of queries, such as cyclic queries and trending  
 1365 queries.

## 1366 REFERENCES

- 1367 Ailon, N., Charikar, M., and Newman, A. 2008. Aggregating inconsistent information: Ranking and cluster-  
 1368 ing. *J. ACM* 55, 5, 23:1–23:27.
- 1369 Amati, G., Amodeo, G., and Gaibisso, C. 2012. Survival analysis for freshness in microblogging search.  
 1370 In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management*  
 1371 *(CIKM'12)*. ACM, New York, 2483–2486.
- 1372 Amodeo, G., Amati, G., and Gambosi, G. 2011. On relevance, time and query expansion. In *Proceedings of*  
 1373 *the 20th ACM International Conference on Information and Knowledge Management (CIKM'11)*. ACM,  
 1374 New York, 1973–1976.

- 1375 Berberich, K., Bedathur, S., Alonso, O., and Weikum, G. 2010. A language modeling approach for tempo-  
 1376 ral information needs. In *Proceedings of the 32nd European conference on Advances in Information*  
 1377 *Retrieval (ECIR'10)*. 13–25.
- 1378 Bian, J., Li, X., Li, F., Zheng, Z., and Zha, H. 2010. Ranking specialization for web search: A divide-and-  
 1379 conquer approach by using topical ranksvm. In *Proceedings of the 19th International Conference on*  
 1380 *World Wide Web (WWW'10)*. 131–140.
- 1381 Choi, J. and Croft, W. B. 2012. Temporal models for microblogs. In *Proceedings of the 21st ACM International*  
 1382 *Conference on Information and Knowledge Management (CIKM'12)*. ACM, New York, 2491–2494.
- 1383 Choi, J., Croft, W. B., and Kim, J. Y. 2012. Quality models for microblog retrieval. In *Proceedings of the*  
 1384 *21st ACM International Conference on Information and Knowledge Management (CIKM'12)*. ACM, New  
 1385 York, 1834–1838.
- 1386 Cohen, W. W., Schapire, R. E., and Singer, Y. 1998. Learning to order things. In *Proceedings of the Conference*  
 1387 *on Advances in Neural Information Processing Systems (NIPS'97)*. 451–457.
- 1388 Dai, N. and Davison, B. D. 2010. Freshness matters: In flowers, food, and web authority. In *Proceedings of*  
 1389 *the 33rd International ACM SIGIR Conference on Research and Development in Information Retrieval*  
 1390 *(SIGIR'10)*. 114–121.
- 1391 Dai, N., Shokouhi, M., and Davison, B. D. 2011. Learning to rank for freshness and relevance. In *Proceedings*  
 1392 *of the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval*  
 1393 *(SIGIR'11)*. 95–104.
- 1394 Dakka, W., Gravano, L., and Ipeirotis, P. G. 2012. Answering general time-sensitive queries. *IEEE Trans.*  
 1395 *Knowl. Data Eng.* 24, 220–235.
- 1396 Dong, A., Chang, Y., Zheng, Z., Mishne, G., Bai, J., Zhang, R., Buchner, K., Liao, C., and Diaz, F. 2010a.  
 1397 Towards recency ranking in web search. In *Proceedings of the 3rd ACM International Conference on*  
 1398 *Web Search and Data Mining (WSDM'10)*. 11–20.
- 1399 Dong, A., Zhang, R., Kolari, P., Bai, J., Diaz, F., Chang, Y., Zheng, Z., and Zha, H. 2010b. Time is of the  
 1400 essence: Improving recency ranking using Twitter data. In *Proceedings of the 19th International Con-*  
 1401 *ference on World Wide Web (WWW'10)*. 331–340.
- 1402 Duan, Y., Jiang, L., Qin, T., Zhou, M., and Shum, H.-Y. 2010. An empirical study on learning to rank of  
 1403 tweets. In *Proceedings of the 23rd International Conference on Computational Linguistics (COLING'10)*.  
 1404 295–303.
- 1405 Efron, M. and Golovchinsky, G. 2011. Estimation methods for ranking recent information. In *Proceedings of*  
 1406 *the 34th International ACM SIGIR Conference on Research and Development in Information Retrieval*  
 1407 *(SIGIR'11)*. 495–504.
- 1408 Efron, M., Organisciak, P., and Fenlon, K. 2012. Improving retrieval of short texts through document expan-  
 1409 sion. In *Proceedings of the 35th International ACM SIGIR Conference on Research and Development in*  
 1410 *Information Retrieval (SIGIR'12)*. ACM, New York, 911–920.
- 1411 Elsas, J. L. and Dumais, S. T. 2010. Leveraging temporal dynamics of document content in relevance ranking.  
 1412 In *Proceedings of the 3rd ACM International Conference on Web Search and Data Mining (WSDM'10)*.  
 1413 1–10.
- 1414 Han, Z., Li, X., Yang, M., Qi, H., Li, S., and Zhao, T. 2012. Hit at trec 2012 microblog track. In *Proceedings*  
 1415 *of Text REtrieval Conference*.
- 1416 Herbrich, R., Graepel, T., and Obermayer, K. 2000. Large margin rank boundaries for ordinal regression. In  
 1417 *Advances in Large Margin Classifiers*, P. J. Bartlett, B. Schölkopf, D. Schuurmans, and A. J. Smola Eds.,  
 1418 115–132.
- 1419 Hüllermeier, E. and Fürnkranz, J. 2010. On predictive accuracy and risk minimization in pairwise label  
 1420 ranking. *J. Comput. Syst. Sci.* 76, 1, 49–62.
- 1421 Joachims, T. 1999. *Advances in Kernel Methods*. 169–184.
- 1422 Joachims, T. 2002. Optimizing search engines using clickthrough data. In *Proceedings of the 8th ACM*  
 1423 *SIGKDD International Conference on Knowledge Discovery and Data Mining (KDD'02)*. 133–142.
- 1424 Jones, R. and Diaz, F. 2007. Temporal profiles of queries. *ACM Trans. Inf. Syst.* 25, 3.
- 1425 Keikha, M., Gerani, S., and Crestani, F. 2011a. Temper: A temporal relevance feedback method. In *Pro-*  
 1426 *ceedings of the 33d European Conference on Advances in Information Retrieval (ECIR'11)*. Springer,  
 1427 436–447.
- 1428 Keikha, M., Gerani, S., and Crestani, F. 2011b. Time-based relevance models. In *Proceedings of the 34th In-*  
 1429 *ternational ACM SIGIR Conference on Research and Development in Information Retrieval (SIGIR'11)*.  
 1430 ACM, New York, 1087–1088.
- 1431 Kulkarni, A., Teevan, J., Svore, K. M., and Dumais, S. T. 2011. Understanding temporal query dynamics.  
 1432 In *Proceedings of the 4th ACM International Conference on Web Search and Data Mining (WSDM'11)*.  
 1433 167–176.

