

Top-k Spatial-keyword Publish/Subscribe Over Sliding Window

Xiang Wang · Ying Zhang · Wenjie Zhang · Xuemin Lin · Zengfeng Huang

the date of receipt and acceptance should be inserted later

Abstract With the prevalence of social media and GPS-enabled devices, a massive amount of *geo-textual* data has been generated in a stream fashion, leading to a variety of applications such as location-based recommendation and information dissemination. In this paper, we investigate a novel real-time top- k monitoring problem over sliding window of streaming data; that is, we continuously maintain the *top-k* most relevant *geo-textual messages* (e.g., geo-tagged tweets) for a large number of *spatial-keyword subscriptions* (e.g., registered users interested in *local events*) simultaneously. To provide the most recent information under controllable memory cost, sliding window model is employed on the streaming geo-textual data. To the best of our knowledge, this is the first work to study top- k spatial-keyword publish/subscribe over sliding window. A novel centralized system, called **Skype** (Top-k Spatial-keyword Publish/Subscribe), is proposed in this paper. In **Skype**, to continuously maintain top- k results for massive subscriptions, we devise a novel indexing structure upon subscriptions such that each incoming message can be immediately delivered on its

arrival. To reduce the expensive top- k re-evaluation cost triggered by message expiration, we develop a novel *cost-based k-skyband* technique to reduce the number of re-evaluations in a cost-effective way. Extensive experiments verify the great efficiency and effectiveness of our proposed techniques. Furthermore, to support better scalability and higher throughput, we propose a distributed version of **Skype**, namely, **DSkype**, on top of **Storm**, which is a popular distributed stream processing system. With the help of fine-tuned subscription/message distribution mechanisms, **DSkype** can achieve orders of magnitude speed-up than its centralized version.

1 Introduction

Recently, with the ubiquity of social media and GPS-enabled mobile devices, large volumes of geo-textual data have been generated in a stream fashion, leading to the popularity of *spatial-keyword publish/subscribe system* (e.g., [25,11,37,24,12]) in a variety of applications such as location-based recommendation and social network. In such a system, each individual user can register her interest (e.g., favorite food or sports) and location as a *spatial-keyword subscription*. A stream of *geo-textual messages* (e.g., e-coupon promotion and tweets with location information) continuously generated by publishers (e.g., local business) are rapidly fed to the relevant users.

The spatial-keyword publish/subscribe system has been studied in several existing work (e.g., [25,11,37]). Most of them are geared towards boolean matching, thus making the size of messages received by users unpredictable. This motivates us to study the problem of top- k spatial-keyword publish/subscribe such that only the top- k most relevant messages are presented to

X. Wang (✉) · W. Zhang · X. Lin · Z. Huang
School of Computer Science and Engineering,
The University of New South Wales
Sydney, Australia
E-mail: xiangw@cse.unsw.edu.au

Y. Zhang
CAI, University of Technology,
Sydney, Australia
E-mail: ying.zhang@uts.edu.au

W. Zhang
E-mail: zhangw@cse.unsw.edu.au

X. Lin
E-mail: lxue@cse.unsw.edu.au

Z. Huang
E-mail: huangzengfeng@gmail.com

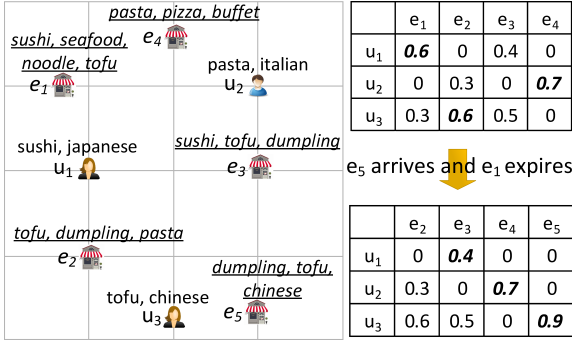


Fig. 1 E-coupon recommendation system

users. Moreover, we adopt the popular sliding window model [4] on geo-textual stream to provide the fresh information under controllable memory usage. In particular, for each subscription, we score a message based on their geo-textual similarity, and the top- k messages are continuously maintained against the update of the sliding window (i.e., message arrival and expiration). Below is a motivating example.

Example 1 Figure 1 shows an example of location-aware e-coupon recommendation system. Three users interested in nearby restaurants are registered with their locations and favorite food, intending to keep an eye on the most relevant e-coupon issued recently. We assume the system only stores the most recent four e-coupons. An e-coupon e will be delivered to a user u if e has the highest score w.r.t. u according to their spatial and textual similarity. Initially, we have four e-coupons, and the top-1 answer of each user is shown in bold in the upper-right table, where the relevance score between user and e-coupon is depicted. When a new e-coupon e_5 arrives and the old e-coupon e_1 expires, the updated results are shown in bottom-right table. Particularly, the top-1 answer of u_1 is replaced by e_3 since e_1 is discarded from the system, while the answer of u_3 is replaced by e_5 , as e_5 is the most relevant to u_3 . The top-1 answer of u_2 remains unchanged.

Challenges. Besides the existing challenges in spatial-keyword query processing [19, 15, 33, 43], our problem presents two new challenges.

The first challenge is to devise an efficient indexing structure for a huge number of subscriptions, such that each message from the high-speed stream can be disseminated immediately on its arrival. The only work that supports top- k spatial-keyword publish/subscribe is proposed by Chen *et al.* [12]. In a nutshell, they first deduce a textual bound for each subscription and then employ DAAT (Document-at-a-time [8]) paradigm to traverse the inverted file built in each spatial node. However, we observe that the continuous top- k monitoring problem is essentially a *threshold-based similarity search problem* from the perspective of message; that

is, a new message will be delivered to a subscription if and only if its score is not less than the current threshold score (e.g., k -th highest score) of the subscription. Consequently, although the DAAT paradigm has been widely used for top- k search (e.g., [17]), it is not suitable to our problem because the advanced threshold-based pruning techniques cannot be naturally integrated under DAAT paradigm.

The second challenge is the top- k re-evaluation problem triggered by frequent message expiration from the sliding window. For example, in Figure 1, the expiration of e_1 invalidates the current top-1 answer (i.e., e_1) of u_1 , and thus the system has to re-compute the new result for u_1 over the sliding window. It is cost-prohibitive to re-evaluate all the affected subscriptions from scratch when a message expires. Some techniques have been proposed to solve this problem (e.g., [41, 29, 7, 31]). Yi *et al.* [41] introduce a $kmax$ strategy, trying to maintain top- k' results, with k' being a value between k and $kmax$, rather than buffering the exact top- k results. Later, Mouratidis *et al.* [29] notice that $kmax$ ignores the dominance relationship between messages, and propose a novel idea to convert top- k maintenance into *partial k -skyband* maintenance to reduce the number of re-evaluations. Nevertheless, they simply use the k -th score of a continuous query (i.e., subscription in our paper) as the threshold of its k -skyband without theoretical underpinnings, which may result in poor performance in practice.

On the other hand, the limited computational resources (e.g., CPU, memory) in a single machine often become the bottleneck when we increase the scale of real-life applications, where millions of active users need to be maintained simultaneously. To alleviate this issue, we extend *Skype* on top of *Storm*¹, an open-source distributed real-time in-memory processing system, to leverage parallel processing such that high throughput can be achieved. *Storm* itself is intrinsically designed to solve real-time stream processing tasks, which therefore best suits our top- k publish/subscribe problem. The main challenge here lies in how to partition and distribute subscriptions and messages such that workload balance and high throughput can be achieved at a small communication cost.

In this paper, we propose a novel centralized system, i.e., *Skype*, to efficiently support top- k **S**patial-**K**eyword **P**ublish/**S**ubscribe over sliding window. Two key modules, *message dissemination module* and *top- k re-evaluation module*, are designed to address the above challenges. Specifically, the message dissemination module aims to rapidly deliver each arriving message to its affected subscriptions on its arrival. We devise efficient subscription indexing techniques

¹ Apache storm project. <http://storm.apache.org/>

which carefully integrate both spatial and textual information. Following the TAAT (Term-at-a-time [9]) paradigm, we significantly reduce the number of non-promising subscriptions for the incoming message by utilizing a variety of spatial and textual pruning techniques. On the other hand, the top- k re-evaluation module is designed to refill the top- k results of subscriptions when their results expire. To alleviate frequent re-evaluations, we develop a novel *cost-based k -skyband* technique which carefully selects the messages to be buffered based on a threshold value determined by a cost model, considering both top- k re-evaluation cost and k -skyband maintenance cost. In addition, to speed-up real-time processing, we follow most of the existing publish/subscribe systems (e.g., [25, 37, 24, 12]) to implement all our indexes in main memory.

To support better scalability beyond *Skype*, we pioneer a novel distributed real-time processing system, namely, *DSkype*, which is a distributed version of *Skype* deployed on top of *Storm*. We propose four different distribution mechanisms, i.e., hashing-based, location-based, keyword-based and prefix-based, to distribute subscriptions and messages to relevant components. Among them, prefix-based technique yields the best overall performance in terms of both throughput and communication cost. For example, it can process nearly 1300 messages per second over 5M subscriptions on a small-size cluster.

Contributions. Our principal contributions are summarized as follows:

- We propose a novel framework, called *Skype*, which continuously maintains top- k geo-textual messages for a large number of subscriptions over sliding window model. To the best of our knowledge, this is the first work to integrate sliding window model into spatial-keyword publish/subscribe system. (Section 4)
- For message dissemination module, we propose both *individual pruning technique* and *group pruning technique* to significantly improve the dissemination efficiency following the TAAT paradigm. (Section 5)
- For top- k re-evaluation module, a novel *cost-based k -skyband* method is developed to determine the best threshold value with in-depth theoretical analysis. It is worth mentioning that our technique is a general approach which can be applied to other continuous top- k problems over sliding window. (Section 6)
- We extend *Skype* on top of *Storm*, a distributed real-time processing environment. By introducing to *Storm* a distribution layer which employs several efficient distribution mechanisms, the distributed version can achieve high throughput with better scalability. As far as we know, this is the first work which

extends top- k publish/subscribe system on top of *Storm*. (Section 7)

- We conduct extensive experiments to verify the efficiency and effectiveness of both *Skype* and its distributed version *DSkype*. It turns out that *Skype* usually achieves up to orders of magnitude improvement compared to its competitors, while *DSkype* achieves further improvement over *Skype* with better scalability and large margin. (Section 8)

2 Related Work

2.1 Spatial-keyword Search

Spatial-keyword search has been widely studied in literatures. It aims to retrieve a set of geo-textual objects based on boolean matching (e.g., [45, 23, 19]) or score function (e.g., [15, 33, 14, 43]) by combining both spatial index (e.g., R-Tree, Quadtree) and textual index (e.g., inverted file). A nice summary of spatial-keyword query processing is available in [13]. Several extensions based on spatial-keyword processing have also been investigated, such as moving spatial-keyword query [21], collective spatial-keyword query [22] and reverse spatial-keyword query [26]. Note that a spatial-keyword search is an ad-hoc/snapshot query (i.e., user-initiated model) while our problem focuses on continuous query (i.e., server-initiated model).

2.2 Publish/Subscribe System

Users register their interest as long-running queries in a publish/subscribe system, and streaming publications are delivered to relevant users whose interests are satisfied. Nevertheless, most of the existing publish/subscribe systems (e.g., [38, 34, 42, 35]) do not consider spatial information. Recently, spatial-keyword publish/subscribe system has been studied in a line of work (e.g., [25, 11, 37, 24, 12]). Among them, [25, 11, 37] study the boolean matching problem while [24] studies the similarity search problem, where each subscription has a *pre-given* threshold. These work are inherently different from ours, and it is non-trivial to extend their techniques to support top- k monitoring.

