

Real-Time Gaining Subscribers by Messaging

Konstantinos Theocharidis
Singapore Management University
Singapore

Hady W. Lauw
Singapore Management University
Singapore

Abstract

The research on how *brands* can *gain subscribers using content* in their *online pages* (e.g., social networks, e-commerce) is a recent and interesting problem. In particular, in each round k features and m users are searched so as the m users to be notified *once* by brand with *messages* containing k features to subscribe to brand page. However, (k, m) -queries have some noticeable *semantic* and *technical* limitations. In this paper, we propose k -queries to gain subscribers in order to *overcome* the limitations of (k, m) -queries. Specifically, we study the realistic problem of *Real-Time Gaining Subscribers by Messaging* (RGSM). RGSM executes in many rounds. In each round, advertiser searches to message the m users who are most interested to *predefined* k features, as in many real-world scenarios, advertiser *already* knows some k features in the system that are similar to their brand content. The focus on RGSM problem is each k -query to be answered in *real-time* over *many* users having *dynamic* preferences across rounds. We deploy the algorithms SCAN, GROUPY, and BUFFER to solve RGSM on *different* and *realistic* query types. Results show that GROUPY is faster than SCAN and that BUFFER is emphatically superior to others; BUFFER solves RGSM in *real-time*.

CCS Concepts

- Information systems → Information retrieval query processing; Social networking sites; Content match advertising.

Keywords

real-time query processing, subscription gain, messaging

ACM Reference Format:

Konstantinos Theocharidis and Hady W. Lauw. 2018. Real-Time Gaining Subscribers by Messaging. In *Proceedings of Make sure to enter the correct conference title from your rights confirmation email (Conference acronym 'XX')*. ACM, New York, NY, USA, 13 pages. <https://doi.org/XXXXXXX.XXXXXXX>

1 Introduction

Finding the proper *keywords* to stimulate the interest of users is a well-known problem [19, 20, 26, 27, 29]. Yet, in social network and e-commerce platforms (hereafter, *online systems*) such *keywords* are *given* since they correspond to specific *online pages* that inherently represent the *preferences* of users. This happens because users *follow* various pages in online systems to get notified of relative posts as they have an interest in such pages [6, 10]; e.g., users that

Permission to make digital or hard copies of all or part 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 bear this notice and the full citation on the first page. Copyrights for components of this work owned by others than the author(s) must be honored. Abstracting with credit is permitted. To copy otherwise, or republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Request permissions from permissions@acm.org.

Conference acronym 'XX, June 03–05, 2018, Woodstock, NY

© 2018 Copyright held by the owner/author(s). Publication rights licensed to ACM.

ACM ISBN 978-1-4503-XXXX-X/18/06

<https://doi.org/XXXXXXX.XXXXXXX>

follow *Nike* page have *Nike* in their preference set. So, in this paper we consider *given* keywords to form *messages* that motivate users to positively react to them. In particular, we aim to gain the *subscription* of notified users by sending messages containing *limited* keywords that they like most; those keywords are a subset of a feature universe (each *feature* corresponds to a specific *online page*) and the *followers* of an online page are the *subscribers* of that page.

Every service that helps *brands* to *gain subscribers* is *important* as it provides a concrete way for a brand *advertiser* to increase the subscribers of brand page in *online systems* and *offline systems*. In this work, we focus on *online systems* as most brands nowadays maintain social network and e-commerce pages in such systems. Yet, note that even classic websites of brands (we consider them as *offline systems*) can also acquire a portion of the new subscribers gained in online systems since brands can advertise their classic websites in their online system pages. The *benefit* of advertisers possessing subscribers in an offline system is the direct engagement with them via e-mail (e.g., *newsletters*) *without having any restriction* in the number of e-mails they can send. However, the number of messages is *limited per day* in online systems; e.g., an advertiser having an account in the social network VK¹ can only send 20 messages per 12 hours to any non-friend VK user (one message per user). In both cases, *subscribers* are very significant for a brand to *grow* its popularity and revenues since they not only *start* the influence propagation process of brand (in online systems) but they are also the *sole* users that receive its e-mails (in offline systems).

Gaining subscribers using content, introduced in [26], is a recent and interesting service for brand advertisers to increase the subscribers of their pages in online systems. Specifically, the work [26] studies the problem of *Gaining Subscribers by Messaging* (GSM) that pertains to (k, m) -queries, one (k, m) -query per round. Authors in that work search in each round both for k features (online pages; keywords) and m non-subscriber users to send messages containing k features to m users aiming to gain the subscription of m users.

In this paper, we study the same service (*gaining subscribers using content*) but we address it in a much more *practical* and *realistic* way. Particularly, we propose the problem of *Real-Time Gaining Subscribers by Messaging* (RGSM) that pertains to k -queries; a k -query finds in each round the m non-subscriber users having the highest preference to *given* k features (online pages; keywords). So, advertiser messages the m users on *predefined* k features. We claim that k -queries overcome the *practical limitations* (*semantic* and *technical*) of (k, m) -queries as we analytically explain as follows.

¹<https://vk.com/>

²Description Codes are: *woman-looking-at-window-6aqjJ0Zfe5o* for picture 1, *person-on-body-of-water-AlfbUlyzhKz0* for picture 2, *woman-wearing-pink-blouse-5fibEXMHRI* for picture 3, *woman-wearing-brown-hat-Hr_uHRfJmo* for picture 4, *selective-focus-photo-of-woman-in-gray-sleeveless-top-DJNoNHP0K_I* for picture 5, *a-mans-hand-holding-onto-a-watch-on-a-dock-w_kQm6F2h54* for picture 6, *person-in-blue-long-sleeve-shirt-showing-left-hand-DfteCt2Wyvc* for picture 7, *person-wearing-silver-round-analog-watch-with-black-leather-strap-ap5lEDv4Z60* for picture 8, and *round-silver-colored-watch-on-rack-during-sunset-Tb38UzCvKCY* for picture 9.

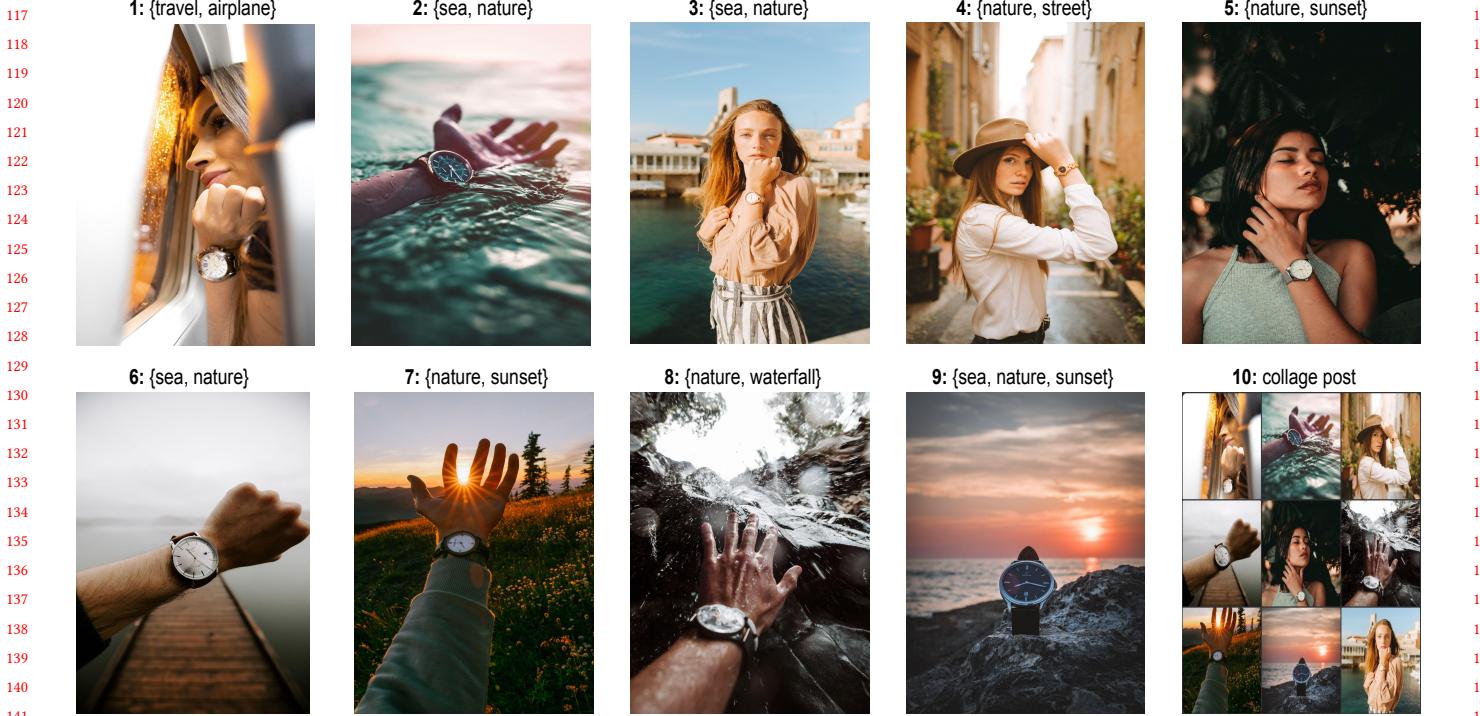


Figure 1: An example depicting 2 cases that mention to the difficulties that advertisers face when the k features are not given. Pictures 1-9 mention to first case while 10th picture mentions to second case. The description of 2 cases lies in the Example 1. All pictures taken from Unsplash by typing in a browser <https://unsplash.com/photos/> followed by the *description code*² for each picture. Each picture (apart from 10) is described by a set of keywords that generally capture the watch-concept of each picture.