- 1434 Laplace, P.-S. 1774. Mémoire sur la probabilité des causes par les évènements. Mémoires de l'Academie  
1435 Royale des Sciences Présentés par Divers Savan., 621–656.
- 1436 Lee, J. H. 1997. Analyses of multiple evidence combination. In *Proceedings of the 20th International ACM*  
1437 *SIGIR Conference on Research and Development in Information Retrieval (SIGIR'97)*. ACM, New York,  
1438 267–276.
- 1439 Li, X. and Croft, W. B. 2003. Time-based language models. In *Proceedings of the 12th ACM International*  
1440 *Conference on Information and Knowledge Management (CIKM'03)*. 69–475.
- 1441 Liang, F., Qiang, R., and Yang, J. 2012. Exploiting real-time information retrieval in the microblogosphere.  
1442 In *Proceedings of the 12th ACM/IEEE-CS Joint Conference on Digital Libraries (JCDL'12)*. 267–276.
- 1443 Liu, S., Liu, F., Yu, C., and Meng, W. 2004. An effective approach to document retrieval via utilizing wordnet  
1444 and recognizing phrases. In *Proceedings of the 27th International ACM SIGIR Conference on Research*  
1445 *and Development in Information Retrieval (SIGIR'04)*. 266–272.
- 1446 Massoudi, K., Tsagkias, M., de Rijke, M., and Weerkamp, W. 2011. Incorporating query expansion and qual-  
1447 ity indicators in searching microblog posts. In *Proceedings of the 32nd European conference on Advances*  
1448 *in Information Retrieval (ECIR'10)*. Springer, 362–367.
- 1449 McCreadie, R., MacDonald, C., Santos, R., and Ounis, I. 2011. University of glasgow at trec 2011: Ex-  
1450 periments with terrier in crowdsourcing, microblog, and web tracks. In *Proceedings of Text REtrieval*  
1451 *Conference*.
- 1452 Metzler, D. and Cai, C. 2011. Usc/isi at trec 2011: Microblog track (notebook version). In *Proceedings of Text*  
1453 *REtrieval Conference*.
- 1454 Ounis, I., MacDonald, C., Lin, J., and Soboroff, I. 2011. Overview of the trec 2011 microblog track. In *Pro-*  
1455 *ceedings of Text REtrieval Conference*.
- 1456 Rijsbergen, C. J. V. 1979. *Information Retrieval* 2nd Ed. Butterworth-Heinemann, Newton, MA.
- 1457 Robertson, S., Walker, S., Jones, S., Hancock-Beaulieu, M., and Gatford, M. 1996. Okapi at TREC-3. 109–126.
- 1458 Robertson, S., Zaragoza, H., and Taylor, M. 2004. Simple bm25 extension to multiple weighted fields. In  
1459 *Proceedings of the 13th ACM International Conference on Information and Knowledge Management*  
1460 *(CIKM'04)*. 42–49.
- 1461 Shaw, J. A., Fox, E. A., Shaw, J. A., and Fox, E. A. 1994. Combination of multiple searches. In *Proceedings of*  
1462 *the 2nd Text REtrieval Conference (TREC-2)*. 243–252.
- 1463 Soboroff, I., Ounis, I., and Lin, J. 2012. Overview of the trec 2012 microblog track. In *Proceedings of Text*  
1464 *REtrieval Conference*.
- 1465 Zhang, W., Liu, S., Yu, C., Sun, C., Liu, F., and Meng, W. 2007. Recognition and classification of noun phrases  
1466 in queries for effective retrieval. In *Proceedings of the 16th ACM International Conference on Informa-*  
1467 *tion and Knowledge Management (CIKM'07)*. ACM, New York, 711–720.
- 1468 Zhang, X., He, B., Luo, T., and Li, B. 2012. Query-biased learning to rank for real-time twitter search.  
1469 In *Proceedings of the 21st ACM International Conference on Information and Knowledge Management*  
1470 *(CIKM'12)*. ACM, New York, 1915–1919.

1471 Received September 2012; revised April 2013; accepted July 2013