The CIQ index proposed by Chen *et al.* [12] is the only close work that supports top- k spatial-keyword publish/subscribe (shown in Figure 2). In CIQ, a Quadtree is used to partition the whole space. Each subscription is assigned to a number of covering cells, forming a disjoint partition of the entire space. In Figure 2, we assume all the subscriptions have the same cell covering, i.e., from c_1 to c_7 . A textual bound (e.g., MinT) is precomputed for each subscription w.r.t. each

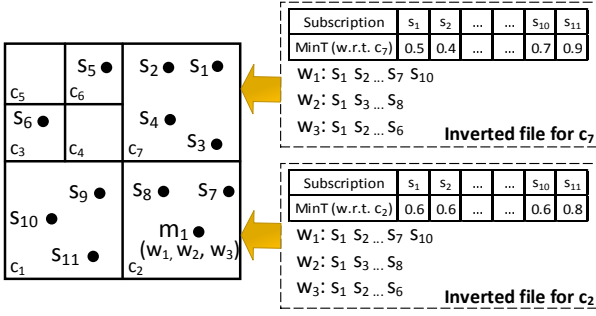


Fig. 2 Example of CIQ index

assigned cell, as shown in the tables where the textual bounds w.r.t. c_2 and c_7 are displayed. An inverted file ordered by subscription id is built to organize the subscriptions assigned to each cell. For a new message (e.g., m_1), CIQ traverses all the inverted files with corresponding cells penetrated by message location (e.g., c_2) in DAAT paradigm, and finds all the subscriptions with textual similarity higher than the precomputed bound as candidates, which are then verified to get final results. However, we notice that DAAT paradigm employed in CIQ cannot integrate some advanced techniques for threshold-based similarity search, given that the nature of our problem is a threshold-based search problem. Contrary to CIQ, our indexing structure is designed for the TAAT paradigm, combined with advanced techniques for threshold-based pruning, thus enabling us to exclude a significant number of subscriptions. Moreover, CIQ indexes each subscription into multiple cells, taking advantage of precomputed spatial bound. However, the gain is limited since the number of covering cells for each subscription cannot be too large; otherwise, it would lead to extremely high memory cost. Thus, we turn to an *on-the-fly* spatial bound computation strategy, where each subscription is assigned to a single cell with finer spatial granularity. Finally, we remark that CIQ integrates a time decay function rather than a sliding window, which, in the worst case, may overwhelm the limited memory.

2.3 Top-k Maintenance Over Sliding Window

One critical problem for top- k maintenance over sliding window is that, when an old element (i.e., message in this paper) expires, we have to recompute the top- k results for the affected continuous queries (i.e., subscriptions in this paper), which is cost-expensive if we simply re-evaluate from scratch. On the flip side, it is also infeasible to buffer all elements and their scores for each individual query to avoid top- k re-evaluation. Several techniques are proposed aiming to identify a trade-off between the number of re-evaluations and the buffer size. In [41], Yi *et al.* introduce a $kmax$ approach.

Rather than maintain exact top- k results, they continuously maintain top- k' results where k' is between k and a parameter $kmax$. However, followed by observation from Mouratidis *et al.* [29], $kmax$ may contain redundant elements due to the overlook of dominance relationship. Thus, Mouratidis *et al.* propose a k -skyband based algorithm to remove redundancy. Since it is very expensive to maintain the *full* k -skyband for each individual query, they only keep elements with scores not lower than the k -th highest score determined by the most recent top- k re-evaluation. We observe that this setting is rather ad-hoc and thus may result in unsatisfactory performance in practice. Böhm *et al.* [7] utilize a delay buffer to avoid inserting the newly-arriving objects with low scores into the k -skyband. However, since each object has to probe query index twice during its life time, their method is not suitable to our problem given the large number of registered queries (i.e., subscriptions). Pripuzic *et al.* [31] propose a probabilistic k -skyband method to drop the data which is unlikely to become top- k results in order to save space and improve efficiency. However, their technique may discard some top- k elements due to its probabilistic nature. In this paper, we propose a novel *cost-based* k -skyband technique to carefully determine the size of k -skyband buffer based on a cost model.

2.4 Distributed Spatial Query Processing

There are a bunch of work studying spatial query processing by utilizing distributed system. Nishimura *et al.* [30] extend HBase² to support multi-dimensional index. Aji *et al.* [1] propose Hadoop-GIS, a distributed data warehouse infrastructure built on top of Hadoop, which provides functionality of spatial data analytics. Later, Eldawy *et al.* [18] develop SpatailHadoop, a full-fledged system which supports various spatial queries by integrating spatial-awareness in each Hadoop layer. Aly *et al.* present an adaptive mechanism on top of Hadoop to partition large-scale spatial data for efficient query processing [2]. Xie *et al.* [40] introduce a system called Simba to provide efficient in-memory spatial analytics by extending Spark SQL engine. All the work above focus on some fundamental spatial queries, such as range query and kNN query, which is inherently different from our top- k spatial-keyword publish/subscribe problem. A very relevant work, called Tornado, which also supports spatial-keyword stream processing on Storm, appears in [27]. Tornado is a general spatial-keyword stream processing system to support both snapshot and continuous queries. However, their main focus is not on the index construction over sub-

² Apache HBase project. <https://hbase.apache.org/>

scription queries, which nevertheless is the main contribution of our paper. Besides, they cannot support the top- k spatial-keyword subscription queries as ours.

On the other hand, many stream processing systems, such as Spark Streaming³, Samza⁴ and Storm, have been developed to support efficient processing of real-time data. Most of them are featured with open-source, low-latency, distributed, scalable and fault-tolerant characteristics. A nice comparison between different stream processing systems can be found in [32]. We choose Storm here mainly because of its simplicity, efficiency, well-documented APIs and very active community⁵. To the best of our knowledge, our work is the first one to support top- k spatial-keyword publish/subscribe in a distributed environment.

3 Preliminary

In this section, we formally present some concepts which are used throughout this paper.

Definition 1 (Geo-textual Message) A geo-textual message is defined as $m = (\psi, \rho, t)$, where $m.\psi$ is a collection of keywords from a vocabulary \mathcal{V} , $m.\rho$ is a point location, and $m.t$ is the arrival time.

Definition 2 (Spatial-keyword Subscription) A spatial-keyword subscription is denoted as $s = (\psi, \rho, k, \alpha)$, where $s.\psi$ is a set of keywords, $s.\rho$ is a point location, $s.k$ is the number of messages that s is willing to receive and $s.\alpha$ is the preference parameter used in the score function.

To buffer the most recent data from geo-textual stream, we adopt a count-based sliding window defined as follows.

Definition 3 (Sliding Window) Given a stream of geo-textual messages arriving in time order, the sliding window \mathcal{W} over the stream with size $|\mathcal{W}|$ consists of most recent $|\mathcal{W}|$ geo-textual messages.

In the following of the paper, we abbreviate *geo-textual message* and *spatial-keyword subscription* as *message* (denoted as m) and *subscription* (denoted as s) respectively if there is no ambiguity. We assume that the keywords in vocabulary \mathcal{V} , as well as the keywords in subscription and message, are sorted in increasing order of their term frequencies. Note that sorting keywords in increasing order of frequency is a widely-adopted heuristic to speed up similarity search [10, 5, 39]. The i -th keyword in s is denoted as

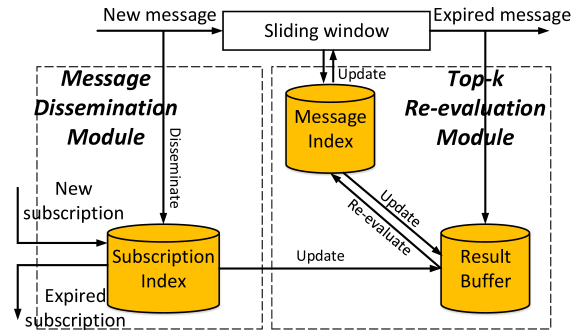


Fig. 3 Framework of Skype

$s.\psi[i]$, and we use $s.\psi[i : j]$ to denote a subset of $s.\psi$, i.e., $\cup_{i \leq k \leq j} \{s.\psi[k]\}$. Particularly, $s.\psi[i :]$ denotes $\cup_{i \leq k \leq |s.\psi|} \{s.\psi[k]\}$. Message m follows the similar notations.

Score function. To measure the *relevance* between a subscription s and a message m , we employ a score function defined as follows:

$$\text{Score}(s, m) = s.\alpha \cdot \text{SSim}(s.\rho, m.\rho) + (1 - s.\alpha) \cdot \text{TSim}(s.\psi, m.\psi) \quad (1)$$

where $\text{SSim}(s.\rho, m.\rho)$ is the spatial proximity and $\text{TSim}(s.\psi, m.\psi)$ is the textual relevance between s and m . Thus, a subscriber can receive messages which are not only close to her location but also fulfil her interest. Meanwhile, the parameter α can be adjusted by subscribers to best satisfy their diverse preferences.

To compute spatial proximity, we utilize Euclidean distance as $\text{SSim}(s.\rho, m.\rho) = 1 - \frac{\text{EDist}(s.\rho, m.\rho)}{\text{MaxDist}}$, where $\text{EDist}(s.\rho, m.\rho)$ is the Euclidean distance between s and m , and MaxDist is maximum distance in the space.

For textual similarity, we employ the well-known *cosine similarity* [28] as $\text{TSim}(s.\psi, m.\psi) = \sum_{w \in s.\psi \cap m.\psi} \text{wt}(s.w) \cdot \text{wt}(m.w)$, where $\text{wt}(s.w)$ and $\text{wt}(m.w)$ are *tf-idf* weights of keyword w in s and m respectively. Note that the weighting vectors of both s and m are normalized to unit length. Also, same as [12], to guarantee the top- k results are *textual-relevant*, a message must contain at least one common keyword with a subscription to become its top- k results.

Problem statement. Given a massive number of spatial-keyword subscriptions and a geo-textual stream, we aim to continuously monitor top- k results for all the subscriptions against the stream over a sliding window \mathcal{W} in real time.

4 Framework

Figure 3 shows the framework of Skype (Top- k Spatial-keyword Publish/Subscribe). We assume our system already has some registered subscriptions. An arriving

³ Apache spark project. <http://spark.apache.org/streaming/>

⁴ Apache samza project. <http://samza.apache.org/>

⁵ <https://github.com/apache/storm>

message will be processed by **message dissemination module**, where a *subscription index* is built to find all the affected subscriptions and update their top- k results. An expired message will be processed by **top- k re-evaluation module**. Specifically, it will check against a *result buffer*, which maintains the top- k results (possibly including some non-top- k results) of all the subscriptions. For the subscriptions that cannot be refilled through result buffer, their top- k results will be re-evaluated from scratch against a *message index* containing all the messages over the sliding window. Note that the message index can be implemented with any existing spatial-keyword index, such as IR-Tree [15] and S2I [33]. Skype can also support subscription update efficiently. A new subscription will be inserted into subscription index, with its top- k results being initialized against message index, while an unregistered subscription will be deleted from both subscription index and result buffer. Note that the subscription index and message index serve different purposes and cannot be trivially combined together.

5 Message Dissemination

In this section, we introduce a novel subscription index, which groups similar subscriptions, to support real-time dissemination against message stream. Specifically, two key techniques, i.e., individual pruning and group pruning, are proposed in Section 5.1 and Section 5.2 respectively, followed by the detailed indexing structure in Section 5.3. Finally, we introduce dissemination algorithm in Section 5.4 and index maintenance in Section 5.5.

5.1 Individual Pruning Technique

For each incoming message m , the key challenge is to determine all the subscriptions whose top- k results are affected. Specifically, we denote the k -th highest score of a subscription s as $\text{kScore}(s)$. Then the top- k results of s need to be updated if $\text{kScore}(s) \leq \text{Score}(s, m)$. In this section, we propose a novel *location-aware prefix filtering* technique to skip an individual subscription efficiently.

5.1.1 Location-aware Prefix Filtering

For ease of exposition, we denote a spatial similarity upper bound between a subscription s and a message m as $\text{SSimUB}(s, \rho, m, \rho)$, which will be discussed in detail in Section 5.1.2. Based on Equation 1, we can derive

a textual similarity threshold for pruning purpose accordingly:

$$\lambda_T(s, \psi, m, \psi) = \frac{\text{kScore}(s)}{1 - s.\alpha} - \frac{s.\alpha}{1 - s.\alpha} \cdot \text{SSimUB}(s, \rho, m, \rho) \quad (2)$$

Then the following lemma claims that if the textual similarity between s and m is less than $\lambda_T(s, \psi, m, \psi)$, we can safely skip s .

Lemma 1 *A message m cannot affect top- k results of a subscription s if $\text{TSim}(s, \psi, m, \psi) < \lambda_T(s, \psi, m, \psi)$.*

Proof It is immediate from Equation 1 and Equation 2.

To utilize Lemma 1, we employ prefix filtering technique, which is widely adopted in textual similarity join problems (e.g., [10, 5, 39]). Prefix filtering is based on the fact that TSim is essentially a vector product; therefore, we can determine the similarity upper bound between two objects by only comparing their prefixes. Before we introduce prefix filtering technique, we first introduce a threshold value for each keyword in s :

$$\text{wtsum}(s, \psi[i]) = \sum_{i \leq j \leq |s.\psi|} \text{wt}(s, \psi[j]) \quad (3)$$

Then we define a *location-aware prefix* as follows.

Definition 4 (Location-aware Prefix) Given a subscription s , a message m and a textual similarity threshold $\lambda_T(s, \psi, m, \psi)$, we use $\text{pref}(s|m) = s.\psi[1 : p]$ to denote the location-aware prefix of s w.r.t. m , where $p = \arg \min_i \{ \text{wtsum}(s, \psi[i+1]) < \lambda_T(s, \psi, m, \psi) \}$.