1.1 Semantic Contributions of k -queries over (k, m) -queries

In many realistic scenarios, the k features are *already* defined by the brand advertiser as we show in the next cases: (a) *John* is a *watch designer* of a new watch brand and wants to get subscribers to his brand online page. He *already* knows that his watch brand has similar properties with k other known watch brands of online system; (b) *Maria* is a *landscape photographer* and she *already* has a portfolio that is possible to be liked by users that have a high preference on k popular pages of online system relative to photography, architecture, or design; (c) *Giorgio* has *already* found a k -size good content to publish to the *news magazine* page of the company in which he works as an advertiser. However, a replacement of prior k -queries with (k, m) -queries of work in [26] can create significant *inconvenience* and *delay*, as we explain in the following example:

EXAMPLE 1. Consider the previous (a)-scenario relative to *John* with the difference that now he uses a (k, m) -query instead of a k -query to advertise his watch brand. So, *John* solves the problem in [26] and assume that the k features returned to his (k, m) -query for current round are the k pages (suppose $k = 3$) “travel is a pleasure”, “we love sea”, and “nature in our planet”. Then, *John* sees that he must find pictures depicting watches having a concept relative to “travel”, “sea”, and “nature”. Suppose that there are 2 cases for *John*: (i) his watches have specific characteristics, (ii) his watches have a variety of characteristics. In both cases, assume that *John* opts to search in

Unsplash platform (that has a plethora of copyright-free pictures) to find suitable watch pictures to a “travel-sea-nature” concept. Based on current content in Unsplash (**realistic facts**), after 100 mouse-scrolls in his watch-keyword search, *John* finds only the 9 pictures depicted in Figure 1 (pictures 1–9); note that the general character of “travel” keyword allows a flexible utilization of pictures that can be relevant to travel. In the scenario of (i)-case and assuming that all the watches designed by *John* have a round-big-shape, he decides to reject the pictures 3, 4, and 7 while being uncertain of the picture 8 due to its watch blur depiction. So, he creates a collage-post including the pictures 1, 2, 5, 6, (8), and 9. In the scenario of (ii)-case, he creates a collage-post like the one depicted in picture 10. In both cases, the created collage-post constitutes the post of current round. However, the big problem highlighted in both cases, is that *John* needs to spend significant time (100 mouse-scrolls in Unsplash) to find proper watch pictures to create a post under the specific k -concept defined by his (k, m) -query. To stress further that, imagine that in any of next rounds, some of the k pages also relate with a “travel-sea-nature” concept, and as it does not make sense the same picture to be used twice, *John* may spend the double time (e.g., 200 mouse-scrolls) in regards to current round till finding a proper watch picture to include in his collage-post. So, overall, and as rounds evolve, it can be cumbersome for advertisers to efficiently utilize the results of (k, m) -queries where the k features are **arbitrary**. Instead, k -queries are **time-saving** as they depend on existing known content (see Example 3 for more details on k -queries).

233 1.2 Technical Contributions of k -queries over 234 (k, m) -queries

235 The technical difference among k -queries and (k, m) -queries is that
236 the latter searches for 2 parameters (k features and m users) compared
237 to only one parameter searched by the former (m users). This fact
238 challenges the *real-time* applicability of (k, m) -queries, especially
239 when it is tested on realistic cases (big user datasets with
240 users having *dynamic* preferences). Nevertheless, authors in [26]
241 do not study *any of the mentioned criteria* relative to their proposed
242 (k, m) -queries, and that restricts the applicability of those queries
243 and as a result their usefulness towards brand advertisers. By
244 contrast, the k -queries we propose in this work are studied in *real-time*
245 over *many* users with *dynamic* preferences. Based on the arguments
246 we provide in previous section (semantic contributions of k -queries
247 over (k, m) -queries), we claim that k -queries have clearly higher
248 chances to be used more in real world than (k, m) -queries. This
249 induces that *several* advertisers may *concurrently* apply k -queries to
250 online system, and so their *real-time* response is *actually* necessary.
251 We stress that *need* of *real-time* applicability in the next example:
252

253 EXAMPLE 2. Consider the popular online system Instagram³. In-
254 stagram may face k -query cases such as: (A) different advertisers
255 representing different pages apply k -queries in similar time periods;
256 (B) a group of advertisers being responsible for the content of the
257 same page, apply simultaneously k -queries to form their content pub-
258 lishing schedule; (C) a single advertiser of a page applies a k -query
259 but selects only a subset of the returned m users, then repeats another
260 k -query, and so on, till forming a set of m users; (D) any combination
261 among (A)–(C). Assuming that Instagram devotes a single machine
262 with serial execution for the new k -query service, and assuming also
263 a logical number of nearby-time k -query requests, Instagram wants
264 to avoid scenarios where several advertisers in queue wait long enough
265 times leading to an impractical k -query service. For instance, if the
266 query reply time is 2 seconds and have 30 advertisers concurrently
267 apply k -queries, the 30th advertiser will wait almost 1 minute to get
268 their m users, which is time-prohibitive for Instagram.
269

270 1.3 Work Contributions

271 Therefore, based on above arguments we discussed in previous
272 two sections, we contend that the RGSM problem (k -queries) we
273 propose in this paper is not only more *practical* and *realistic* than
274 GSM problem [26] ((k, m) -queries) but is also more *technically*
275 *challenging* due to the actual necessity of real-time applicability.
276

277 The *main contributions* of this paper are the following:

278 **Problem.** We propose and study the RGSM problem, the *first* one
279 that relates with k -queries for gaining subscribers using content in
280 *real-time* over *many* users with *dynamic* preferences. RGSM is an
281 important problem as it provides to brands a realistic and efficient
282 way to increase their subscribers.

283 **Algorithms.** We develop three RGSM algorithms, named SCAN,
284 GROUPY, and BUFFER. We consider SCAN as *baseline*, GROUPY as *advanced*
285 *baseline*, and BUFFER as *state-of-the-art* for the RGSM problem.

286 **Experiments.** We conduct a thorough experimental evaluation on
287 six different and realistic query types (*case studies*) by providing
288 scalable and sensitivity analysis results of RGSM algorithms.

289 ³<https://instagram.com/>

291 2 Problem Statement

292 We define the **RGSM problem** as follows:

293 **PROBLEM 1.** Given an online system depicting a social network or
294 e-commerce platform with $|V|$ users, a feature universe L , a weighted
295 feature set F_v of size $|L|$ capturing the preferences of each user v , a
296 number of m notification messages, each one, containing k features
297 corresponding to an entry of a queries set Q , and a number of rounds
298 n , find in *real-time* in each round t which m non-subscriber users
299 to message so as to maximize the total subscription gain SG over n
300 rounds. A user that receives a message cannot be notified again. The
301 SG of a user v in round t ($v.SG$) is computed as the sum of values
302 (counters) present in round t at the k positions of F_v . The preferences of
303 users are *dynamic*; the counter of any feature in each F_v may increase
304 over rounds. So, the goal of RGSM is the following:

$$305 SG = \operatorname{argmax}_m \sum_{t=1}^n SG_t(k, m) \text{ s.t. } k \text{ fixed, } m \text{ sought} \quad (1)$$

306 The typical scenario in RGSM is that different advertisers *simulta-*
307 *neously* apply k -queries to the online system for finding their
308 m users; so, the response to those queries should be in *real-time*.
309 We consider that each k -query belongs to a *fixed* set Q since the
310 k -feature(keyword)-combinations that interest the advertiser are
311 *concrete*. The reason is that each advertiser uploads content for
312 an online page that has a specific purpose (e.g., relates with *sport*
313 *fashion* or *business economics* or *diet products*, etc.). For clarity, we
314 provide an example about how Q is computed for k -queries, below:

315 EXAMPLE 3. Assume an advertiser that runs a brand page about
316 “*sport fashion*” in an online system. Advertiser observes that the product
317 characteristics of their brand look similar to the ones that are present
318 in posts of Nike, Adidas, Puma, and Reebok pages of the same online
319 system. Yet, advertiser knows that only $k = 3$ features are available
320 to take place in a message. So, in this case, Q includes 3-out-of-4
321 feature combinations. Namely, $Q = \{(Nike, Adidas, Puma), (Nike,$
322 $Adidas, Reebok), (Nike, Puma, Reebok), (Adidas, Puma, Reebok)\}$.
323 Depending on current brand content, advertiser selects the proper
324 k -query $\in Q$ to request online system for finding the proper m users.

325 3 RGSM Algorithms

326 Before proceeding to the description of RGSM algorithms, we
327 present a summary of them that outlines their connection at a
328 high level. SCAN relies its execution on a priority queue without
329 storing any information about updated SG of users over rounds.
330 GROUPY incorporates this information that SCAN misses by storing all
331 non-notified users in a distributed structure and uses a candidate
332 pruning technique based on the maximum *likes* of users per day.
333 BUFFER converts the distributed structure of GROUPY to a linear one
334 that stores a limited number of users with highest SG at that time,
335 exploiting also the candidate pruning technique of GROUPY.
336

337 3.1 SCAN Algorithm

338 Algorithm SCAN presents the SCAN method to solve RGSM.

339 The logic of SCAN is to maintain a priority queue to find the m
340 users with highest SG in each round and marked them as notified
341 for next rounds. In specific, in each round t , advertiser first sets a
342 k -query to the online system (Line 3). Then, for each non-notified
343

349 user v , SCAN computes their SG and may update the priority queue
 350 Q_1 (Lines 4–10). Note that Q_1 is ascending-sorted on $v.SG$ meaning
 351 that its *top* entry denotes the entry with the lowest SG . When all
 352 non-notified users of round t are examined, Q_1 contains the users
 353 to notify in round t and SG increases accordingly (Lines 11–12).

Algorithm SCAN

```

356   Input :  $V, F_v, Q, k, m, n$ 
357   Output :  $SG$  // total subscription gain over  $n$  rounds
358   1  $SG = 0; notified = \emptyset$ ; // notified users over  $n$  rounds
359   2 for  $t = 1, \dots, n$  do
360     Advertiser applies a  $k$ -query  $\in Q$ ; // each  $F_v$  may have now
361     higher values relative to  $k$  features
362      $Q_1 = \emptyset$ ; // pr. queue with entries asc-sorted on  $v.SG$ 
363     for each user  $v \in V \setminus notified$  do
364       Compute  $v.SG$  based on  $k$  positions of  $F_v$ ;
365       if  $|Q_1| < m$  then  $Q_1.push((v, v.SG))$ ;
366       else if  $|Q_1| = m$  and  $v.SG > top.SG$  then
367          $Q_1.pop((top, top.SG))$ ;
368          $Q_1.push((v, v.SG))$ ;
369     for each user  $v \in Q_1$  do
370        $SG = SG + v.SG$ ;  $notified.insert(v)$ ;
371   13 return  $SG$ ;
  
```

3.2 GROUPY Algorithm

Algorithm GROUPY presents the GROUPY method for RGSM.

GROUPY utilizes a *group structure* (GS) in its operation. GS stores for each *query* a *fixed number* (*groupsN*) of *equal-size* (*groupS*) group triples (*gp-triples*). Each *gp-triple* is a map of a *groupId* to triples of the form $(v.userId, v.SG, v.lastUpdate)$ sorted in descending order on $v.SG$ that is the SG of user v to the k keywords of *query*; $v.lastUpdate$ denotes the last round t in which $v.SG$ was computed. Initially, $groupsN * groupS = |V|$, but over rounds less than $|V|$ users are stored in GS due to skipping of *notified* users. Yet, a user v has the same *groupId* for each *query* throughout the whole execution of algorithm. Namely, distribution of users in groups is done once by serially scanning of V and is the same for each *query*.

Along with GS, GROUPY uses a cache (*Qcache*) that comprises the queries stored in GS. However, in real world, *Qcache* cannot contain all the queries of advertiser, so we also implement a cache-replacement policy that replaces the least-frequent *query* present in cache with the new received one; in case of tie, the least-recent *query* is replaced. *Qcache* plays an important role in the operation of GROUPY since in case of replacement, we load all the non-notified users to $GS[query]$, but in opposite case, we do not modify at all the content of $GS[query]$; in both cases, we save time.

Further, GROUPY leverages a *beneficiary value* (*benVal*) in its operation that depicts the maximum number of *likes* (clicking the *like* icon) that any user can yield to the posts published by any feature within the period of a whole day.