The following lemma claims that location-aware prefix is sufficient to decide whether a message can be top- k result of a subscription.

Lemma 2 *Given a subscription s and a message m , $\text{pref}(s|m) \cap m.\psi = \emptyset$ is sufficient to skip s regarding m .*

Proof Since $\text{pref}(s|m) \cap m.\psi = \emptyset$, $\text{TSim}(s, \psi, m, \psi) \leq \sum_{p+1 \leq i \leq |s.\psi|} \text{wt}(s, \psi[i]) \cdot 1.0 < \lambda_T(s, \psi, m, \psi)$, where p is defined in Definition 4 and 1.0 is the maximum weight for keyword in m . Then the lemma holds immediately based on Lemma 1.

Example 2 Figure 4 shows an example of location-aware prefix, with 3 registered subscriptions and 3 incoming messages. The underlined value to the right of each keyword corresponds to its weight, and we do not normalize the keyword weight for simplicity. Assuming $\text{SSimUB}(s_1, \rho, m_1, \rho) = 0.98$, then $\lambda_T(s_1, \psi, m_1, \psi) = \frac{0.7}{1-0.6} - \frac{0.6}{1-0.6} \cdot 0.98 = 0.28$. Thus, $\text{pref}(s_1|m_1) = \{w_1, w_2, w_3\}$. Since $m_1.\psi \cap \text{pref}(s_1|m_1) = \emptyset$, we can skip s_1 w.r.t. m_1 .

Definition 8 (Spatial Upper Bound) Given a subscription s inside a cell c and a message m outside c , the spatial upper bound, denoted as $\text{SSimUB}(s, \rho, m, \rho)$ is computed as $1.0 - \frac{\text{mindist}(s, \rho, c, r) + \text{mindist}(m, \rho, c, r)}{\text{MaxDist}}$, where $\text{mindist}(s, \rho, c, r)$ and $\text{mindist}(m, \rho, c, r)$ are the same as Definition 6 and 7.

Following the example in Figure 4, by combining both inner and outer distance, we can get a tighter spatial upper bound between s_2 and m_2 as $\text{SSimUB}(s_2, \rho, m_2, \rho) = 1 - \frac{0.25+0.3}{1.6} = 0.65$.

Note that the inner spatial bound can be pre-computed and materialized, while the outer spatial bound has to be computed on-the-fly as it is relevant to the location of an arriving message. However, the computation cost of $\text{SSimUB}_{\text{out}}(m, \rho, c, r)$ is not expensive since we only need to compute this value against each leaf cell. Finally, we remark that when s and m are within the same cell, both $\text{SSimUB}_{\text{in}}(s, \rho, c, r)$ and $\text{SSimUB}_{\text{out}}(m, \rho, c, r)$ are always 1.0.

Example 5 An example is shown in Figure 4. If we assume the $\text{SSimUB}(s_2, \rho, m_2, \rho) = 1.0$, we have $\lambda_T(s_2, m_2) = \frac{0.6}{1-0.5} - \frac{0.5}{1-0.5} \cdot 1.0 = 0.20$, and $\text{pref}_+(s_2|m_2) = \{w_1, w_2\}$. Thus, s_2 cannot be skipped w.r.t. m_2 . However, if we utilize the inner spatial bound and outer spatial bound together, we have $\text{SSimUB}(s_2, \rho, m_2, \rho) = 1 - \frac{0.25+0.3}{1.6} = 0.65$, $\lambda_T(s_2, m_2) = 0.55$, and $\text{pref}_+(s_2|m_2) = \{w_1\}$. In this case, we can safely skip s_2 w.r.t. m_2 .

5.1.3 Bound Estimation for Unseen Keywords

Since we employ TAAT paradigm to visit inverted file, we can estimate a textual upper bound for unseen keywords. If this upper bound plus the textual similarity that has already been computed is still less than the required threshold, we can safely skip s . The textual upper bound between the unseen keywords of s and m can be computed as follows:

$$\begin{aligned} \text{TSimUB}(s, \psi[i:], m, \psi[j:]) = \min \Big\{ & \text{wtsum}(s, \psi[i]) \\ & \times \text{maxwt}(m, \psi[j]), \text{wtsum}(m, \psi[j]) \times \text{maxwt}(s, \psi[i]) \Big\} \end{aligned} \quad (5)$$

where i and j are starting positions of unseen keywords. Then the following theorem claims we can skip a subscription by utilizing the textual upper bound.

Theorem 2 Given a subscription s , a message m and their textual similarity threshold $\lambda_T(s, \psi, m, \psi)$, assuming we have already computed the partial similarity between $s, \psi[1:i]$ and $m, \psi[1:j]$, denoted as $\text{TSim}(s, \psi[1:i], m, \psi[1:j])$, then $\text{TSim}(s, \psi[1:i], m, \psi[1:j]) + \text{TSimUB}(s, \psi[i+1:], m, \psi[j+1:]) < \lambda_T(s, \psi, m, \psi)$ is sufficient to skip s .

Proof As $\text{TSim}(s, \psi, m, \psi) \leq \text{TSim}(s, \psi[1:i], m, \psi[1:j]) + \text{TSimUB}(s, \psi[i+1:], m, \psi[j+1:])$, we have $\text{TSim}(s, \psi, m, \psi) < \lambda_T(s, \psi, m, \psi)$. The theorem holds immediately from Lemma 1.

Example 6 In Figure 4, consider that we are currently disseminating m_3 . Based on the dissemination algorithm to be discussed later in Section 5.4, we need to traverse the inverted lists in cell c_3 (where s_3 resides) for all the keywords in m_3, ψ one by one. We first check the inverted list of w_2 since w_2 is the 1st keyword of m_3 . Assuming $\lambda_T(s_3, \psi, m_3, \psi) = 0.48$, we cannot skip s_3 since $w_2 \in \text{pref}_+(s_3|m_3) = \{w_1, w_2\}$. However, since w_2 is the 2nd keyword in s_3, ψ , we can compute $\text{TSim}(s_3, \psi[1:2], m_3, \psi[1:1]) = 0.4 \cdot 0.7 = 0.28$, and $\text{TSimUB}(s_3, \psi[3:], m_3, \psi[2:]) = \min(0.4 \cdot 0.3, 0.5 \cdot 0.3) = 0.12$. Because $0.28 + 0.12 < \lambda_T(s_3, \psi, m_3, \psi) = 0.48$, we can immediately skip s_3 .

5.2 Group Pruning Technique

After applying individual pruning technique, many subscriptions can be skipped without the need to compute their exact similarity w.r.t. a message. To further enhance the performance, we propose a novel *Group Pruning Technique* such that we can skip a group of subscriptions without the need to visit them individually. To begin with, we first define *subscription-dependent prefix* for a message.

Definition 9 (Refined Sub-dependent Prefix)

Given a message m , a subscription s and $\lambda_T(s, \psi, m, \psi)$, we use $\text{pref}_+(m|s) = m, \psi[1:p]$ to denote the refined subscription-dependent prefix of m w.r.t. s , where $p = \arg \min_i \{ \text{maxwt}(s, \psi[j]) \times \text{wtsum}(m, \psi[i+1]) < \lambda_T(s, \psi, m, \psi) \}$ with $s, \psi[j] = m, \psi[i+1]$.

The following lemma claims the refined subscription-dependent prefix is sufficient to determine whether a message could be top- k result of a subscription.

Lemma 4 Given a subscription s and a message m , $\text{pref}_+(m|s) \cap s, \psi = \emptyset$ is sufficient to skip s regarding m .

Proof Since $\text{pref}_+(m|s) \cap s, \psi = \emptyset$, $\text{TSim}(s, \psi, m, \psi) \leq 0 + \text{wtsum}(m, \psi[p+1]) \cdot \text{maxwt}(s, \psi[j]) < \lambda_T(s, \psi, m, \psi)$, where p is defined in Definition 9 and $s, \psi[j] = m, \psi[p+1]$. Then the lemma holds immediately based on Lemma 1.

Let us denote the posting list of keyword w in cell c as $\text{plist}(c, w)$, which contains all the subscriptions having w and residing in c . Then based on Lemma 4, for a subscription s in $\text{plist}(c, w)$, if $s, w \notin \text{pref}_+(m|s)$, we can safely skip s . Further, if this holds for a group of subscriptions on $\text{plist}(c, w)$, we can safely skip the whole group as follows.

Lemma 5 Given a message m , a keyword $w \in m.\psi$, a posting list $plist(c, w)$ and a group of subscriptions \mathcal{G} inside $plist(c, w)$, $\max_{s \in \mathcal{G}} \{\maxwt(s.w)\} \cdot wtsum(m.w) < \min_{s \in \mathcal{G}} \{\lambda_T(s.\psi, m.\psi)\}$ is sufficient to skip the whole group \mathcal{G} .

Proof For each $s \in \mathcal{G}$, $\maxwt(s.w) \cdot wtsum(m.w) < \lambda_T(s.\psi, m.\psi)$ holds which indicates $s.w \notin \text{pref}_+(m|s)$ according to Definition 9. Thus, s can be skipped based on Lemma 4, and therefore \mathcal{G} can be skipped immediately.

The left side of the inequality in Lemma 5 can be computed in $O(1)$ time since we can materialize $\max_{s \in \mathcal{G}} \{\maxwt(s.w)\}$ for each group. However, for the right side, it would be quite inefficient if we compute it on the fly for each new message. To avoid this, we propose a lower bound for $\min_{s \in \mathcal{G}} \{\lambda_T(s.\psi, m.\psi)\}$ which can be computed in constant time. In the following, we first present the subscription grouping strategy and then introduce the details of the lower bound deduction.

5.2.1 α -Partition Scheme

Intuitively, we should group subscriptions with similar $\lambda_T(s.\psi, m.\psi)$ such that we can get a tighter textual threshold for the group. We first let $\text{SSimUB}(s.\rho, m.\rho) = \text{SSimUB}_{out}(m.\rho, c.r)$. Note that we compute $\text{SSimUB}(s.\rho, m.\rho)$ by utilizing $\text{SSimUB}_{out}(m.\rho, c.r)$ only in order to make $\text{SSimUB}(s.\rho, m.\rho)$ independent of a specific subscription. Therefore, it is observed from Equation 2 that, for the computation of $\lambda_T(s.\psi, m.\psi)$, only $\frac{kScore(s)}{1-s.\alpha}$ and $\frac{s.\alpha}{1-s.\alpha}$ are dependent on s while $\text{SSimUB}(s.\rho, m.\rho)$ is irrelevant to s . For simplicity, we denote $\frac{kScore(s)}{1-s.\alpha}$ as $kScore^*(s)$ and $\frac{s.\alpha}{1-s.\alpha}$ as $s.\alpha^*$ respectively. Then, we partition subscriptions into groups based on their α^* values, such that the subscriptions inside a group have similar α^* values. We employ a quantile-based method to partition the domain of α^* to ensure that each group has similar number of subscriptions. Then, we can skip the whole group \mathcal{G} as stated in the following theorem.

Theorem 3 Given a group \mathcal{G} generated by α -partition in a posting list $plist(c, w)$, we denote $\min_{s \in \mathcal{G}} \{kScore^*(s)\}$ as $kScore^*(\mathcal{G})$ and $\max_{s \in \mathcal{G}} \{s.\alpha^*\}$ as $\mathcal{G}.\alpha^*$. then $\max_{s \in \mathcal{G}} \{\maxwt(s.w)\} \cdot wtsum(m.w) < kScore^*(\mathcal{G}) - \mathcal{G}.\alpha^* \cdot \text{SSimUB}_{out}(m.\rho, c.r)$ is sufficient to skip the whole group \mathcal{G} .

Proof It is obvious that for any subscription s in \mathcal{G} , we have $\max_{s' \in \mathcal{G}} \{\maxwt(s'.w)\} \cdot wtsum(m.w) < kScore^*(\mathcal{G}) - \mathcal{G}.\alpha^* \cdot \text{SSimUB}(c.r, m.\rho) \leq kScore^*(s) - s.\alpha^* \cdot \text{SSimUB}(c.r, m.\rho) = \lambda_T(s.\psi, m.\psi)$. Thus, we have $\max_{s' \in \mathcal{G}} \{\maxwt(s'.w)\} \cdot wtsum(m.w) <$

$\min_{s' \in \mathcal{G}} \{\lambda_T(s'.\psi, m.\psi)\}$. Combined with Lemma 5, the theorem holds immediately.

Time complexity. The condition checking in Theorem 3 takes $O(1)$ time, since we can precompute the values of $kScore^*(\mathcal{G})$ and $\mathcal{G}.\alpha^*$.