The logic of GROUPY is to avoid updating GS as much as possible and in case of updates they should be efficient and meaningful. As meaningful considered the updates relative to users that could be included in the next m -users (if not taking place in current m -ones) and also can lead to an efficient pruning of total candidate users in GS. Namely, GROUPY addresses a k -query of current round t by

maximum-exploit the GS content of previous rounds and efficiently-and-minimum-updating GS for current round.

Algorithm GROUPY

```

407   Input :  $V, F_v, Q, k, m, n$ 
408   Output :  $SG$  // total subscription gain over  $n$  rounds
409   Param. :  $groupsN, groupS, cacheS, benVal$ 
410   1  $SG = 0; notified = \emptyset$ ; // notified users over  $n$  rounds
411   2  $Qcache = \emptyset; GS = \emptyset$ ; // entry: query to gp-triples
412   3 for  $t = 1, \dots, n$  do
413     Advertiser applies a  $k$ -query  $\in Q$ ; // each  $F_v$  may have now
414     higher values relative to  $k$  features
415     Update Qcache and (may) GS for query based on groupsN,
416     groupS, cacheS, and notified;
417      $Q_1 = \emptyset$ ; // priority queue with entries asc-sorted first on
418     v.lastUpdate and second on v.SG
419     Select the first  $m$  triples (one per group) from  $GS[query]$ 
420     whose users  $\notin$  notified to fill  $Q_1$ ;
421     if  $GS[query]$  already updated for round  $t$  then
422       Sequentially examine each group in  $GS[query]$  to update
423        $Q_1$  until  $top.SG \geq v.SG$  for current user  $v$  in group;
424       // examination skips notified users
425       for each user  $v \in Q_1$  do
426          $SG = SG + v.SG$ ;  $notified.insert(v)$ ;
427     else if  $GS[query]$  not updated for round  $t$  then
428        $minQ_2 = \emptyset; maxQ_2 = \emptyset$ ; // priority queues with entries
429       asc-sorted and desc-sorted (respectively) first on v.SG and
430       second on v.lastUpdate
431        $updates = \emptyset$ ; // entry: user to pair of updates
432       Fill  $minQ_2$ ,  $maxQ_2$ , and updates from  $Q_1$ ;
433       Sequentially benVal-examine each group in  $maxQ_2$  to
434       may update  $minQ_2$  and updates; // examination skips
435       notified users
436       Repeat the benVal-examine process of Line 16 for each
437       group  $\in GS[query] \setminus maxQ_2$ ;
438       for each user  $v \in updates \setminus minQ_2$  do
439         Update  $GS[query]$  for pair updates of  $v$ ;
440       for each user  $v \in minQ_2$  do
441          $SG = SG + v.SG$ ;  $notified.insert(v)$ ;
442     return  $SG$ ;
  
```

In more detail, GROUPY handles a k -query in each round t (Line 4) by first updating *Qcache* and GS (Line 5). *Qcache* is always updated while GS gets updated only when *query* is not found in GS (*query* either first time applied or previously replaced by other ones). In case of updating GS, all non-notified users along with their SG and *lastUpdate* values relative to round t are inserted to $GS[query]$.

After updating *Qcache* and GS, GROUPY uses a priority queue Q_1 to store m non-notified users (along with their SG and *lastUpdate* values), each one belonging to a different group in $GS[query]$, who updated their SG more recently than others; in case of tie, users with higher SG are selected (Lines 6–7). The reason of such selection is to quickly find a strong m -candidate-set since distributed-and-recently-updated users can have a good pruning effect to other candidates in GS. In case $GS[query]$ updated, Q_1 has updated values relative to round t (Line 8), else, it has non-updated values relative to previous rounds of t (Line 12).

465 In the former case (Line 8), GROUPY simply examines each user-
 466 triple in a group of GS until a user v is found that cannot modify
 467 the content of Q_1 . Then, GROUPY examines the next group in GS , and
 468 so on (Line 9). When modifications in Q_1 are over, notified users of
 469 round t are found and increase the total SG (Lines 10–11).

470 In the latter case (Line 12), GROUPY moves the content of Q_1 (after
 471 updating $v.SG$ for all triples in Q_1) to the priority queues $minQ_2$ and
 472 $maxQ_2$ (inverse to each other) as also to the $updates$ map (Lines 13–
 473 15). $maxQ_2$ is used for examination (having the user with highest
 474 SG as its *top* entry) while $minQ_2$ is used for modifications as it keeps
 475 the final m -users (*top* entry corresponds to the user with lowest SG).
 476 The map $updates$ stores, for each user that incurs a modification
 477 to $minQ_2$, a pair entry of user previous $\langle SG, lastUpdate \rangle$ values
 478 (present in $GS[query]$ from earlier rounds) and user current $\langle SG,$
 479 $lastUpdate \rangle$ values. If a user v inserted to $updates$ it means that v
 480 is a *strong candidate* (managed to modify $minQ_2$) to take place in
 481 final set of m -users of current round or close next ones.

482 After initializing $minQ_2$, $maxQ_2$, and $updates$, GROUPY proceeds
 483 to an examination process based on $benVal$ to may update the
 484 content of $minQ_2$ and $updates$, by first examining the groups of
 485 $maxQ_2$ (Line 16) and then the other groups of $GS[query]$ (Line 17). The
 486 reason for such examination ordering is that it is expected that
 487 the users with high SG (present in $maxQ_2$) can reduce the total
 488 number of processed users needed over all groups in $GS[query]$ in
 489 order to find the final m users for messaging in round t .

490 Function BenValExamine presents the *benVal*-examine process
 491 that is utilized by GROUPY and BUFFER algorithms.

493 **Function** BenValExamine

494 **Input** : top, v, t, bV, MS (GS or BS), $query$

495 **Output** : $minQ_2, updates$ // modified or unchanged

496 1 $dif_1 = t - v.lastUpdate;$

497 2 $val_1 = v.SG + \lceil dif_1 * bV - ((dif_1 - 1)/2.0) * bV \rceil;$

498 3 **if** $top.SG > val_1$ then // top of $minQ_2$

499 4 $dif_2 = t - \text{last } t'$ where $query$ inserted to MS ;

500 5 $val_2 = v.SG + \lceil dif_2 * bV - ((dif_2 - 1)/2.0) * bV \rceil;$

501 6 **if** $top.SG > val_2$ then // top of $minQ_2$

502 7 **break**; // process next group of $GS[query]$ or complete
 the examination of buffer in $BS[query]$

503 8 **else**

504 9 **continue**; // process next user in same group of
 $GS[query]$ or next triple of buffer in $BS[query]$

505 10 Check whether v modifies $minQ_2$ and $updates$;

506 11 **return** $minQ_2, updates$;

510 Its input relates with top (top entry of $minQ_2$), v (currently ex-
 511 amined user in GS for GROUPY or BS for BUFFER algorithm we ex-
 512 plain later), t (current round), bV (depicts $benVal$ for brevity), MS
 513 (main structure; GS for GROUPY or BS for BUFFER), and $query$ for
 514 the query of round t . Its output relates with any modifications to
 515 $minQ_2$ and $updates$. In short, BenValExamine computes a *first* ben-
 516 eficiary value val_1 depicting a *maximum* possible increased $v.SG$
 517 considering the *delta* among current round t and last round where
 518 $v.SG$ updated (Line 2). In case that val_1 is at least equal to $top.SG$
 519 (Line 3) then the *actual* $v.SG$ may incur modifications to $minQ_2$
 520 and $updates$ (Line 10). Otherwise, BenValExamine computes a *second*
 521 beneficiary value val_2 depicting a *maximum* possible increased
 522

523 $v.SG$ considering the *delta* among current round t and smallest
 524 round where $v'.SG$ updated over any $v' \in MS[query]$ (Line 5); that
 525 smallest round equals to the last round where $query$ inserted to
 526 MS (either as first time applied or after replacement). In case that
 527 val_2 is less than *top* entry in $minQ_2$ then GROUPY proceeds with
 528 the examination of next group in GS (Lines 6–7), as its SG -order
 529 verifies that no user in same group with v can incur any changes
 530 to the final m -users, and so no reason to compute their updated
 531 values for round t . Else, GROUPY proceeds with the examination of
 532 next user in same group of GS (Lines 8–9).

533 When *benVal* examination finishes, GROUPY updates $GS[query]$
 534 only for the users stored in $updates$ but not in $minQ_2$ (Lines 18–19). The
 535 effective operation of BenValExamine entails that *few users* are
 536 found in $updates$ and that saves update time to $GS[query]$. Distribu-
 537 tion of users in groups also helps updates to apply to small segments
 538 (groups) and that decreases further the update time; that is why GS
 539 formatted in a distributed way. Finally, the users kept in $minQ_2$ are
 540 the final m -users to notify in round t and they respectively increase
 541 the total SG (Lines 20–21).

3.3 BUFFER Algorithm

Algorithm BUFFER presents the BUFFER method for RGSM.

Algorithm BUFFER

Input : V, F_v, Q, k, m, n
 Output : SG // total subscription gain over n rounds
 Param. : $bufS, benVal$

1 $SG = 0; notified = \emptyset$; // notified users over n rounds

2 $BS = \emptyset$; // entry: $query$ to triples

3 **for** $t = 1, \dots, n$ **do**

4 Advertiser applies a k -query $\in Q$; // each F_v may have now
 higher values relative to k features

5 $minQ_2 = \emptyset$; // priority queue with entries asc-sorted first on
 $v.SG$ and second on $v.lastUpdate$

6 $updates = \emptyset$; // entry: user to pair of updates

7 **if** $query \in BS$ then // buffer already loaded

8 Fill $minQ_2$ with the first m triples of $BS[query]$ whose
 users \notin $notified$;

9 Insert to $updates$ the m triples of Line 8;

10 Sequentially *benVal-examine* each triple (after m -th
 triple) of $BS[query]$ to may update $minQ_2$ and $updates$;
 // examination skips *notified* users

11 **if** *benVal-examination* of Line 10 completed **then**

12 Repeat Lines 18–21 of GROUPY; // BS for GS

13 **else if** $query \notin BS$ or *benVal-examine* not over **then**

14 Find the $bufS + m$ users with the highest SG as SCAN
 does it for m users in its Lines 4–10;

15 Clear $BS[query]$ and move the $bufS$ users with lowest
 SG to $BS[query]$; // load buffer

16 **for** each user v of rest m users **do** // highest SG

17 $SG = SG + v.SG; notified.insert(v)$;

18 **return** SG ;

We explained the execution of GROUPY in an elaborate way since BUFFER can be considered as a *parameterized* version of GROUPY but it works emphatically faster than GROUPY. In fact, BUFFER was created by addressing the bottlenecks of GROUPY and it constitutes a much simpler implementation despite its impressive performance.

581 After analysis, we found that the overhead of GROUPY is attributed
 582 to two main things. First and foremost, the distributed format of GS
 583 does not usually find quickly an effective enough *m-candidate-set*
 584 and this entails several users to be processed till examination of
 585 all groups in GS to complete (Lines 13–17). A secondary overhead
 586 related to the first is the time required to fill the priority queue
 587 Q_1 (Lines 6–7). The reason is that as rounds evolve, the use of Q_1
 588 to form an initial *m-candidate-set* seems not crucial in algorithm
 589 effectiveness due to the user distribution in GS.

590 The workaround to both prior mentioned bottlenecks of GROUPY
 591 is, instead of using many groups, BUFFER uses just a single group
 592 (*buffer*) to store and modify the content of each *query*. Now, BUFFER (i)
 593 uses a *buffer* (noted by *BS*) as main structure and cache concurrently,
 594 (ii) avoids the problems relative to user distribution, (iii) does not
 595 need an initial *m-candidate-set* to start its operation, and (iv) retains
 596 the small segment that is time-saving to updates application.

597 In more detail, BUFFER uses a *buffer* of size *bufS* to store a *query*
 598 content to *BS[query]* either when *query* is first time applied (*query* \notin
 599 *BS*) or when *benVal*-examine process did not finish (did not meet a
 600 *break* condition; Line 7 in *BenValExamine*) as shown in Line 13. The
 601 *query* content of *BS[query]* comprises triples of the form $\langle v.userId, v.SG, v.lastUpdate \rangle$
 602 sorted in descending order on *v.SG* and carries
 603 the same semantics we mentioned for GROUPY. The *loading* (caching)
 604 of *buffer* is done as follows: First, a set of *bufS+m* users are selected
 605 having the highest *SG* at the moment (round *t*) that advertiser applies
 606 the *k*-query (Line 14). Then, the *bufS* users of them with the
 607 lowest *SG* form the *buffer* (along with their *SG* and *lastUpdate*
 608 values) that init or replaces a prior *buffer* in *BS[query]* (Line 15).

609 The structure, meaning, and purpose of using *minQ₂* and *updates*
 610 is identical to GROUPY, yet note that now the *m-candidate-set* is found
 611 much faster and that contributes to much less users (notification
 612 candidates for round *t*) to be processed. The reason is that instead of
 613 distributed users, users in *buffer* are now the really strong ones (having
 614 the highest *SG* over all non-notified users) and that helps the
 615 *benVal*-examine process to complete much more quickly. However,
 616 as rounds evolve, *benVal*-examine process can be incomplete (not
 617 all possible candidates are considered) over the fixed users present
 618 in *buffer*, and so *buffer* needs to be reloaded with another portion
 619 of users. Yet, in case of reloading, the *m* users of round are already
 620 found and modify total *SG* and *notified* (Lines 16–17) as BUFFER
 621 exploits the sorting needed to form a new *buffer* for *BS[query]*.

622 Last, the examination for finding the *m*-users to notify in round
 623 *t* is much faster due to the absence of Q_1 and $maxQ_2$ (Lines 7–10).

625 4 Experimental Evaluation

626 4.1 Experimental Setup

628 **Implementation.** We wrote code in C++ and we did experiments
 629 on an AMD Ryzen 5 4600U CPU @2.1 GHz machine with 16GB
 630 RAM running Linux Ubuntu 20.04.3 LTS 64-bit.

631 **Data.** We sampled data from the social network VK⁴ that expresses
 632 the Russian version of Facebook. In particular, we sampled 2M users

633 ⁴VK (<https://vk.com/>) represents the Russian version of Facebook in terms of scale,
 634 functionalities, variety of topics, user accounts, brand pages, etc. It has a much more
 635 flexible and unrestricted API (<https://dev.vk.com/en/reference>) than rest social net-
 636 works. Further, according to Wikipedia, VK had been the 16th most visited website in
 637 the world and at the moment it has more than 800M users. All these factors make it
 638 very suitable as a social network data source for research-oriented purposes.

Table 1: The 27 categories of VK in alphabetical order.

VK Categories	
Animals	Auto_motor Beauty_health Celebrity
Cities_countries	Communication_Services
Consumer_Services	Culture_art Education
Entertainment	Finance_and_Insurance Food_recipes
Hobbies	Home_renovation Internet Job_search
Media	Medicine Music Products_stores
Professional_Services	Relationship_family
Restaurants	Social_and_public_organizations
Sport	Tourism_and_Leisure Transportation_Services

Table 2: The 6 query types based on 12 VK categories.

Type	Name	Categories
A	Beautician	Beauty_health Medicine
B	Traveler	Cities_countries Tourism_and_Leisure
C	Educator	Culture_art Education
D	Economist	Finance_and_Insurance Job_search
E	Chef	Food_recipes Restaurants
F	Interior Designer	Home_renovation Products_stores

661 having preferences (F_v) over 540 features of VK that is the sum of
 662 the 20 most popular features (pages) from each category of VK; VK
 663 has in total 27 categories (see Table 1). Initially (round 0), $F_v(f)$
 664 contains the aggregate number of likes that user v yield to any post
 665 published by f (in the page of f) over the years 2010–2017 of VK.
 666 We then use a period of 500 days, where the first day (round 1)
 667 maps to 1st January 2018, to consider any post published in VK by
 668 any f in that period and increases the counters in each $F_v(f)$.

669 **Query Types.** To capture the *realistic* and *dynamic* change of
 670 $F_v(f)$ over rounds, we opted to form 6 different *query types* (as
 671 shown in Table 2) to evaluate the RGSM performance. Specifically,
 672 Q includes 20 (3-out-of-6) *k*-query feature combinations to take
 673 place for each query type, where $k = 3$ and the number 6 maps to
 674 a set comprising the 3 (this value is irrelevant to the value of k)
 675 most popular features from each pair of relative categories. E.g.,
 676 for the query type A (Beautician), the mentioned set contains the
 677 three most popular features in category “Beauty_health” and the
 678 three most popular features in category “Medicine”; 6 elements in
 679 total (6 corresponds to the value 4 in Example 3). This respectively
 680 holds for the rest query types in Table 2. For more details on Q , we
 681 remind readers to see the Example 3.

682 We stress that we consider VK as a big social network and not as
 683 a single dataset. Instead, we consider each query type as a different
 684 dataset; each query type represents a different *case study* since
 685 $F_v(f)$ changes do not correlate at all among different query types.
 686 In other words, each query type enables a distinct RGSM evaluation
 687 besides the fact that all evaluations are done in VK. We believe that
 688 the experiments we provide on 6 separate and unrelated query types
 689 (case studies; datasets) are sufficient to prove our contributions.

690 **Cache.** The cache of GROUPY is denoted by *Qcache* while for BUFFER
 691 the *buffer* (stored per *BS[query]*) operates as cache itself. To be
 692 realistic enough to the presence of massive data in online systems,
 693 we decide a 25% *cache decrease policy* for both algorithms. For
 694 GROUPY, this means that *Qcache* can store only the 75% of the queries
 695 stored in GS. For BUFFER, the effect is that the size of *buffer* is 7.5K

that equals to the group size of GROUPY decreased by 25%; we remind that *buffer* operates as a single group of GROUPY, so the maximum size that *buffer* can have equals to *groupS* (group size of GROUPY).

In more detail, the 25% *cache decrease policy* applies differently in GROUPY and BUFFER due to their different implementation. Since GS stores much more data than BS per query, we have 25% queries decrease in $|Q|$ (set of all k -queries) contained in Qcache that benefit GS; $|Qcache| = 75\%|Q|$ (we denote $|Qcache|$ by *cacheS*). On the contrary, BS does not store much data per query, so $|Q|$ does not affect BUFFER. For that, we apply 25% decrease for the size of *buffer* relative to each query stored in BS ($BS[query]$). As *buffer* continuously updated, it plays the role of cache for BUFFER; $bufS = 75\%groupS$ (we denote the size of *buffer* by *bufS*).

Parameters. The default parameters used in our experiments are as follows: $|V| = \{1M, 1.5M, 2M\}$, $|L| = 540$, $m = 20$, $k = 3$, $n = 500$, $|Q| = 20$, $benVal = 10$. For GROUPY, $groupsN = \{100, 150, 200\}$ respectively for $|V|$, $groupS = 10K$, $cacheS = 15$. For BUFFER, $bufS = 7.5K$. Also, the sequence of k -queries $\in Q$ are selected *uniformly-at-random* (with a different initialization seed) in each experiment; that enhances further the independence among experiments we present in this work (additionally to the notion of aforementioned *case studies*).

Algorithms. Besides the SCAN, GROUPY, and BUFFER algorithms shown in Section 3, there is no other work that can be used to solve RGSM. We analytically justify this argument in Sections 5 and 6. In the former, we explain why our prior GSM solutions [26] are not meaningful for the RGSM problem. In the latter, we additionally discuss other related works that are also not close to RGSM purposes. We remind that SCAN is the *baseline*, GROUPY is the *advanced baseline*, and BUFFER is the *state-of-the-art* method to solve the RGSM problem.

SG Values. Our RGSM algorithms do not actually differ in their achieved SG values⁵. In particular, the optimal SG values per round are derived from SCAN, and the SG values per round of GROUPY and BUFFER are almost identical to the ones of SCAN, as we show in some of our experiments later. Any differences in SG values are either minor (due to the *benVal* parameter used in GROUPY and BUFFER) or derived from SG-ties. Based on that, we stress that the RGSM algorithms in this work are exclusively evaluated on their *efficiency* (running time) as their *efficacy* (SG values) can be considered identical.

4.2 Experimental Approach

We present our experimental approach in the next two sections; the current section operates as preamble to familiarize the readers for the context that will follow.

In the first section (4.3), we evaluate the performance of RGSM algorithms in *scalable datasets* only on *running time* since all methods produce almost identical SG; any differences among found m users over rounds are either minor or derived from SG-ties. Each experiment relates with 500 rounds, so we opted to include only the necessary information (*running time*) for ease of readability.

In the second section (4.4), we present the *actual SG* values of algorithms along with their *running time* per round. This section illustrates the *sensitivity analysis* of algorithms' performance based on different selection values for k and m parameters. Although SG values do not differ among algorithms (as in previous section),

⁵The SG value per round is computed as the sum of separate SG values for each one of the m users selected to be notified in that round.

we show them for completeness and verification reasons; the 100 rounds per experiment here enable better the SG values depiction.

4.3 Performance: Scalability Results

Figure 2 presents the *running time* (depicting how fast a k -query replied including the update time in the algorithm structures) over 500 rounds of SCAN, GROUPY, and BUFFER algorithms in 1M, 1.5M, and 2M users for query type A. Figure 3 does it for query type B; Figure 4 does it for query type C; Figure 5 does it for query type D; Figure 6 does it for query type E; and Figure 7 does it for query type F.

We see that in all cases, BUFFER is significantly faster than SCAN and GROUPY. On average, for the query types A, B, C, and F, BUFFER is roughly at least 2 and 2.5 orders of magnitude faster than GROUPY and SCAN, respectively. Even for the query types D and E (in which F_v is not updated as often as in other query types and so *benVal-examine* process is not similarly effective), BUFFER has several actual real-time responses. When BUFFER performs similarly to SCAN in some rounds, this means that a new *buffer* loaded at such points.