5.2.2 Early Termination Within Group

When a group \mathcal{G} cannot be skipped given a message, we have to check each subscription in it. To avoid this, we propose an *early termination technique* to early stop within a group when the group cannot be skipped totally. To enable early termination, for each group \mathcal{G} in $plist(c, w)$, we sort the subscriptions in \mathcal{G} by their $kScore^*$ values increasingly. For each subscription s in \mathcal{G} , we denote the subscriptions with $kScore^*$ not less than $kScore^*(s)$ as $\mathcal{G}[s] = \{s' \in \mathcal{G} | kScore^*(s') \geq kScore^*(s)\}$, and maintain two statistics $\maxwt(\mathcal{G}[s])$ and $\mathcal{G}[s].\alpha^*$ w.r.t. keyword w as follows:

$$\maxwt(\mathcal{G}[s]) = \max_{s' \in \mathcal{G}[s]} \{\maxwt(s'.w)\} \quad (6)$$

$$\mathcal{G}[s].\alpha^* = \max_{s' \in \mathcal{G}[s]} \{s'.\alpha^*\} \quad (7)$$

Then we can employ early termination as follows.

Theorem 4 Given a group \mathcal{G} inside a posting list $plist(c, w)$, and assuming \hat{s} is the subscription with smallest position in \mathcal{G} such that the following inequality holds: $\maxwt(\mathcal{G}[\hat{s}]) \cdot wtsum(m.w) < kScore^*(\hat{s}) - \mathcal{G}[\hat{s}].\alpha^* \cdot \text{SSimUB}_{out}(m.\rho, c.r)$, then there is no need to check the subscriptions after \hat{s} (including \hat{s} itself).

Proof For any subscription s' after \hat{s} , the following inequalities hold: $kScore^*(\hat{s}) \leq kScore^*(s')$, $\maxwt(\mathcal{G}[\hat{s}]) \geq \maxwt(\mathcal{G}[s'])$, $\mathcal{G}[\hat{s}].\alpha^* \geq \mathcal{G}[s']. \alpha^*$. Thus, $\maxwt(s'.w) \cdot wtsum(m.w) \leq \maxwt(\mathcal{G}[s']) \cdot wtsum(m.w) \leq \maxwt(\mathcal{G}[\hat{s}]) \cdot wtsum(m.w) < kScore^*(\hat{s}) - \mathcal{G}[\hat{s}].\alpha^* \cdot \text{SSimUB}_{out}(m.\rho, c.r) \leq kScore^*(s') - \mathcal{G}[s']. \alpha^* \cdot \text{SSimUB}_{out}(m.\rho, c.r) \leq kScore^*(s') - s'.\alpha^* \cdot \text{SSimUB}_{out}(m.\rho, c.r) = \lambda_T(s'.\psi, m.\psi)$. Based on Definition 9, we know that $s'.w \notin \text{pref}_+(m|s')$. Thus s' can be skipped based on Lemma 4. Thus, the theorem holds immediately.

Time complexity. To speed-up the real-time processing, we precompute $\maxwt(\mathcal{G}[s])$ and $\mathcal{G}[s].\alpha^*$ and store them with each subscription in the group \mathcal{G} . The condition checking in Theorem 4 can be efficiently computed in $O(\log|\mathcal{G}|)$ time with a binary search method.

5.2.3 Cell-based Pruning

Besides the above group pruning technique, we notice that for some cells which are far away from the location of an arriving message, we can safely skip the whole cell.

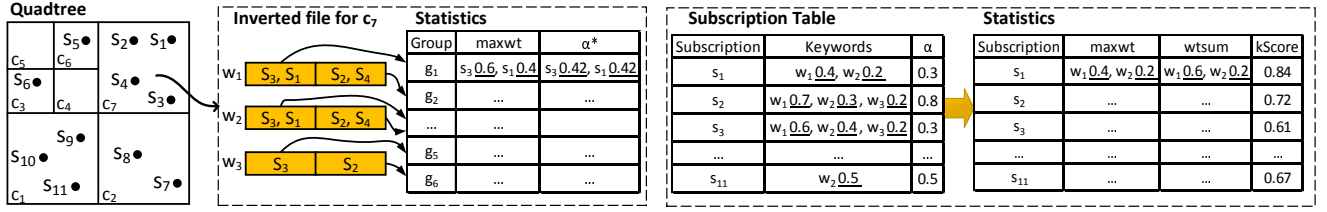


Fig. 5 Subscription index

Specifically, for each subscription s within a cell c , we can derive a spatial similarity threshold as follows:

$$\lambda_S(s, \rho) = \frac{\text{kScore}(s)}{s, \alpha} - \frac{1 - s, \alpha}{s, \alpha} \quad (8)$$

where we assume the textual similarity achieves the largest value, i.e., 1. Then we can reach the following lemma.

Lemma 6 Given a cell c , if $\min_{s \in c} \lambda_S(s, \rho) > \text{SSimUB}_{out}(m, \rho, c, r)$, we can safely skip all the subscriptions in cell c .

Proof For $\forall s \in c$, we have $\lambda_S(s, \rho) = \frac{\text{kScore}(s)}{s, \alpha} - \frac{1 - s, \alpha}{s, \alpha} > \text{SSimUB}_{out}(m, \rho, c, r)$. Thus, $\text{kScore}(s) > (1 - s, \alpha) + s, \alpha \cdot \text{SSimUB}_{out}(m, \rho, c, r) \geq \text{Score}(s, m)$. Thus, m cannot be top- k results of any s in c .

5.3 Subscription Index

Relying on all the techniques discussed above, our subscription index is essentially a **Quadtree** structure integrated with inverted file in each leaf cell, as shown in Figure 5. For each registered subscription, we store its detailed information and relevant statistics in a subscription table, and insert it into a leaf cell of **Quadtree** based on its spatial location. Note that in **Quadtree**, we only store the subscription id referring to its detailed information in subscription table. Within each leaf cell, an inverted file is built upon all the subscriptions inside the cell. Then each posting list in inverted file is further partitioned into groups based on the subscription preference α^* to enable group pruning. Each group is also associated with some statistics mentioned above. Finally, to facilitate early termination, the subscriptions within each group are ordered based on their kScore^* .

5.4 Dissemination Algorithm

Algorithm 1 shows our message dissemination algorithm. We follow a *filtering-and-verification* paradigm, where we first generate a set of candidate subscriptions (Lines 1-1), and then compute the exact scores to determine the truly affected ones, with the updated results being disseminated accordingly (Line 1). Specifically,

we first initialize an empty map \mathcal{R} to store candidates with their scores (Line 1). Then the **maxwt** and **wtsum** values for all the keywords in the arriving message m are computed for later use (Line 1). For each leaf cell c surviving from cell pruning (Line 1), we first compute $\text{SSimUB}_{out}(m, \rho, c, r)$ and then traverse the inverted file in cell c following a TAAT manner. For each group \mathcal{G} encountered in $\text{plist}(c, w)$ (Line 1), we skip \mathcal{G} if group pruning can be applied (Line 1); otherwise, we identify \hat{s} for early termination based on Theorem 4 (Line 1 and Line 1). For each surviving subscription s , we employ location-aware prefix filtering (Line 1) and bound estimation for unseen keywords (Line 1) to skip it as early as possible. For the surviving subscriptions, we store the accumulated textual similarity so far w.r.t. m in \mathcal{R} , while for the skipped subscriptions, we set $\mathcal{R}[s]$ to negative infinity (Line 1). Finally, for each subscription in \mathcal{R} with $\mathcal{R}[s] > 0$, we verify it and update its top- k results if needed (Line 1). Note that when verifying a candidate s , we only need to compute the exact spatial similarity to get the final score because the textual similarity, i.e., $\mathcal{R}[s]$, has already been computed. The statistics relevant to pruning techniques are also updated in Line 1.

5.5 Index Maintenance

Our indexing structure can also support subscription update efficiently. For a new subscription s , we first find the leaf cell containing its location, and then insert it into the inverted file with $O(|s, \psi| \cdot \log |\mathcal{G}|)$ cost. Note that the statistics mentioned above need to be updated accordingly. For an expired subscription, we simply delete it from index and update the statistics if necessary.

6 Top-k Re-evaluation

In this section, we present the details of top- k re-evaluation module. We first introduce some background knowledge for k -skyband in Section 6.1. Then we present our *cost-based skyband* technique in detail in Section 6.2. In the following of this paper, we denote

Algorithm 1: MessageDissemination(m)

```

Input :  $m$  : a new incoming message
1  $\mathcal{R} := \emptyset$  /* A candidate map */;
2 for  $1 \leq i \leq |m.\psi|$  do
3    $\lfloor$  Compute  $\text{maxwt}(m.\psi[i])$  and  $\text{wtsum}(m.\psi[i])$ ;
4 for each leaf cell  $c$  in the Quadtree do
5   if  $c$  is skipped by Lemma 6 then /* Cell pruning */
6      $\lfloor$  Continue;
7   Compute outer spatial bound  $\text{SSimUB}_{out}(m.\rho, c.r)$ ;
8   for  $1 \leq i \leq |m.\psi|$  do
9      $w := m.\psi[i]$ ;
10    for each group  $\mathcal{G}$  in  $\text{plist}(c, w)$  do
11      if  $\max_{s \in \mathcal{G}} \{\text{maxwt}(s.w)\} \cdot \text{wtsum}(m.w) <$ 
12         $\text{kScore}^*(\mathcal{G}) - \mathcal{G}.\alpha \cdot \text{SSimUB}_{out}(m.\rho, c.r)$ 
13        then /* Group pruning based on Theorem 3 */
14           $\lfloor$  Continue;
15      Identify  $\hat{s}$  based on Theorem 4;
16      for each subscription  $s$  in group  $\mathcal{G}$  do
17        if  $s == \hat{s}$  then /* Early termination based on Theorem 4 */
18           $\lfloor$  Break;
19        if  $s \in \mathcal{R}$  &&  $\mathcal{R}[s] == -\infty$  then
20           $\lfloor$  Continue;
21        Compute  $\lambda_T(s.\psi, m.\psi)$  based on both inner and outer spatial bounds;
22        if
23           $\text{maxwt}(m.w) \cdot \text{wtsum}(s.w) < \lambda_T(s.\psi, m.\psi)$ 
24          then /* Theorem 1 */
25             $\lfloor$  Continue;
26        if  $s \in \mathcal{R}$  then
27           $\lfloor$   $\mathcal{R}[s] := \text{wt}(s.w) \cdot \text{wt}(m.w)$ ;
28        else
29           $\mathcal{R} := \mathcal{R} \cup s$ ;
30           $\mathcal{R}[s] := \text{wt}(s.w) \cdot \text{wt}(m.w)$ ;
31         $\text{pos} :=$  the position of keyword  $w$  in  $s.\psi$ ;
32        if  $\mathcal{R}[s] + \text{TSimUB}(s.\psi[\text{pos} + 1 : ], m.\psi[i + 1 : ]) < \lambda_T(s.\psi, m.\psi)$  then
33          /* Theorem 2 */
34           $\lfloor$   $\mathcal{R}[s] := -\infty$ ;
35 VerifyAndUpdate( $\mathcal{R}$ ) and disseminate updated results to corresponding subscriptions;

```

the k -skyband buffer (either fully or partially) of a subscription s as $s.\mathcal{A}$ for simplicity, and the exact top- k results are denoted as $s.\mathcal{A}_k$. Meanwhile, $s.k$ is denoted as k if it is clear from context.

6.1 K-Skyband

The idea of utilizing k -skyband to reduce the number of re-evaluations for top- k queries over a sliding window is first proposed in [29]. In particular, for a given subscription s , only the messages in its corresponding k -skyband can appear in its top- k results over the sliding window, thus being maintained. Following are formal definitions of *dominance* and *k-skyband*.

Algorithm 2: TopkRe-evaluation(m)

```

Input :  $m$  : an expired message
1 for each subscription  $s$  whose  $k$ -skyband buffer contains  $m$ 
2   do
3     Delete  $m$  from  $s.\mathcal{A}$ ;
4     if  $|s.\mathcal{A}| < k$  then
5       Compute the best  $s.\theta$  based on our cost model;
6       Retrieve
7        $B := \{m | m \in \mathcal{W} \ \&\& \ \text{Score}(s, m) \geq s.\theta\}$ ;
8        $s.\mathcal{A} := k\text{-skyband of } B$ ;
9     Extract  $s.\mathcal{A}_k$  from  $s.\mathcal{A}$ ;

```

Definition 10 (Dominance) A message m_1 dominates another message m_2 w.r.t. a subscription s if both $\text{Score}(s, m_1) \geq \text{Score}(s, m_2)$ and $m_1.t > m_2.t$ hold.

Definition 11 (k -skyband) The k -skyband of a subscription s , denoted as $s.\mathcal{A}$, contains a set of messages which are dominated by less than k other messages.

Instead of keeping k -skyband over all the messages in the sliding window, which is cost-prohibitive, Mouratidis *et al.* [29] maintain a *partial k-skyband*. Specifically, they only maintain the messages with score not lower than a threshold $s.\theta$, where $s.\theta$ is the $\text{kScore}(s)$ after the most recent top- k re-evaluation for s and remains unchanged until next re-evaluation is triggered. However, as our experiments suggest, the method in [29] may result in expensive computational cost due to the improper selection of $s.\theta$.

To alleviate the above problem, we propose a novel *cost-based k-skyband* technique, which judiciously selects a best threshold $s.\theta$ for the k -skyband maintenance of each subscription. To start with, we present an overview of our top- k re-evaluation algorithm in Algorithm 2. For each subscription s containing the expired message m , if the size of $s.\mathcal{A}$ after deleting m is less than k , we need to re-evaluate its top- k results from scratch. Specifically, we first compute a proper threshold $s.\theta$ based on our cost model (Line 2), and then re-compute k -skyband buffer $s.\mathcal{A}$ based on B , which contains all the messages with score at least $s.\theta$ (Line 2 and Line 2) ⁶. Note that B can be computed by utilizing message index. Finally, we extract top- k results from $s.\mathcal{A}$ (Line 2). The key challenge here is to estimate a best threshold $s.\theta$, which will be discussed in the following in detail. We remark that we use the term *re-evaluation* to refer in particular to the top- k re-computation against message index.

6.2 Cost-based K-Skyband

The general idea of our cost-based k -skyband model is to select a best threshold $s.\theta$ for each subscription such

⁶ The same technique in [29] is used to compute k -skyband.

that the overall cost defined in the cost model can be minimized. The following theorem guarantees that, as long as we maintain a partial k -skyband over all the messages with score not lower than $s.\theta$, we can extract top- k results from partial k -skyband safely when some message expires.

Theorem 5 *Given a subscription s , let $kScore_{last}(s)$ be the $kScore(s)$ after the most recent top- k re-evaluation for s . We always have $s.A_k \subseteq s.A$ if the following conditions hold: (1) $|s.A| \geq k$; (2) $s.A$ is a partial k -skyband which is built over all the messages with score at least $s.\theta$ in the sliding window, where $0 \leq s.\theta \leq kScore_{last}(s)$.*

Proof We prove it by contradiction. Assuming there exists a message $m \in s.A_k$ while $m \notin s.A$, then we discuss two possible cases: (1) $Score(s, m) \geq s.\theta$; (2) $Score(s, m) < s.\theta$. For the first case, since $m \notin s.A$, m must be dominated by more than k messages in $s.A$, which indicates it cannot be top- k results, i.e., $m \notin s.A_k$. For the second case, at least k messages in $s.A$ must have a higher score than m because $|s.A| \geq k$ and all the messages in $s.A$ have score at least $s.\theta$. Thus, m still cannot be top- k results. Thus, the original assumption does not hold, which immediately indicates $s.A_k \subseteq s.A$.

Thus, based on Theorem 5, we can safely extract top- k results from k -skyband buffer $s.A$ when $|s.A| \geq k$; when $|s.A| < k$, we have to re-evaluate from message index.

Our cost-based k -skyband model, based on Theorem 5, aims to find the best $s.\theta$ such that the overall cost can be minimized for each subscription. We mainly consider two costs. The first one is *k -skyband maintenance cost*, denoted as $C_{sm}(s)$, which is triggered upon message arrival and expiration. The second one is *top- k re-evaluation cost*, denoted as $C_{re}(s)$, which is triggered when some message expires and the top- k results can no longer be retrieved from k -skyband buffer. We aim to estimate the expected overall cost w.r.t. each message update, i.e., message arrival and message expiration, each of which we assume occurs with probability $\frac{1}{2}$ as the window slides. To simplify the presentation, we denote as $\text{prob}(s.\theta)$ the probability that the score between a random message and a subscription s is at least $s.\theta$. We may immediately derive $\text{prob}(s.\theta)$ for a given $s.\theta$ from historical data, assuming the score follows previous distribution. The details of these two costs are presented in the following respectively.

6.2.1 K -Skyband Maintenance Cost

The maintenance of k -skyband is triggered when the following two types of updates happen, both with prob-

ability $\frac{1}{2} \cdot \text{prob}(s.\theta)$, where $\frac{1}{2}$ is the probability of message arrival or message expiration due to the count-based sliding window, and $\text{prob}(s.\theta)$ is the probability that the score between a random message and s is at least $s.\theta$. Please note that if the independence assumption does not hold for messages, the above probabilities cannot be estimated accurately, and we may resort to utilizing historical data for the estimation.

The first type of update is triggered when a message m with score at least $s.\theta$ arrives. Apart from the insertion of m into $s.A$, the dominance counters of all the messages in $s.A$ with score not higher than $Score(s, m)$ will increase by 1, and the messages with dominance counter equal to k will be evicted. Since we implement our k -skyband buffer with a linked list sorted by $Score(s, m)$. The above operations can be processed in $O(|s.A|)$ time with a linear scan. The next challenge is to estimate $|s.A|$. Based on the independence assumption between score dimension and time dimension, the expected number, i.e., $|s.A|$, of messages in the partial k -skyband is $k \cdot \ln(\frac{|\mathcal{W}| \cdot \text{prob}(s.\theta)}{k})$ [44], where $|\mathcal{W}|$ is the size of sliding window. Please note that if the independence assumption does not hold, the worst case space complexity will be $|\mathcal{W}| \cdot \text{prob}(s.\theta)$.

The second type of update occurs when an old message m among the k -skyband buffer of s expires. In this case, we only need to delete m from $s.A$ in $O(|s.A|)$ time. Note that m does not dominate any remaining messages and therefore the dominance counters of the remaining messages are not affected. Finally, we get the total cost of k -skyband maintenance as follows:

$$\begin{aligned} C_{sm}(s) &= \frac{1}{2} \cdot \text{prob}(s.\theta) \cdot |s.A| + \frac{1}{2} \cdot \text{prob}(s.\theta) \cdot |s.A| \\ &= \text{prob}(s.\theta) \cdot k \cdot \ln\left(\frac{|\mathcal{W}| \cdot \text{prob}(s.\theta)}{k}\right) \end{aligned} \quad (9)$$

6.2.2 Top- k Re-evaluation Cost

The top- k re-evaluation cost can be formalized as:

$$C_{re}(s) = C_{topk}(s) \cdot \frac{1}{\mathbb{Z}(s)} \quad (10)$$

where $C_{topk}(s)$ is the average top- k computation cost over message index for subscription s , and $\mathbb{Z}(s)$ is the expected number of message updates that is required to trigger top- k re-evaluation, i.e., leading to $|s.A| < k$. The value of $C_{topk}(s)$ can be estimated by the average of previous top- k computation cost against message index. The remaining issue is how to estimate $\mathbb{Z}(s)$, which is non-trivial.

To solve this problem, we model the streaming updating process as a simple *random walk*. A random walk is a stochastic sequence RW_n , with RW_0 being the starting position, defined by $RW_n = \sum_{i=1}^n X_i$ where

X_i are independent and identically distributed random variables (i.e., i.i.d.). The random walk is *simple* if $\text{prob}(X_i = 1) = p$, $\text{prob}(X_i = -1) = q$ and $\text{prob}(X_i = 0) = r$, where $p + q + r = 1$. We map the estimation of $\mathbb{Z}(s)$ into a simple random walk as follows. We model the change of k -skyband buffer $s.\mathcal{A}$ w.r.t. each message update as an i.i.d. variable X_i . X_i is set to 1 when the size of $s.\mathcal{A}$ is increased by 1 at i -th step, while X_i is set to -1 when the size of $s.\mathcal{A}$ is decreased by 1. When the size of $s.\mathcal{A}$ does not change, X_i is set to 0. Unfortunately, it is difficult to estimate the probability of $\text{prob}(X_i = 1)$ and $\text{prob}(X_i = -1)$ for each message update, due to the eviction of messages by dominance relationship. For example, for a new message, the size of $s.\mathcal{A}$ may decrease rather than increase due to the eviction of messages with dominance counter reaching k . To address this problem, rather than estimating $\mathbb{Z}(s)$ for $s.\mathcal{A}$ maintenance, we estimate $\mathbb{Z}'(s)$ for $s.\mathcal{A}'$, which contains all the messages with score not lower than $s.\theta$. Specifically, when we maintain $s.\mathcal{A}'$, we do not consider the dominance relationship between messages for each message update, and thus the messages dominated by k (or more) messages are not evicted. Clearly, $s.\mathcal{A}'$ is a superset of $s.\mathcal{A}$, i.e., $s.\mathcal{A} \subseteq s.\mathcal{A}'$. The following theorem guarantees that $\mathbb{Z}(s)$ is equal to $\mathbb{Z}'(s)$.

Theorem 6 *The expected number of message updates that is required to trigger top-k re-evaluation for $s.\mathcal{A}$ maintenance is the same as that for $s.\mathcal{A}'$ maintenance, i.e., $\mathbb{Z}(s) = \mathbb{Z}'(s)$.*

Proof To show $\mathbb{Z}(s)$ is equal to $\mathbb{Z}'(s)$, it is sufficient to prove that $|s.\mathcal{A}| < k$ if and only if $|s.\mathcal{A}'| < k$ at any point. Initially, we have $s.\mathcal{A} \cup s.\mathcal{A}_{\text{evict}} = s.\mathcal{A}'$ and $s.\mathcal{A} \cap s.\mathcal{A}_{\text{evict}} = \emptyset$, where $s.\mathcal{A}_{\text{evict}}$ is a set of messages evicted by messages in $s.\mathcal{A}$. We prove that after each message update, $s.\mathcal{A} \cup s.\mathcal{A}_{\text{evict}} = s.\mathcal{A}'$ always holds. As to the arrival of a new message m , if m is inserted into $s.\mathcal{A}'$, it will also be inserted into $s.\mathcal{A}$ since m will not be dominated by any existing message due to its freshness; and vice versa. Note that some messages may be evicted from $s.\mathcal{A}$ to $s.\mathcal{A}_{\text{evict}}$ due to dominance relationship. Thus, $s.\mathcal{A} \cup s.\mathcal{A}_{\text{evict}} = s.\mathcal{A}'$ holds. As to the expiration of an old message m , (1) If $m \in s.\mathcal{A}$, m will also expire from $s.\mathcal{A}'$, while $s.\mathcal{A}_{\text{evict}}$ does not change. (2) If $m \notin s.\mathcal{A}$, m will expire from both $s.\mathcal{A}_{\text{evict}}$ and $s.\mathcal{A}'$. Thus, $s.\mathcal{A} \cup s.\mathcal{A}_{\text{evict}} = s.\mathcal{A}'$ still holds. Therefore, when $|s.\mathcal{A}| < k$ occurs, $s.\mathcal{A}_{\text{evict}}$ must be empty because there is no k messages in $s.\mathcal{A}$ that can dominate any message in $s.\mathcal{A}_{\text{evict}}$. Thus, $|s.\mathcal{A}'| = |s.\mathcal{A}| < k$ holds. Contrarily, when $|s.\mathcal{A}'| < k$ occurs, it is immediate that $|s.\mathcal{A}| < k$ because $s.\mathcal{A}$ is always a subset of $s.\mathcal{A}'$. Therefore, the theorem holds.

Based on the above theorem, we turn to estimate $\mathbb{Z}'(s)$, which is much easier. Now the probability distribution of X_i can be estimated as:

tribution of X_i can be estimated as:

$$\text{prob}(X_i) = \begin{cases} \frac{1}{2} \cdot \text{prob}(s.\theta) & \text{if } X_i = 1, \\ \frac{1}{2} \cdot \text{prob}(s.\theta) & \text{if } X_i = -1, \\ 1 - \text{prob}(s.\theta) & \text{if } X_i = 0. \end{cases} \quad (11)$$

We denote the initial size of $s.\mathcal{A}'$ as $|s.\mathcal{A}'_{\text{init}}|$. Now, the estimation of $\mathbb{Z}'(s)$ is equivalent to a well-known random walk problem, namely *Monkey at the cliff with reflecting barriers* [20]. Specifically, we set the starting position RW_0 as $|s.\mathcal{A}'_{\text{init}}|$ and the destination position as $k - 1$; the i.i.d. variable X_i is defined as Equation 11; and the reflecting barrier is set as $2 \cdot RW_0$. By applying some mathematical reduction based on the property of random walk [20], we get the following result.

$$\begin{aligned} \mathbb{Z}'(s) = & \frac{2 \cdot (|s.\mathcal{A}'_{\text{init}}| - k + 1) \cdot |s.\mathcal{A}'_{\text{init}}|}{\text{prob}(s.\theta)} \\ & + \frac{(|s.\mathcal{A}'_{\text{init}}| - k + 1) \cdot (|s.\mathcal{A}'_{\text{init}}| - k + 2)}{\text{prob}(s.\theta)} \end{aligned} \quad (12)$$

where $|s.\mathcal{A}'_{\text{init}}|$ can be estimated as $\text{prob}(s.\theta) \cdot |\mathcal{W}|$. Thus, the top- k re-evaluation cost in Equation 10 can be estimated by replacing $\mathbb{Z}(s)$ with $\mathbb{Z}'(s)$. Based on Equation 9 and Equation 10, we get our final cost model:

$$\mathcal{C}(s) = \mathcal{C}_{sm}(s) + \mathcal{C}_{re}(s) \quad (13)$$

To minimize Equation 13 where the only variable is $s.\theta$, we employ an incremental estimation algorithm similar to gradient descent [3] to compute the best value of $s.\theta$.