Regarding the performance of GROUPY, we observe that in almost all cases is clearly faster than SCAN (on average) till the middle of rounds (round 250). After that point, it remains better than SCAN (on average) but its superiority declines as rounds evolve. At rounds where GROUPY is less efficient than SCAN, it happens due to cache replacements. The general behavior of GROUPY is attributed to the not so effective *m-candidate-set* due to the distributed format of GS. In first rounds, that set works adequately as *benVal-examine* process can easily find the strong users (with high SG) that separate over others. Yet, as rounds evolve, the deficiencies of a random distribution slow down the responses of GROUPY.

Last, a specific remark on scalable user sizes, is that we see that BUFFER increases even more its speedup over competitors as dataset grows, while GROUPY performs arbitrarily. The reason is that we initialize the random query process with a different seed each time; GROUPY is affected by that due to its randomized distribution effects.

4.4 Performance: Sensitivity Analysis Results

Figure 8 presents the *running time* and SG results of SCAN, GROUPY, and BUFFER algorithms in 1M users for query type A and (a) $k = 2$, $m = 10$, (b) $k = 2$, $m = 20$, (c) $k = 3$, $m = 10$, (d) $k = 3$, $m = 20$. Figure 9 does it for query type B; Figure 10 does it for query type C; Figure 11 does it for query type D; Figure 12 does it for query type E; and Figure 13 does it for query type F. In all figures, the colors *red*, *green*, and *black* (of solid lines) depict the *running time* of SCAN, GROUPY, and BUFFER (as in Figures 2–7), while the colors *blue*, *cyan*, and *purple* (of dashed lines) depict their SG values, respectively. A zoom of **400% scale** is advised to see the SG results analytically.

We remind that when $k = 3$ then $|Q| = 20$ and *cacheS* = 15. Yet, when $k = 2$ then $|Q| = 15$ and *cacheS* = 11 that maps to 25% loss of $|Q|$. The cache size of BUFFER (*bufS*) remains the same in all cases since as stated the cache of BUFFER is irrelevant to $|Q|$ changes.

All results make sense under various aspects. First, the SG values of all algorithms are almost identical as expected (that is why we did not show such results in Figures 2–7). Second, as m grows over a fixed k , both *time* and SG increase since more users involved that require more processing but concurrently yield higher efficacy (SG). Third, the SG values get more compact as k grows due to

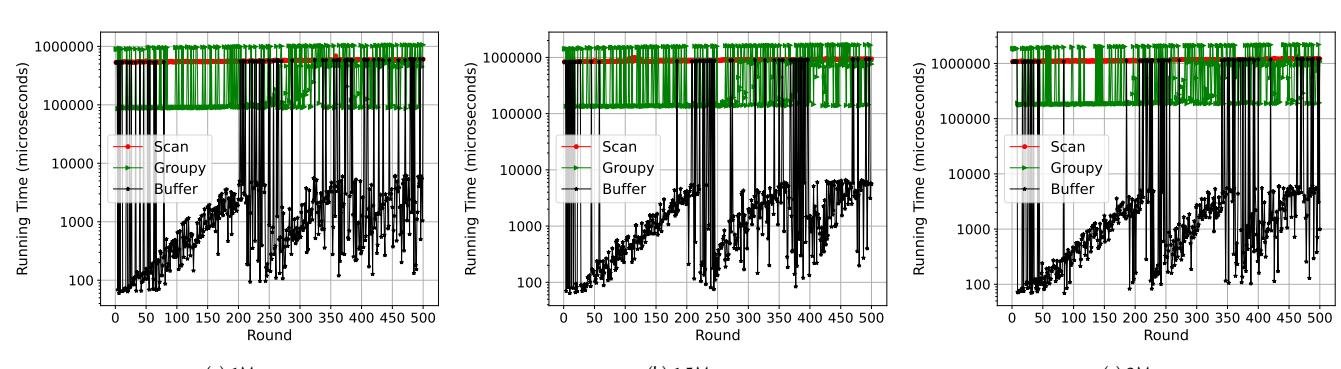


Figure 2: Running time of SCAN, GROUPY, and BUFFER for *Beautician* (type A) over 500 rounds in 1M, 1.5M, and 2M users.

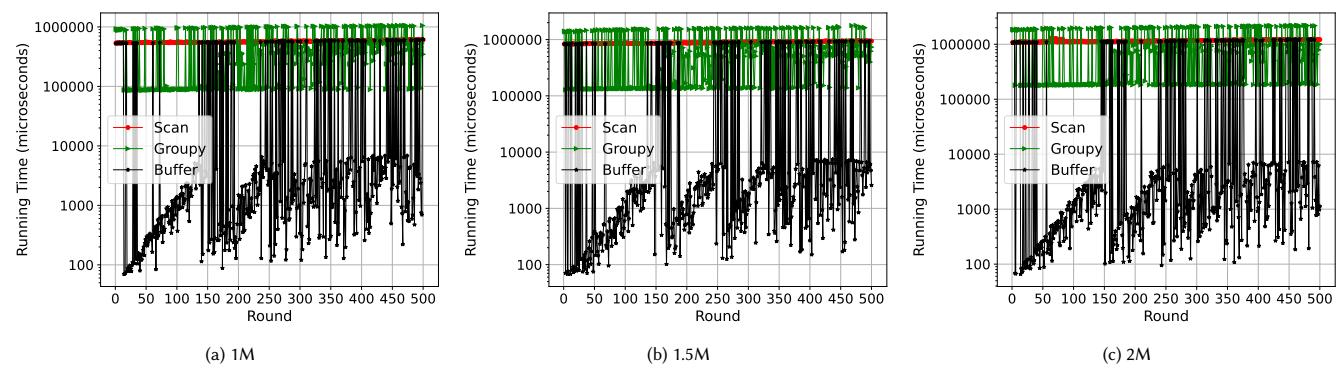


Figure 3: Running time of SCAN, GROUPY, and BUFFER for *Traveler* (type B) over 500 rounds in 1M, 1.5M, and 2M users.

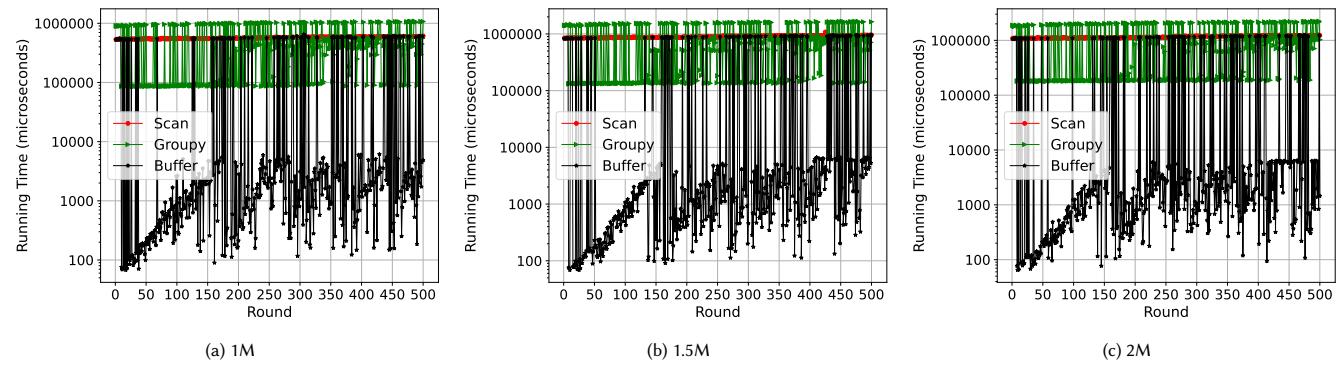


Figure 4: Running time of SCAN, GROUPY, and BUFFER for Educator (type C) over 500 rounds in 1M, 1.5M, and 2M users.

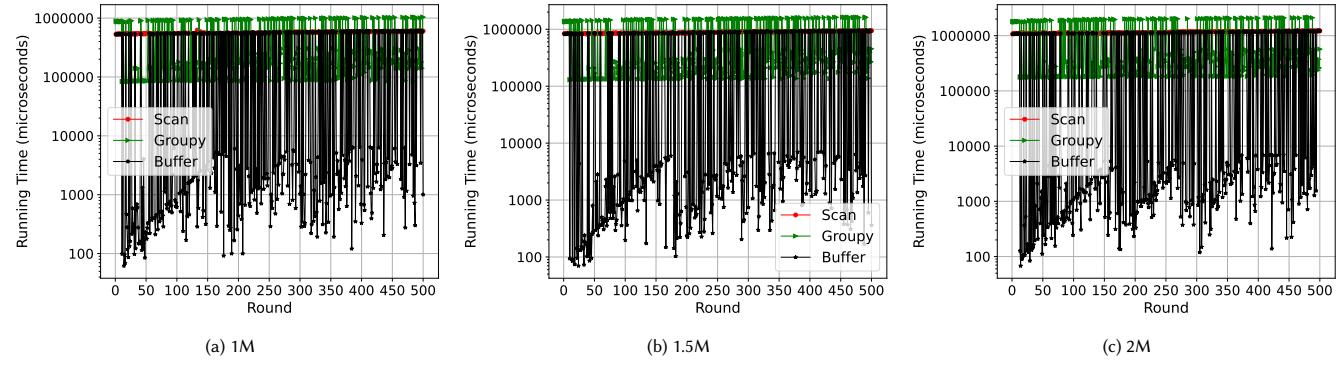


Figure 5: Running time of SCAN, GROUPY, and BUFFER for Economist (type D) over 500 rounds in 1M, 1.5M, and 2M users.

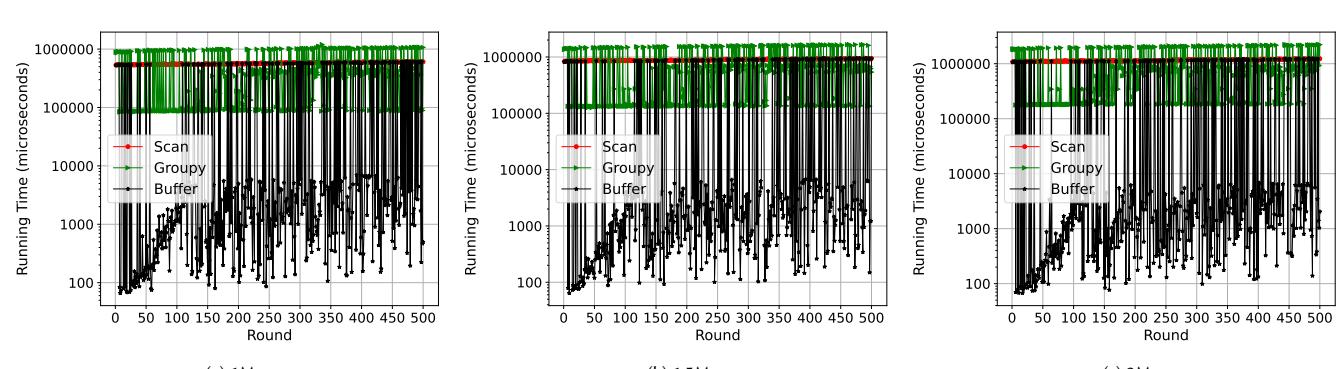


Figure 6: Running time of SCAN, GROUPY, and BUFFER for Chef (type E) over 500 rounds in 1M, 1.5M, and 2M users.