Remark. To accommodate our cost-based skyband model with the message dissemination algorithm, we need to replace $\text{kScore}(s)$ in Section 5 with $s.\theta$ such that any message with score not lower than $s.\theta$ will be considered to possibly affect the top- k results of s . Moreover, since our dominance definition simply depends on the 2-dimensional score-time space while is irrelevant to the exact score function, our technique can be easily applied to other top- k monitoring problems with different score functions.

6.3 Discussions

Initialization of incoming subscriptions. The initialization of a new subscription s can be processed in a similar way to Algorithm 2, where we regard the initial size of $s.\mathcal{A}$ as 0 and execute Lines 2-2 in Algorithm 2 sequentially.

Time-based sliding window model. Our techniques discussed above can also be extended to support time-based sliding window model, where only the messages within a recent time period are maintained. Unlike count-based sliding window whose size is constant, the

size of time-based sliding window, i.e., $|\mathcal{W}|$, can change at any time due to the volatile message workload. To estimate $|\mathcal{W}|$, we assume that the message workload does not change significantly in the near future. Then we can estimate $|\mathcal{W}|$ by the historical message workload from a recent period. Another difference is that the probability of message arrival (resp. expiration) cannot be regarded as $\frac{1}{2}$ trivially as indicated in Equation 9 and Equation 11, because the number of message arrival and the number of message expiration are possibly rather different in each timestamp. To alleviate this issue, we resort to estimating the above probabilities based on the relative proportion of message arrival and expiration within a recent time period. Then the probabilities (e.g., $\frac{1}{2}$) in Equation 9 and Equation 11 are updated accordingly. We also conduct experiments to verify the efficiency of our techniques under time-based sliding window in Section 8.

7 Distributed Processing

In this section, we introduce DSkype, a distributed top- k spatial-keyword publish/subscribe system built on top of Storm. We first touch some background knowledge about Storm in Section 7.1, followed by the detailed system framework in Section 7.2. Four novel distribution mechanisms are discussed in Section 7.3, which manage to partition the subscriptions and messages to multiple bolt instances for parallel processing. The maintenance issue is finally discussed in Section 7.4. To the best of our knowledge, this is the first work to extend top- k spatial-keyword publish/subscribe system to a distributed environment.

7.1 Storm Background

Storm is a distributed, fault-tolerant and general-purpose stream processing system. Unlike Hadoop⁷ which is mainly designed to process batch tasks, Storm is designed to process streaming data continuously and endlessly. There are three key abstractions in Storm: **spout**, **bolt** and **topology**. A spout is a source of streams, which reads input stream from external resources, such as Twitter API⁸. A bolt is a processing unit responsible for data processing, which handles any number of input streams and produces any number of new output streams. A topology is a network of spouts and bolts, with each directed edge in the network representing a bolt subscribing to the output stream of some other spout or bolt. Essentially, topology defines the working flow of a real-time computation task, which is

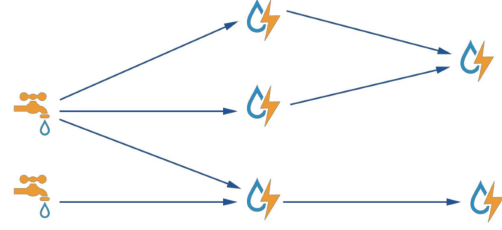


Fig. 6 A simple storm flow

similar to a MapReduce job [16]. Storm employs various **stream groupings** techniques⁹, such as shuffle grouping and fields grouping, to specify for each bolt instance which streams it should receive as input. Figure 6 depicts a simple Storm working flow where there are two spouts and five bolts connected by directed edges. Note that each spout/bolt can have many parallel-running instances/tasks.

7.2 Framework

Figure 7 shows the topology of our DSkype, which contains five main components:

- **Subscription/Message spouts.** Subscription spouts receive new subscription request while message spouts collect message stream from external source, e.g., Twitter API. The incoming streams are then subscribed by other components of the topology.
- **Distribution bolts.** Distribution bolts receive streams from spouts and navigate them downwards to the subscription bolts according to some carefully designed distribution mechanisms (Section 7.3), aiming to achieve good workload balance. This component is critical to the overall communication cost and throughput of our system. Note that the distribution bolts will also route new messages to the message bolts to ensure that the sliding window is always up-to-date.
- **Subscription bolts.** Subscription bolts partition the subscription index and result buffer among multiple parallel-running tasks. A new subscription from a distribution bolt is inserted into one or more subscription bolts, and a new message is processed simultaneously among multiple bolts. Note that distribution bolts and subscription bolts together correspond to message dissemination module in the centralized version discussed in Section 4.
- **Message bolts.** Message bolts maintain the sliding window in a distributed manner, and each bolt contains part of the sliding window. A message index (e.g., IR-Tree, S2I) is built over the messages residing in each bolt. The top- k re-evaluation request

⁷ Apache Hadoop project. <https://hadoop.apache.org/>

⁸ <https://dev.twitter.com/rest/public>

⁹ <http://storm.apache.org/releases/0.10.0/Concepts.html>

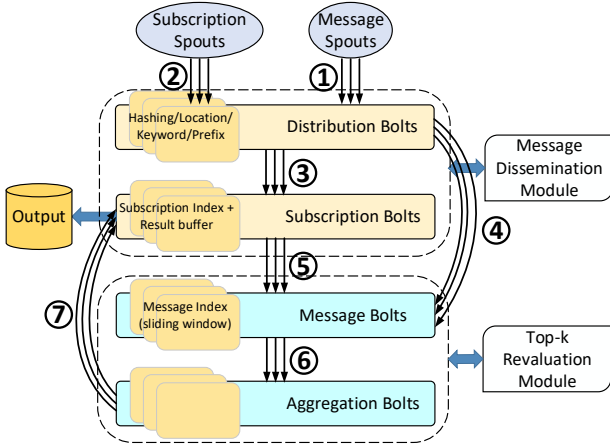


Fig. 7 DSkye topology (solid arrows indicate stream flow between components)

for a subscription s issued by subscription bolts will be processed concurrently among all message bolts, each generating a partial message buffer consisting of all the messages with score at least $s.\theta$. Note that each message is stored in only one message bolt.

- **Aggregation bolts.** Aggregation bolts are introduced to aggregate the partial message buffer generated by message bolts. Then the final k -skyband buffer is computed and forwarded back to the subscription bolts where s resides. Note that message bolts and aggregation bolts together form the counterpart to the top- k re-evaluation module in its centralized version.

All the stream groupings in the topology are summarized in Table 1. We remark that the result buffer can be easily swapped to any persistent state, such as Memcached¹⁰ and HDFS [36], to support various applications.

Working flow. When a new message m is digested by a message spout, it will be delivered to a distribution bolt (stream ①, see Figure 7). The distribution bolt then will navigate m to some of subscription bolts (stream ③) such that it can be processed against local subscription index in a parallel manner. The distribution bolt will also disseminate m to a message bolt (stream ④) to keep the sliding window therein up-to-date. When m expires, the top- k re-evaluations triggered by m in the subscription bolts will be emitted to message bolts (stream ⑤) and then to aggregation bolts (stream ⑥), where the k -skyband buffers of all affected subscriptions will be re-computed and forwarded back to the subscription bolts (stream ⑦). Similarly, a new subscription s will be firstly delivered from a subscription spout to a distribution bolt (stream ②), and then routed to one or more subscription bolts (stream ③)

for indexing based on the distribution mechanism. The initial result of s will be computed from message bolts (stream ⑤) and aggregation bolts (stream ⑥) accordingly and forwarded back to subscription bolt (stream ⑦) where s is indexed. An unregistered subscription will be simply deleted from all its residing subscription bolts.

Challenges. As the number of subscriptions increases, the subscription bolts become the main bottleneck of our system. Meanwhile, the communication cost between distribution bolts and subscription bolts dominates all the other communication cost as we increase the number of subscription bolts. Thus, the key challenge in DSkye is to develop an efficient distribution mechanism to assign subscriptions and messages only to some inevitable subscription bolts, such that both small communication cost and high throughput can be realized while still guaranteeing the correctness of our algorithms. Furthermore, the distribution mechanism should be able to handle workload balance, since both the subscription workload and message workload in real life are extremely biased regarding keywords and locations. At last, the distribution mechanism should be light-weighted without consuming many CPU and memory resources.

7.3 Distribution Mechanism

In this section, we present several novel, efficient and light-weighted distribution mechanisms, which can be integrated into distribution bolts. For the ease of exposition, we assume we already have a set of existing subscriptions \mathcal{S} and a random message m sampled from message stream. We denote the number of subscription bolts as N_{sb} , with each bolt identified by a partition index ranging from 0 to $N_{sb}-1$. A distribution mechanism aims to partition \mathcal{S} into N_{sb} subscription bolts and navigate the message m to relevant subscription bolts for top- k dissemination. In the following, we propose four different distribution methods, namely *hashing-based*, *location-based*, *keyword-based* and *prefix-based*, respectively.

7.3.1 Hashing-based Method

Hashing-based method partitions the subscriptions based on a uniform hashing function defined as follows:

$$h(s) = s.id \mod N_{sb} \quad (14)$$

where $s.id$ is a unique id assigned to each subscription, and $h(s)$ is the bolt index where s should be allocated.

Analysis. Since each subscription is allocated to only one bolt, the *replication ratio* of subscriptions is 1. The replication ratio here indicates the number of times a

¹⁰ <http://storm.apache.org/releases/0.10.0/Trident-tutorial.html>

Table 1 Stream grouping methods. ID in first column corresponds to ID in Figure 7.

ID	Source	Destination	Grouping	Description
①	message spouts	distribution bolts	shuffle grouping	Each message is distributed to only one distribution bolt randomly.
②	subscription spouts	distribution bolts	shuffle grouping	Each subscription is distributed to only one distribution bolt randomly.
③	distribution bolts	subscription bolts	all/direct grouping	Subscriptions/messages are distributed to subscription bolts with either all grouping or direct grouping (depending on the distribution mechanism employed).
④	distribution bolts	message bolts	direct grouping	Each message is distributed to only one message bolt based on its unique id to keep the sliding window up-to-date.
⑤	subscription bolts	message bolts	all grouping	Each top- k re-evaluation request is distributed to all message bolts for parallel processing.
⑥	message bolts	aggregation bolts	files grouping	The partial result buffers from the same subscription are distributed to the same aggregation bolt.
⑦	aggregation bolts	subscription bolts	direct grouping	The aggregated k -skyband buffer is forwarded back to one or more subscription bolts where this subscription resides.

subscription has been stored in the system. Note that we ignore the replication ratio of messages, because it is always 1 regardless of the distribution mechanisms. Meanwhile, for any new message, it needs to be distributed to all the subscription bolts to ensure the correctness. Thus, the average *communication cost* of each message is N_{sb} . Note that we only consider the communication cost between distribution bolts and subscription bolts w.r.t. each message since it is dominant. For example, on a cluster with 32 subscription bolts, the communication cost between distribution bolts and subscription bolts account for more than 90% of total communication cost.

7.3.2 Location-based Method

Hashing-based method is simple and can achieve very good workload balance, because the number of subscriptions in each subscription bolt is nearly the same by the nature of uniform hashing. However, it does not take the location factor into consideration. Intuitively, distributing subscriptions with high spatial similarity into the same bolt can lead to lower *AMP* cost and thus higher throughput, since we can acquire better spatial bounds as discussed in Section 5.1.2. On the other hand, it is also pivotal to balance subscription and message workloads among the subscription bolts. To this end, we propose a cost-based spatial partition schema using KD-Tree [6], where the leaf nodes of KD-Tree form a disjoint partition of the whole space. For each leaf node, we allocate a subscription bolt to process all the subscriptions whose locations are inside the leaf node. Formally, given any node, denoted as nd , in KD-Tree, we estimate its cost by:

$$C(nd) = N(nd) \times p(nd) \quad (15)$$

where $N(nd)$ is the number of subscriptions whose locations are inside nd and $p(nd)$ is the probability that

Algorithm 3: SpatialPartition($nd, depth$)

Input : nd : a node in KD-Tree
 $depth$: the depth of nd
Output : A spatial partition consisting of all the leaf nodes in KD-Tree

```

1 if  $depth == maxdepth$  then
2   return;
3 if  $depth$  is even then /* split by x-coordinate */
4   Find a value  $x$  from the x-coordinates of all the
   subscriptions in  $nd$ , which leads to minimum
    $|C(nd_1) - C(nd_2)|$ , where  $nd_1$  and  $nd_2$  are two
   child nodes split based on  $x$ ;
5 else /* split by y-coordinate */
6   Find a value  $y$  from the y-coordinates of all the
   subscriptions in  $nd$ , which leads to minimum
    $|C(nd_1) - C(nd_2)|$ , where  $nd_1$  and  $nd_2$  are two
   child nodes split based on  $y$ ;
7 Assign subscriptions in  $nd$  to  $nd_1$  and  $nd_2$  based on
   the splitting value;
8 SpatialPartition( $nd_1, depth + 1$ );
9 SpatialPartition( $nd_2, depth + 1$ );
```

a random incoming message falls inside nd . Note that $p(nd)$ can be easily estimated from historical message workloads. The cost-based spatial partition algorithm is depicted in Algorithm 3, which is very similar to the original KD-Tree construction algorithm. The key difference is that, unlike KD-Tree which selects a line halving all the points along x-axis or y-axis alternately, our algorithm tries to find a splitting line which minimizes the cost difference between two children (Lines 3 and 3) in order to achieve workload balance. We limit the total number of partitions (i.e., leaf nodes) by setting a *maxdepth* value (Line 3).

Analysis. The time complexity of Algorithm 3 is bounded by $O(|\mathcal{S}| \times \log |\mathcal{S}|)$ as we need to sort subscriptions in \mathcal{S} by x-coordinate and y-coordinate respectively beforehand and then do a divide-and-conquer partition. Each subscription will be assigned to only one subscription bolt containing its location. Thus, the replication ratio is 1. On the other hand, each incoming message

needs to be delivered to all subscription bolts to guarantee the algorithm correctness. Thus, the average communication cost is N_{sb} .

7.3.3 Keyword-based Method

Both hashing-based and location-based methods have to send each message to all the subscription bolts, which results in high communication overhead especially when we increase the number of subscription bolts. To alleviate this issue, we propose a novel keyword-based partition mechanism, which can reduce the communication cost significantly at the cost of small subscription replications. The general idea is that, each subscription bolt only accounts for a subset of keywords; thus, each subscription s (resp. message m) will be distributed only to the subscription bolts whose keyword sets overlap with $s.\psi$ (resp. $m.\psi$). To start with, similar to Section 7.3.2, we estimate the processing cost of a subscription bolt with keyword set W_i as follows:

$$C(W_i) = N(W_i) \times p(W_i) \quad (16)$$

where $N(W_i)$ is the number subscriptions whose keywords overlap with W_i and $p(W_i)$ is the probability that a random incoming message contains at least one keyword from W_i . Note that $p(W_i)$ can be estimated from historical message workloads. We then define *variance* to measure the workload balance as follows:

$$Var = \frac{1}{N_{sb}} \times \sum_{i=0}^{N_{sb}-1} (C(W_i) - C_\mu)^2 \quad (17)$$

where C_μ is the average cost over all keyword sets.

We are now ready to present the keyword partition problem: given a keyword vocabulary, denoted as \mathcal{V} , where keywords are ordered by their frequencies, we aim to partition \mathcal{V} into N_{sb} subsets, i.e., $W_0, W_1, \dots, W_{N_{sb}-1}$, each covering a number of consecutive keywords, such that the following conditions are satisfied: (1). $\forall 0 \leq i < j \leq N_{sb} - 1, W_i \cap W_j = \emptyset$, (2). $\cup_{0 \leq i \leq N_{sb}-1} W_i = \mathcal{V}$, (3). Var is minimized.

Since the search space is exponentially large, we resort to an efficient heuristic algorithm to solve this problem, which is demonstrated in Algorithm 4. The idea is to recursively partition the keyword set into two subsets by carefully selecting a splitting keyword to balance the cost between two subsets. The parameter *maxdepth* here (Line 4) is used to control the number of partitions we need. Once we get keyword partitions, a subscription s will be allocated to the i^{th} bolt as long as $s.\psi \cap W_i \neq \emptyset$. Similarly, a message m will also be distributed to the i^{th} bolt if $m.\psi \cap W_i \neq \emptyset$.

Example 7 Figure 8 shows an example of keyword-based method. Assume there are 10 keywords in the vocabulary, and there are four subscription bolts. Based

Algorithm 4: KeywordPartition($W, depth$)

Input : W : current keyword set
 $depth$: the depth of recursion

Output : A keyword partition

- 1 **if** $depth == maxdepth$ **then**
- 2 **return**;
- 3 Find a splitting keyword w from W which leads to minimum $|C(W_1) - C(W_2)|$, where W_1 contains all the words on the left of w (inclusive) while W_2 contains all the keywords on the right of w (exclusive);
- 4 KeywordPartition($W_1, depth + 1$);
- 5 KeywordPartition($W_2, depth + 1$);

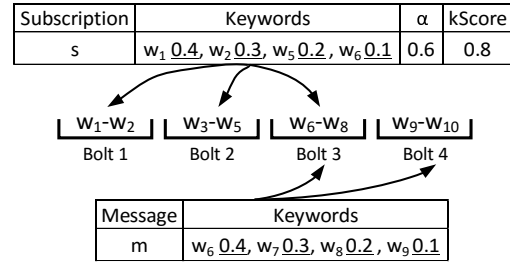


Fig. 8 Keyword-based example

on the keyword-based distribution method, subscription s is allocated to bolt 1, bolt 2 and bolt 3 while message m is distributed to bolt 3 and bolt 4.

Analysis. The time complexity of Algorithm 4 is $O(|\mathcal{V}| \times maxdepth + |\mathcal{S}|)$ if we use a linear scan to find the splitting keyword (Line 4). Each subscription is distributed to at most $|s.\psi|$ bolts and each message is distributed to at most $|m.\psi|$ bolts. Thus, the replication ratio is at most $|s.\psi|$ while the average communication cost is at most $\min(|m.\psi|, N_{sb})$.

7.3.4 Prefix-based Method

Compared to both hashing-based and location-based methods, keyword-based method can reduce the communication cost significantly, especially when we have a large number of subscription bolts (i.e., large N_{sb}), since each message is only distributed to keyword-overlapping bolts. However, the pitfall is that we have to duplicate subscriptions among different subscription bolts to ensure the correctness of our algorithms, which often leads to poor throughput as shown in the experimental part. To further reduce the subscription replications as well as improve throughput, we propose a light-weighted prefix-based method, which only distributes subscriptions based on their keyword prefixes, rather than all the keywords. However, we cannot employ the same prefix as defined in Definition 4, because the value of $SSimUB(s, \rho, m, \rho)$ is not available beforehand. To overcome this issue, we define a new textual similar-

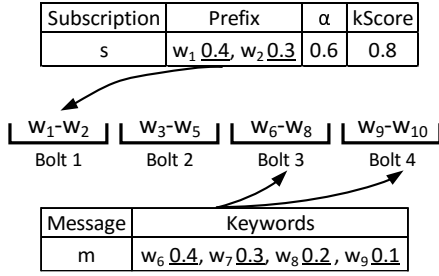


Fig. 9 Prefix-based example

ity threshold $\lambda_T(s.\psi)$ as follows:

$$\lambda_T(s.\psi) = \frac{\text{kScore}(s)}{1 - s.\alpha} - \frac{s.\alpha}{1 - s.\alpha} \cdot 1.0 \quad (18)$$

where we always assume $\text{SSimUB}(s.\rho, m.\rho)$ to be 1.0 regardless of the actual location of the incoming message. Then, we propose a *loose prefix*, denoted as $\text{pref}(s)$, which is defined as follows.

Definition 12 (Loose Prefix) Given a subscription s and a textual similarity threshold $\lambda_T(s.\psi)$, we use $\text{pref}(s) = s.\psi[1 : p]$ to denote the loose prefix of s , where $p = \arg \min_i \{\text{wtsum}(s.\psi[i + 1]) < \lambda_T(s.\psi)\}$.

We remark that this definition is similar to Definition 4, except that we replace $\lambda_T(s.\psi, m.\psi)$ with $\lambda_T(s.\psi)$, such that the loose prefix $\text{pref}(s)$ is irrelevant to the message. The following lemma guarantees the correctness of our distribution mechanism based on loose prefix.

Lemma 7 *Our algorithm is correct as long as we distribute each subscription s (resp. message m) only to the subscription bolts whose keyword sets intersect with $\text{pref}(s)$ (resp. $m.\psi$).*

Proof Based on the distribution method, it is immediate to conclude that if $\text{pref}(s) \cap m.\psi \neq \emptyset$ (and thus m might be a candidate of s), m must be distributed to at least one subscription bolt where s also resides. On the other hand, if $\text{pref}(s) \cap m.\psi = \emptyset$, m must not be a candidate of s . This can be easily proved based on a similar deduction from the proof of Lemma 2. Thus, there is no need to distribute the subscription s to the subscription bolts whose keyword sets overlap with $s.\psi - \text{pref}(s)$.

Example 8 Following the same example in Figure 8, Figure 9 depicts an example of prefix-based method. From Equation 18, we get $\lambda_T(s.\psi) = \frac{0.8}{1-0.6} - \frac{0.6}{1-0.6} = 0.5$, and therefore $\text{pref}(s) = \{w_1, w_2\}$. Thus, subscription s is only allocated to bolt 1 while message m is still distributed to bolt 3 and bolt 4. It is obvious that the computation between s and m is avoided since m cannot be the top- k results of s .

Analysis. The replication ratio is determined by $|\text{pref}(s)|$ based on Lemma 7, which is usually smaller

Table 2 Summary of distribution mechanisms (worst case)

Method	Replica. ratio	Comm. cost
Hashing-based	1	N_{sb}
Location-based	1	N_{sb}
Keyword-based	$ s.\psi $	$\min(m.\psi , N_{sb})$
Prefix-based	$ \text{pref}(s) $	$\min(m.\psi , N_{sb})$

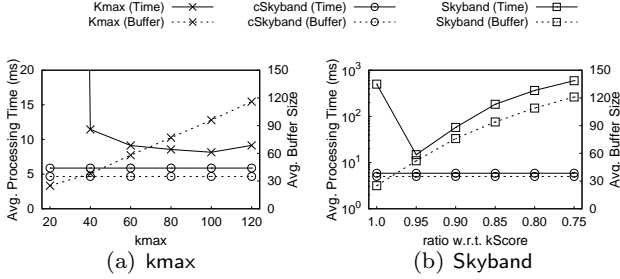
than $|s.\psi|$. Meanwhile, the average communication cost is bounded by $\min(|m.\psi|, N_{sb})$, which is the same as keyword-based method. As shown in the experiments, the prefix-based method not only can reduce the replication ratio but also can improve the system throughput with a large margin compared to keyword-based method. Note that Algorithm 4 is also used here to get keyword partitions, except that the value of $N(W_i)$ in Equation 16 needs to be recomputed since we only use the loose prefix of subscription. Besides, a subscription may need to be reallocated among subscription bolts when $\text{pref}(s)$ changes due to the updating of $\text{kScore}(s)$. We employ a lazy reallocation strategy, where the reallocation is triggered only when $\text{pref}(s)$ needs to cover more keywords.

Remark. We remark that all the distribution mechanisms discussed above are light-weighted indexes employed in the distribution bolts in order to facilitate the distribution of subscriptions and messages. This is different from the subscription index built in each subscription bolt, which aims to accelerate top- k dissemination regarding each incoming message.

Discussions. We summary all the distribution mechanisms discussed above in Table 2, where we report the replication ratio of subscriptions and average communication cost w.r.t. each message, both of which are the main factors of system throughput. From the table, we notice that when $|m.\psi| < N_{sb}$, both keyword-based and prefix-based methods are likely to benefit a lot from the large reduction in communication cost, and are expected to have higher throughput; however, when $|m.\psi| > N_{sb}$, both keyword-based and prefix-based methods may not perform very well because they cannot reduce communication cost by a large margin, while suffered from the extra cost triggered by replicates. In Section 8.2, we have conducted detailed experiments to evaluate the performances of different algorithms under different settings, which further verify our analysis above. Furthermore, workload balance is also a main factor contributing to the system throughput. The hashing-based method can achieve best balance due to the nature of uniform hashing. However, for all the other three methods, we observe that the workload balance can also be well maintained, since we take into consideration both the subscription workload and message workload during distribution.