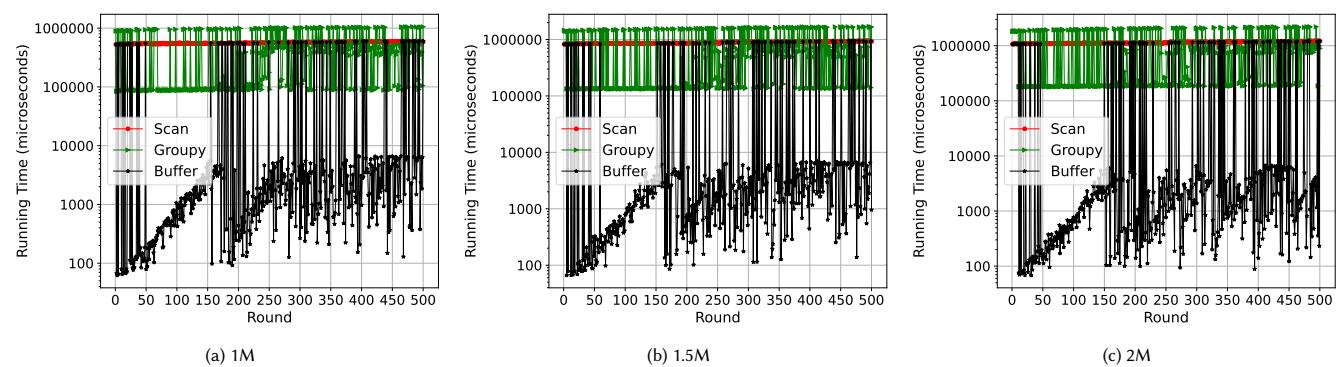


Figure 7: Running time of SCAN, GROUPY, and BUFFER for *Interior Designer* (type F) over 500 rounds in 1M, 1.5M, and 2M users.

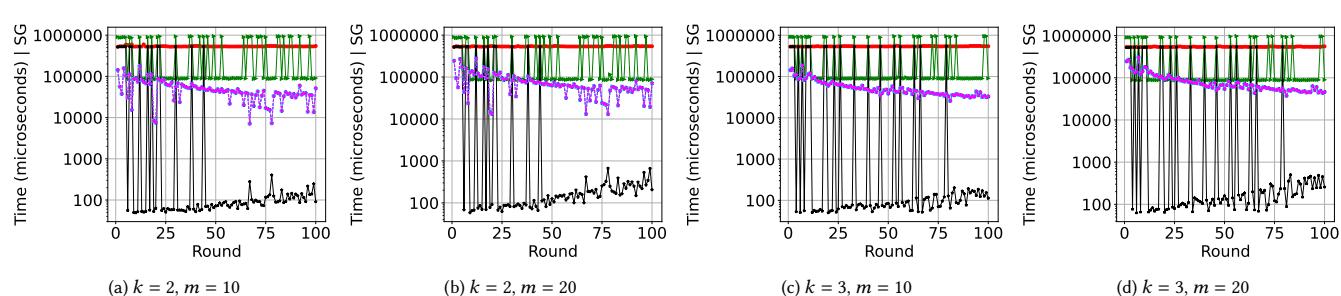


Figure 8: Running time and SG results of SCAN, GROUPY, and BUFFER for Beautician (type A) over 100 rounds in 1M users.

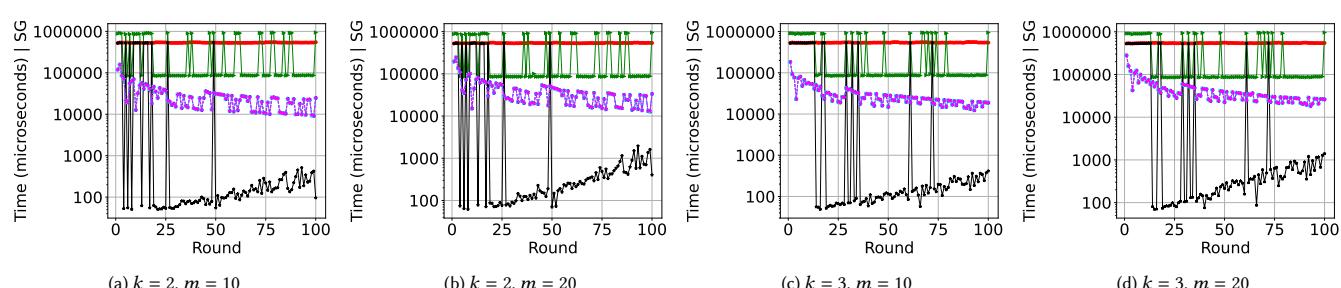
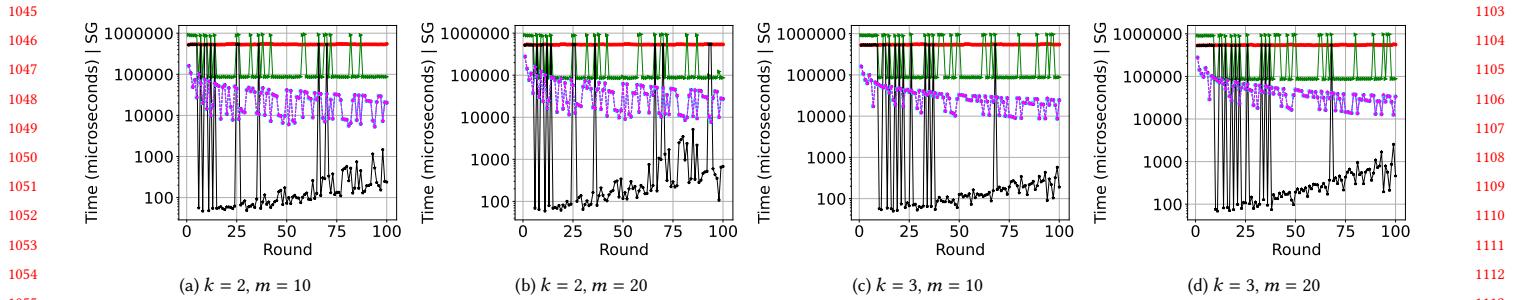
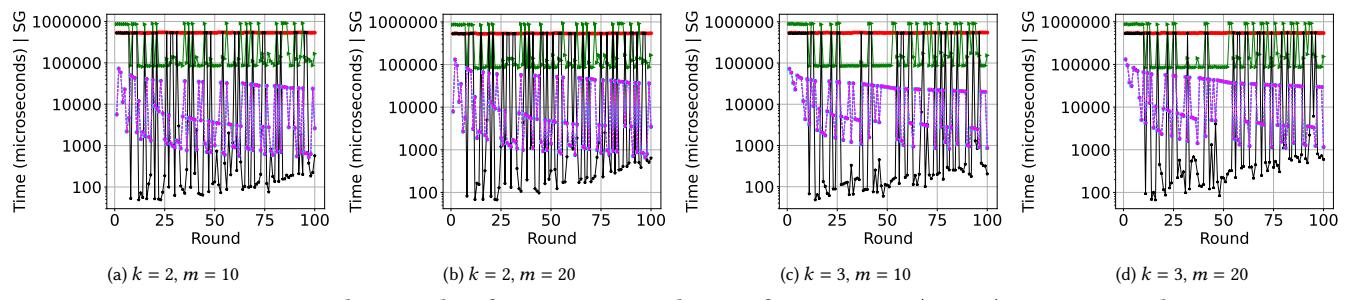
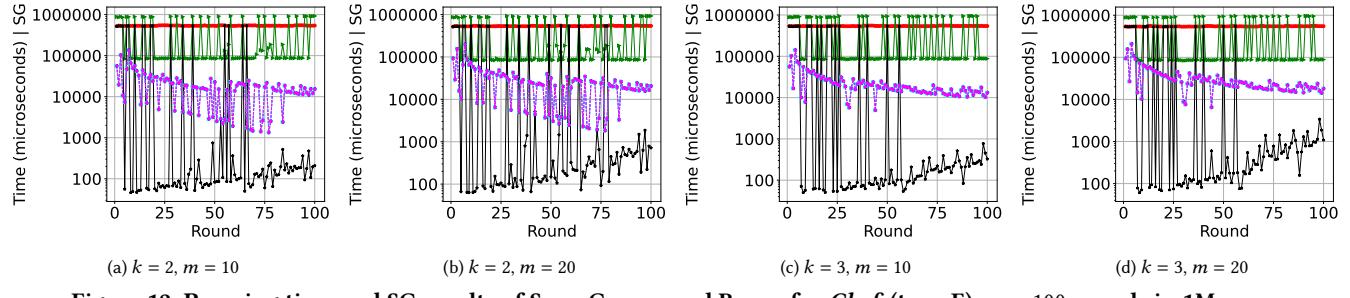
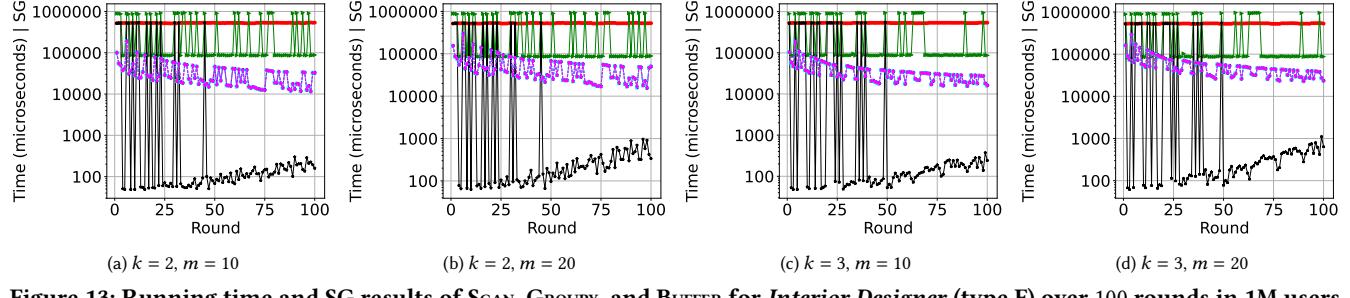


Figure 9: Running time and SG results of SCAN, GROUPY, and BUFFER for Traveler (type B) over 100 rounds in 1M users.

Figure 10: Running time and SG results of SCAN, GROUPY, and BUFFER for *Educator* (type C) over 100 rounds in 1M users.Figure 11: Running time and SG results of SCAN, GROUPY, and BUFFER for *Economist* (type D) over 100 rounds in 1M users.Figure 12: Running time and SG results of SCAN, GROUPY, and BUFFER for *Chef* (type E) over 100 rounds in 1M users.Figure 13: Running time and SG results of SCAN, GROUPY, and BUFFER for *Interior Designer* (type F) over 100 rounds in 1M users.

the more competitive SG value of each notified user. Fourth, more different queries (larger $|Q|$ derived from higher k) usually induce more *buffer* initializations for *BUFFER* but less *buffer* replacements, while usually incur better cache utilization for *GROUPY*; any performance fluctuations depend on the random seed of each query type as we also mentioned for Figures 2–7. Last, all *running time* results are in line with the ones discussed in previous section.

5 GSM Work

Here, we elaborate on basic aspects of our prior work that proposed the mentioned (k, m) -queries to address the *Gaining Subscribers by Messaging* (GSM) problem in [26]. We do that to clarify the difference between the RGSM problem (k -queries) we study in this work and the GSM problem ((k, m) -queries) we studied before [26].

We remind the definition of **GSM problem** as follows:

1161 PROBLEM 2. Given a social network $G = (V, E)$ with $|V|$ users
 1162 and $|E|$ edges, a feature universe L , a weighted feature set F_v of size $|L|$
 1163 capturing the preferences of each user v , a budget k , a limited number
 1164 of m notification messages, a similarity threshold d , and a number of
 1165 rounds n , find in each round t what k content features to publish and
 1166 which m users to notify (via messages containing the k features) so as
 1167 to maximize the cumulative subscription gain SG over n rounds:
 1168

$$1169 \quad SG = \operatorname{argmax}_{k,m} \sum_{t=1}^n SG_t(k, m) \quad \text{s.t. } k \text{ and } m \text{ sought} \quad (2)$$

$$1170$$

1171 Same to RGSM (Section 2), $SG_t(k, m)$ denotes the *subscription gain*
 1172 of m users over k features for round t ; $SG_t(k, m)$ is defined as
 1173 a weighted sum depicting the aggregate preference of m users
 1174 to k features in round t . Note also that although F_v in GSM can
 1175 express *dynamic preferences* as in RGSM, in our GSM work [26]
 1176 we considered only *static preferences* (fixed F_v ; once tuned never
 1177 change). Still, as in RGSM, a user can be notified *only once* in GSM.
 1178

1179 We continue by describing the execution of **SUBSTITUTE** algorithm
 1180 that solves the GSM problem (**SUBSTITUTE** is the best and optimized
 1181 version of all other algorithms studied in [26]). Algorithm **SUBSTITUTE**
 1182 presents its execution for GSM.

Algorithm **SUBSTITUTE**

```

1183   Input :  $V, E, L, F_v, k, m, d, n$ 
1184   Output :  $SG, R$  // total subscription gain and reach (out-degree) of
1185              $m$  users over  $n$  rounds
1186
1187   1  $SG, R = 0$ ;  $notified = \emptyset$ ; // notified users over  $n$  rounds
1188   2 for each  $c_f$  from  $k$ -feature-combinations-of- $L$  do
1189     3    $fC[c_f] = \emptyset$ ; // the  $m$  users with highest  $SG$  for  $c_f$ 
1190     4    $selV[c_f].insert((v, v.SG))$ ; // desc-sorted on  $v.SG$ 
1191
1192   5 for  $t = 1, \dots, n$  do
1193     6    $roundBest = \emptyset$ ; // entry:  $\langle users, features, SG, R \rangle$ 
1194     7   for each  $c_f$  from  $k$ -feature-combinations-of- $L$  do
1195       8     if users in  $notified$  not contained in  $fC[c_f]$  then
1196       9       Check whether  $fC[c_f]$  updates  $roundBest$  using  $d$ 
1197       10      for  $fC[c_f].SG$  comparisons and  $fC[c_f].R$  as second
1198                   filter when needed;
1199       11      continue; // proceed with the next  $c_f$ 
1200       12       $uC = \emptyset$ ; // the  $m$ -size user combinations for  $c_f$ 
1201       13      Enrich  $uC$  the minimum the possible by using
1202                    $SG$ -comparison-optimizations based on  $d$  over a gradual
1203                   processing of users in  $selV[c_f]$ ;
1204       14      Compute  $fC[c_f]$  by considering all entries in  $uC$ ;
1205       15      Repeat Line 9 to possibly update  $roundBest$ ;
1206
1207     16      $SG = SG + roundBest.SG$ ;  $R = R + roundBest.R$ ;
1208     17     for each  $v \in roundBest.users$  do  $notified.insert(v)$ ;
1209     18     Delete from each  $selV[c_f]$  all users in  $notified$ ;
1210
1211   19 return  $SG, R$ ;

```

1211 The main logic of **SUBSTITUTE** is in each round t to compute
 1212 $roundBest$ (includes the (k, m) solution) by constantly updating
 1213 it via the processing of each feature combination c_f . In particular,
 1214 **SUBSTITUTE** first initializes $fC[c_f]$ and computes $selV[c_f]$ for each
 1215 c_f (Lines 2–4); the former intended to store the m -size user combi-
 1216 nation that yields the highest SG for c_f , while the latter stores a
 1217 portion of users in V (based on n value) in a descending order on

1218 their personal SG values ($v.SG$) for c_f . Then, it starts the processing
 1219 of each round for the computation of *roundBest* (Line 5).

1220 There are two ways for each c_f to update *roundBest*. The first
 1221 and time-saving way is the *notified* users of previous rounds to not
 1222 included in $fC[c_f]$, and so $fC[c_f]$ should not be computed again;
 1223 in that case, $fC[c_f]$ uses its existing value to may update *roundBest*
 1224 (Lines 8–9). The second and time-costly way is **SUBSTITUTE** to generate
 1225 all the necessary m -size user combinations for comparison so as to
 1226 find the m users with highest SG for c_f to be stored in $fC[c_f]$, and
 1227 use that $fC[c_f]$ to may update *roundBest* (Lines 11–14).

1228 For $fC[c_f]$ to update *roundBest*, it entails that either $fC[c_f]$
 1229 yields clearly higher SG than *roundBest* based on the similarity
 1230 threshold d , or $fC[c_f].SG$ and *roundBest.SG* are similar (based on
 1231 d) but the reach of $fC[c_f]$ (total out-degree in G of m users in
 1232 $fC[c_f]$) is higher than of *roundBest* ($fC[c_f].R > roundBest.R$).
 1233

1234 Regarding the second case about updating *roundBest* (Lines 11–
 1235 14) that is the main costly part of the algorithm, the SG -comparison-
 1236 optimizations relate with the processing of $selV[c_f]$. Specifically,
 1237 **SUBSTITUTE** gradually examines which lower-ordered users in $selV[c_f]$
 1238 can *substitute* higher-ordered users in $selV[c_f]$ until no more sub-
 1239 stitutions are possible. The meaning of a *substitution* is that a lower-
 1240 ordered user can generate (after substitution) an m -size user combi-
 1241 nation to include in uC (Line 11) that is necessary to be included to
 1242 avoid false misses during the total comparison over uC for finding
 1243 the m users with highest SG for $fC[c_f]$; uC operates as a *knowl-*
 1244 *edge base* that expands with each new substitution. However, to
 1245 improve the efficiency of **SUBSTITUTE**, in [26] we applied additional
 1246 optimizations (related with the notion of *invisibility* in [26]) that
 1247 enable some false misses to take place, and so the optimal SG value
 1248 per round (and overall) to be affected. In short, a user is marked
 1249 *invisible* in $selV[c_f]$ when it is estimated that the user (besides
 1250 the fact that the user can substitute others) cannot yield a strong
 1251 enough m -size user combination to include in uC ; so, that user does
 1252 not expand uC and it is ignored in the rest processing of $selV[c_f]$.

1253 Having described the **SUBSTITUTE** algorithm, we emphasize that
 1254 **SUBSTITUTE** cannot meaningfully apply to solve the RGSM problem
 1255 (we study in this paper), and in case of adaptation, it would be
 1256 similar to **SCAN** algorithm we presented for RGSM in Section 3.
 1257 Specifically, in RGSM the k is fixed, so **SUBSTITUTE** will have a single
 1258 c_f to consider. Also, the RGSM problem does not consider the simi-
 1259 larity threshold parameter d , so all the optimizations of **SUBSTITUTE**
 1260 based on d are not necessary for RGSM. Further, RGSM does not
 1261 use a second filter (like the *reach* used in GSM). Last, RGSM con-
 1262 siders dynamic preferences (F_v change over rounds) instead of
 1263 static preferences used in GSM. After considering all such details,
 1264 for **SUBSTITUTE** to solve RGSM it needs to just recompute $selV[c_f]$
 1265 (where c_f is fixed) in each round (due to dynamic change of F_v)
 1266 and take the m highest-ordered users in $selV[c_f]$ to compute the
 1267 SG of current round. That functionality of $selV[c_f]$ is modeled by
 1268 the priority queue Q_1 used in **SCAN** (Section 3).

1269 Finally, it is important for readers to understand that the techni-
 1270 cal differences among GSM and RGSM problems derive from their
 1271 different applicability. In more detail, the (k, m) -queries used in
 1272 GSM, help an advertiser to *explore* what content would be promis-
 1273 ing when the advertiser does not have concrete ideas about what
 1274 content to publish in their brand page. By the time the advertiser

1275

1277 finds what kind of content works well for their page (after running
 1278 several (k, m) -queries) then the advertiser may stop issuing
 1279 (k, m) -queries and may start issuing k -queries used in RGSM by
 1280 focusing on specific aspects of content (fixed k). We claim that the
 1281 real-time applicability is not so critical for the *exploration* nature
 1282 of (k, m) -queries, and so longer response times to such queries are
 1283 permissible for GSM purposes. In contrast, k -queries are much more
 1284 frequent in real world (as explained in Section 1), and that does
 1285 requests for their real-time response that it is feasible to be achieved
 1286 due to the search of a single parameter (m) in RGSM instead of two
 1287 parameters (k, m) that are sought in GSM. To enhance further the
 1288 chances for k -queries to apply under real-time settings, we do not
 1289 use in RGSM the discussed parameters d and R that are used in GSM.
 1290 To conclude, the more complex modeling around (k, m) -queries
 1291 leaves their real-time response an *open research challenge* whereas
 1292 the more frequent but simplified k -queries allow their real-time
 1293 study *in a feasible way*, and that is the focus of this paper.

1294 6 Related Work

1295 *Gaining subscribers* is valuable to brands for *Influence Maximization*
 1296 (IM) purposes; we describe why in the next paragraph. IM asks
 1297 for the k users who can maximize the influence of a given post
 1298 in a social network, and is a very popular problem that has been
 1299 intensely studied the last two decades. The work in [11] was the
 1300 one that first formally studied the IM problem and subsequent
 1301 works tried to improve its efficacy and efficiency [1–5, 7, 15, 23].
 1302 The work in [16] presents a detailed survey on IM and discusses
 1303 several of its variants. One of such variants is the *Content-Aware*
 1304 *Influence Maximization* (CAIM) problem that we introduced and
 1305 studied in [10]. CAIM can be considered as the inverse variant of
 1306 IM, since instead of seeking for k users to promote a given post,
 1307 in CAIM we search for k features to tune an influential post that
 1308 starts its propagation from the social network page subscribers of
 1309 major-advertised brand in that post. CAIM paved the way for next
 1310 works [12, 24, 26] to achieve content-based influence in a network
 1311 via looking for the proper features to form the relative content.

1312 *Gaining subscribers* is crucial both for IM and CAIM. In the former,
 1313 the k sought influential users to initiate the propagation of a
 1314 given post, are more feasibly and beneficially found in the page sub-
 1315 scribers set of interested-to-get-advertised brand. Having more sub-
 1316 scribers entails that even more influential and loyal initial adopters
 1317 can be found in the subscribers set. In the latter, the post propa-
 1318 gation takes place only in case that there are some subscribers in
 1319 the brand page. All such subscribers are utilized in the propagation
 1320 process as all of them operate as initial adopters; so, the more the
 1321 page subscribers the higher the propagation of post in a network.

1322 In addition, *gaining subscribers* is vital when the *network topology*
 1323 of online systems is partially agnostic [9, 13, 22]; e.g., it is observed
 1324 that many social networks provide limited access to their data. In
 1325 such cases, the subscribers of respective brand are the *sole* users
 1326 for analysis and advertising since any or all connections beyond
 1327 their level may not be available; so, all the influence techniques are
 1328 exclusively deployed and depending on those subscribers which are
 1329 always accessible even if that access requests a manual retrieval.

1330 The first endeavor to *gain subscribers using content* is our prior
 1331 work in [26] that we analytically discussed in Section 5 in order to

1332 clarify the difference among that work and RGSM. Moreover, other
 1333 works relative to RGSM problem are the studies [19, 20, 27, 29];
 1334 [29] implements a system that learns how to extract *keywords* from
 1335 web pages for advertisement targeting; [27] shows that extracted
 1336 *keywords* from different types of meta-data are better suited for
 1337 content-based recommendation than the manually assigned ones;
 1338 [20] analyzes the performance of search engine queries comprising
 1339 *keywords* relative to gender and branded terms; [19] improves the
 1340 performance of e-commerce session-based recommendation by
 1341 generating *keywords* entirely from the click sequence in the current
 1342 session. Yet, all previous works stimulate the interest of users by
 1343 searching the proper keywords to do that whereas RGSM relies on
 1344 given keywords. So, their methods are not suitable for RGSM since
 1345 they are scheduled based on a different objective function.

1346 Our *gaining subscribers* problems (GSM in [26] and RGSM here)
 1347 can also naturally and meaningfully utilize the components de-
 1348 ployed in our recently published works [24, 25]. Specifically, the
 1349 work [24] adaptively finds (based on social network feedback) influ-
 1350 ential features to tune the content of a post in each round, aiming
 1351 to maximize the cumulative influence of posts in a social network
 1352 over all rounds. Features are randomly selected in each round from
 1353 a set of *non-eliminated* features. Instead of random selection, fea-
 1354 tures could be selected by the advertiser of brand page for GSM
 1355 or RGSM purposes without affecting the functionality of learners
 1356 developed in [24]. In addition, the work [25] finds how similar are
 1357 two features; we remind that each feature corresponds to a social
 1358 network page in all our relevant works [10, 24–26], and current
 1359 one. So, highly similar features derived from [25] could enrich the
 1360 search feature space for GSM and the capacity of Q for RGSM.

1361 Finally, it still makes sense to refer some *complementary studies*
 1362 that can stimulate further the interest of m non-subscriber users
 1363 in GSM and RGSM problems; further stimulation relates with any
 1364 actions that additionally boost the intention of notified users (found
 1365 m users in GSM and RGSM) to subscribe to the brand page of respec-
 1366 tive advertiser. Such complementary studies pertain to *personalized*
 1367 *influence maximization* works [8, 14, 17, 18, 21] that maximize the
 1368 influence on *targeted* users based on their preferences, and to *active*
 1369 *friending* works in social networks [28, 30]. Active friending
 1370 is a recommendation strategy that guides the interested user to
 1371 systematically approach their *specific* friending targets; referred
 1372 works maximize the probability of friending targets to accept the
 1373 invitation of interested user (become social friends with user).

1374 7 Conclusion

1375 We proposed how brands can utilize the *existing content* in their
 1376 pages to online systems (social networks, e-commerce) to increase
 1377 their *subscribers* in *real-time*. Specifically, we studied the novel
 1378 RGSM problem and deployed the algorithms SCAN, GROUPY, and
 1379 BUFFER to solve it over many and scalable users with dynamic pref-
 1380 erences. Results showed that, overall, BUFFER *actually* solves RGSM
 1381 in *real-time*, and this fact aligns with the targets of this work.

1382 8 Acknowledgments

1383 This research/project is supported by the National Research Foun-
 1384 dation, Singapore under its AI Singapore Programme (AISG Award
 1385 No: AISG2-RP-2021-020).

References

- [1] Christian Borgs, Michael Brautbar, Jennifer Chayes, and Brendan Lucier. 2014. Maximizing Social Influence in Nearly Optimal Time. In *SODA*. 946–957.
- [2] Wei Chen, Chi Wang, and Yajun Wang. 2010. Scalable Influence Maximization for Preventive Viral Marketing in Large-Scale Social Networks. In *KDD*. 1029–1038.
- [3] Wei Chen, Yifei Yuan, and Li Zhang. 2010. Scalable Influence Maximization in Social Networks under the Linear Threshold Model. In *ICDM*. 88–97.
- [4] Yi-Cheng Chen, Wen-Chih Peng, and Suh-Yin Lee. 2012. Efficient algorithms for influence maximization in social networks. *KAIS* 33, 3 (2012), 577–601.
- [5] Suqi Cheng, Huawei Shen, Junming Huang, Wei Chen, and Xueqi Cheng. 2014. IMRank: Influence Maximization via Finding Self-Consistent Ranking. In *SIGIR*. 475–484.
- [6] Kevin Curran, Sarah Graham, and Christopher Temple. 2011. Advertising on Facebook. *International Journal of E-Business Development* 1, 1 (2011), 26–33.
- [7] Amit Goyal, Francesco Bonchi, and Laks V. S. Lakshmanan. 2011. A Data-Based Approach to Social Influence Maximization. *PVLDB* 5, 1 (2011), 73–84.
- [8] Jing Guo, Peng Zhang, Chuan Zhou, Yanan Cao, and Li Guo. 2013. Personalized Influence Maximization on Social Networks. In *CIKM*. 199–208.
- [9] Thibaut Horel and Yaron Singer. 2015. Scalable Methods for Adaptively Seeding a Social Network. In *WWW*. 441–451.
- [10] Sergei Ivanov, Konstantinos Theodoridis, Manolis Terrovitis, and Panagiotis Karras. 2017. Content Recommendation for Viral Social Influence. In *SIGIR*. 565–574.
- [11] David Kempe, Jon Kleinberg, and Éva Tardos. 2003. Maximizing the Spread of Influence through a Social Network. In *KDD*. 137–146.
- [12] Ansh Khurana, Alvis Logins, and Panagiotis Karras. 2020. Selecting Influential Features by a Learnable Content-Aware Linear Threshold Model. In *CIKM*. 635–644.
- [13] Paul Lagrée, Olivier Cappé, Bogdan Cautis, and Silviu Maniu. 2019. Algorithms for Online Influencer Marketing. *TKDD* 13, 1 (2019), 1–31.
- [14] Jong-Ryul Lee and Chin-Wan Chung. 2014. A query approach for influence maximization on specific users in social networks. *TKDE* 27, 2 (2014), 340–353.
- [15] Jure Leskovec, Andreas Krause, Carlos Guestrin, Christos Faloutsos, Jeanne VanBriesen, and Natalie Glance. 2007. Cost-Effective Outbreak Detection in Networks. In *KDD*. 420–429.
- [16] Yuchen Li, Ju Fan, Yanhao Wang, and Kian-Lee Tan. 2018. Influence Maximization on Social Graphs: A Survey. *TKDE* 30, 10 (2018), 1852–1872.
- [17] Yuchen Li, Dongxiang Zhang, and Kian-Lee Tan. 2015. Real-Time Targeted Influence Maximization for Online Advertisements. *PVLDB* 8, 10 (2015), 1070–1081.
- [18] Qi Liu, Zheng Dong, Chuanren Liu, Xing Xie, Enhong Chen, and Hui Xiong. 2014. Social marketing meets targeted customers: A typical user selection and coverage perspective. In *ICDM*. 350–359.
- [19] Yuanxing Liu, Zhaochun Ren, Wei-Nan Zhang, Wanxiang Che, Ting Liu, and Dawei Yin. 2020. Keywords Generation Improves E-Commerce Session-Based Recommendation. In *WWW*. 1604–1614.
- [20] Partha Mukherjee and Bernard J. Jansen. 2014. Performance analysis of keyword advertising campaign using gender-brand effect of search queries. *Electronic Commerce Research and Applications* 13, 2 (2014), 139–149.
- [21] Hung T. Nguyen, Thang N. Dinh, and My T. Thai. 2016. Cost-aware Targeted Viral Marketing in billion-scale networks. In *INFOCOM*. 1–9.
- [22] Lior Seeman and Yaron Singer. 2013. Adaptive Seeding in Social Networks. In *FOCS*. 459–468.
- [23] Youze Tang, Xiaokui Xiao, and Yanchen Shi. 2014. Influence Maximization: Near-Optimal Time Complexity Meets Practical Efficiency. In *SIGMOD*. 75–86.
- [24] Konstantinos Theodoridis, Panagiotis Karras, Manolis Terrovitis, Spiros Skiadopoulos, and Hady W. Lauw. 2024. Adaptive Content-Aware Influence Maximization via Online Learning to Rank. *TKDD* 18, 6, Article 146 (2024).
- [25] Konstantinos Theodoridis and Hady W. Lauw. 2024. Community Similarity based on User Profile Joins. In *EDBT*. 572–583.
- [26] Konstantinos Theodoridis, Manolis Terrovitis, Spiros Skiadopoulos, and Panagiotis Karras. 2022. A Content Recommendation Policy for Gaining Subscribers. In *SIGIR*. 2501–2506.
- [27] Christian Wartena, Wout Slakhorst, and Martin Wibbels. 2010. Selecting Keywords for Content Based Recommendation. In *CIKM*. 1533–1536.
- [28] De-Nian Yang, Hui-Ju Hung, Wang-Chien Lee, and Wei Chen. 2013. Maximizing acceptance probability for active friending in online social networks. In *KDD*. 713–721.
- [29] Wen-tau Yih, Joshua Goodman, and Vitor R. Carvalho. 2006. Finding Advertising Keywords on Web Pages. In *WWW*. 213–222.
- [30] Yapu Zhang, Jianxiong Guo, Wenguang Yang, and Weili Wu. 2020. Targeted activation probability maximization problem in online social networks. *IEEE Transactions on Network Science and Engineering* 8, 1 (2020), 294–304.

1451
1452
1453
1454
1455
1456
1457
1458
1459
1460
1461
1462
1463
1464
1465
1466
1467
1468
1469
1470
1471
1472
1473
1474
1475
1476
1477
1478
1479
1480
1481
1482
1483
1484
1485
1486
1487
1488
1489
1490
1491
1492
1493
1494
1495
1496
1497
1498
1499
1500
1501
1502
1503
1504
1505
1506
1507
1508