Table 3 Datasets Statistics

Datasets	TWEETS	GN	YELP
# of msg	12.7M	2.2M	1.6M
Vocabulary size	1.7M	208K	85K
Avg. # of msg keywords	9	7	37
Size in GB	2.26	0.3	1.04

**Fig. 10** Tuning baseline algorithms

7.4 Maintenance

Storm provides a graphical interface to monitor the workload of each spout/bolt¹¹. We fork a background process to access the interface and monitor the workload of each component automatically and periodically. When some component becomes overloaded, we may simply increase its parallelism, while decreasing its parallelism if it becomes idle. The subscription bolts discussed in the above distribution mechanisms can be easily further partitioned or merged according to the system workload. For example, for the location-based method, we can simply split a leaf node if it becomes overloaded.

8 Experiments

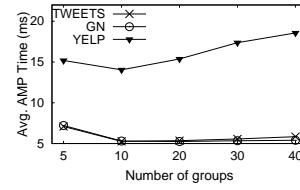
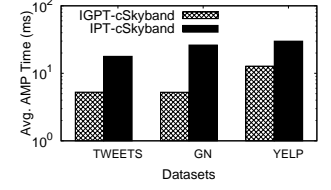
8.1 Centralized Evaluations

In this section, we conduct extensive experiments to verify the efficiency and effectiveness of *Skype* in a single machine. All experiments are implemented in C++, and conducted on a PC with 3.4GHz Intel Xeon 2 cores CPU and 32GB memory running Red Hat Linux. Following the typical setting of publish/subscribe systems (e.g., [24,12]), we assume indexes fit in main memory to support real-time response.

8.1.1 Experimental Setup

As this is the first work to study top- k spatial-keyword publish/subscribe over sliding window, we extend previous work [12] to sliding window. We implement and compare the following algorithms. For message dissemination module:

¹¹ <http://storm.apache.org/releases/0.10.0/STORM-UI-REST-API.html>

**Fig. 11** Vary # groups**Fig. 12** Pruning methods

- **CIQ**. The subscription index proposed in [12]¹².
- **IGPT**. The subscription index proposed in our paper, which combines both **Individual** and **Group** Pruning Techniques.

For top- k re-evaluation module:

- **Skyband**. The k -skyband algorithm proposed in [29].
- **kmax**. The kmax algorithm proposed in [41].
- **cSkyband**. The cost-based k -skyband algorithm proposed in our paper.

Note that our *Skype* algorithm is the combination of IGPT and cSkyband. We apply IR-Tree [15] to index messages.

Datasets. Three datasets are deployed for experimental evaluations. *TWEETS* is a real-life dataset collected from Twitter [25], containing 12.7M tweets with geo-locations from 2008 to 2011. *GN* is obtained from the US Board on Geographic Names¹³ in which each message is associated with a geo-location and a short text description. *YELP* is obtained from Yelp¹⁴, which contains user reviews and check-ins for thousands of businesses. The statistics of three datasets are summarized in Table 3.

Subscription workload. We generate top- k subscriptions based on the above datasets. For each dataset, 1M geo-textual messages are randomly selected. For each selected message, we randomly pick j keywords as subscription keywords with $1 \leq j \leq 5$. The weight of each keyword is computed according to *tf-idf* weighting scheme¹⁵. The subscription location is the same as message location. For each subscription, the preference parameter α is randomly selected between 0 and 1, while the default value of k , i.e., the number of top- k results, is set to 20.

Message workload. Our simulation starts when the sliding window with default size of 1M is full and continuously runs for 100,000 arriving messages over the sliding window.

We report the average processing time, including average arriving message processing time (i.e., **AMP**) and average expired message processing

¹² The time decay function and related index in CIQ are removed to adapt to our problem.

¹³ <http://geonames.usgs.gov>

¹⁴ <http://www.yelp.com/>

¹⁵ <https://en.wikipedia.org/wiki/Tfidf>

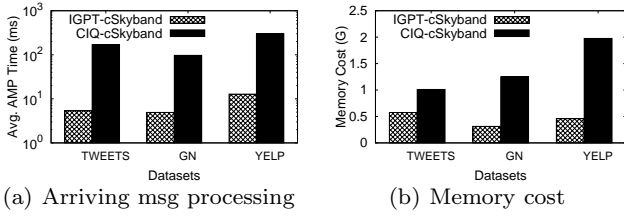


Fig. 13 Compare dissemination algorithms

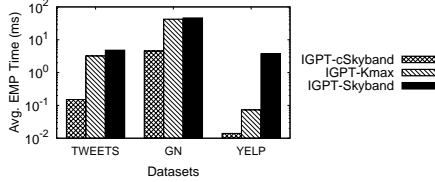


Fig. 14 Compare re-evaluation algorithms

time (i.e., **EMP**), as well as the index size. By default, the number of α -partition groups is set to 10. The maximum number of subscriptions that can be stored in each cell is set to 1000.

8.1.2 Experimental Tuning

Tuning kmax and Skyband. In the first set of experiments, we tune the performance of both kmax and Skyband techniques in Figure 10 on *TWEETS* dataset, where IGPT algorithm is employed for message dissemination. For better understanding, we evaluate average processing time (denoted as solid line) and average buffer size (denoted as dotted line) in the same figure. We also show the results of our IGPT-cSkyband algorithm under default settings which remains unchanged. For kmax algorithm (Figure 10(a)), we vary kmax from 20 to 120. It is noticed that a small kmax leads to high top- k re-evaluation cost while a large kmax results in high message dissemination cost and buffer maintenance cost. We set 60 as the default kmax value since it strikes a good trade-off between performance and buffer size (i.e., memory cost). For Skyband algorithm, we vary the threshold score $s.\theta$ from $1.0 \times kScore$ to $0.75 \times kScore$ where the smaller value leads to larger buffer size. It is noticed that when the ratio is 1.0 which is the original setting in [29], the performance of Skyband is poor due to the frequent re-evaluations. For comparison fairness, $s.\theta$ is set to its sweet point $0.95 \times kScore$ for Skyband in the following experiments. It is worth mentioning that our cSkyband always outperforms Skyband under all settings because cSkyband can tune a best threshold for each individual subscription based on the cost model while there is no sensible way to tune Skyband for millions of subscriptions. The similar trends are also observed on other datasets.

Vary the number of groups in α -partition. Figure 11 reports the *AMP* time of Skype algorithm

Table 4 Average buffer size of different re-evaluation algorithms (Unit: 1)

Algorithm	<i>TWEETS</i>	<i>GN</i>	<i>YELP</i>
IGPT-cSkyband	35	33	28
IGPT-Kmax	58	58	59
IGPT-Skyband	52	58	56

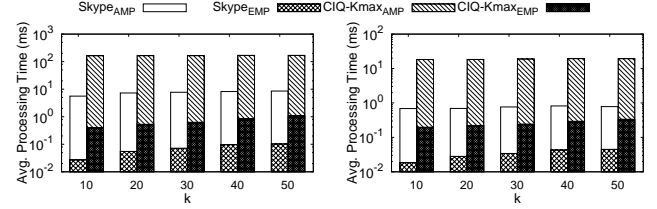


Fig. 19 Zipf distribution workload

Fig. 20 Performance of time-based sliding window

against three datasets where the number of groups varying from 5 to 40. It is shown that we can achieve a good trade-off between the group filtering effectiveness and group checking costs when the number of groups is set to 10, which is used as default value in the following experiments.

Effect of pruning techniques. In this experiment, we compare the *AMP* time of different pruning techniques employed in our message dissemination module. Specifically, we compare IGPT with IPT, which only employs individual pruning technique in Figure 12. We observe that IGPT algorithm can achieve at least three times improvement compared with IPT over all the datasets, which verifies the efficiency of our group pruning techniques. This is mainly because the group pruning technique can skip the whole group without the need to check individual subscription, and can terminate early within a group. In the following experiments, we always use IGPT as our dissemination algorithm.

8.1.3 Performance Evaluation

Compare message dissemination algorithms. In this experiment, we compare the performance of different dissemination algorithms. Specifically, we compare CIQ and IGPT with cSkyband being top- k re-evaluation algorithm. As shown in Figure 13(a), our algorithm can achieve about 10 times faster than CIQ algorithm, due to the benefit of individual pruning technique and group pruning technique. On the other hand, as shown in Figure 13(b), even if we need to maintain some additional statistics, the memory cost of our subscription index is much smaller than that of CIQ, since our algorithm only indexes each subscription into single cell, rather than multiple cells.

Compare top- k re-evaluation algorithms. In this experiment, we compare the performance of different top- k re-evaluation strategies combined with our

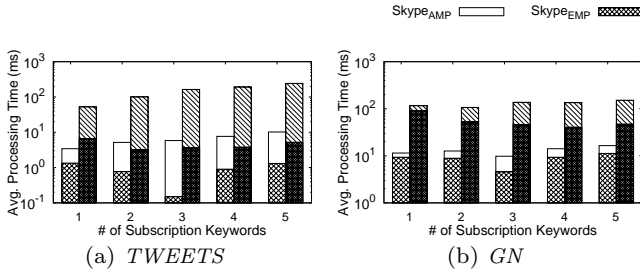


Fig. 15 Effect of number of subscription keywords

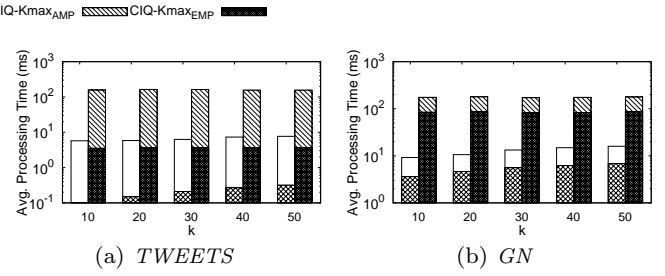
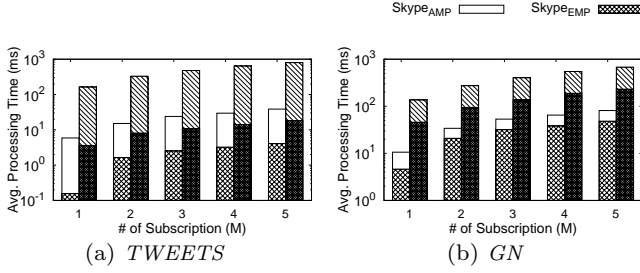
Fig. 16 Effect of number of top- k results

Fig. 17 Effect of number of subscriptions

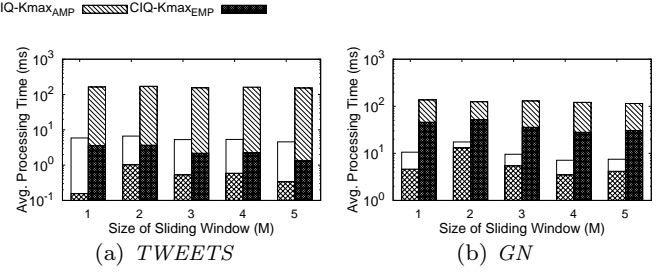


Fig. 18 Effect of sliding window size

IGPT algorithm. Specifically, we compare k max algorithm [41], k -skyband algorithm [29] and our cost-based k -skyband algorithm, which are denoted as IGPT-Kmax, IGPT-Skyband, IGPT-cSkyband respectively. The average *EMP* time is reported in Figure 14. We observe that our cSkyband algorithm can achieve about 4-20 times improvement compared to the second best algorithm. This is mainly due to the adaptiveness of our cost model which can tune a best threshold for each subscription. Table 4 demonstrates the average buffer size of each algorithm. Our algorithm maintains much fewer number of messages than other competitors due to the advantage of our cost model. In the following experiments, we only compare our Skype algorithm (i.e., IGPT-cSkyband) with CIQ-Kmax, which performs best among all the baselines.

Effect of number of subscription keywords. We assess the effect of number of subscription keywords in Figure 15. We notice that the *AMP* time increases as we vary the number of keywords from 1 to 5. This is obvious since more candidates will be encountered during traversing posting lists when the number of keywords is large. As to the *EMP* time, we observe that the selectivity is low and fewer messages are *relevant* at initial, thus leading to high cost. With increasing number of keywords, the selectivity increases, thus reducing the number of re-evaluations accordingly. Finally, when the number of keyword reaches 4 or 5, a message is less likely to have a high score w.r.t. a subscription due to the smaller weight assigned to each subscription keyword on average, resulting in the increase of *EMP*

time. The overall processing time increases decently for a large number of keywords.

Effect of number of top- k results. In this set of experiments, we analyse the effect of number of top- k results, i.e., k , in Figure 16. For *AMP* time, as we increment k from 10 to 50, the average *kScore* of subscriptions decreases; therefore, an arriving message is more likely to influence more subscriptions, leading to high *AMP* time in our algorithm. Meanwhile, a large k usually results in high *EMP* time, because the subscriptions with low selectivity are more likely to expire and incur top- k re-evaluations. Besides, the k -skyband maintenance cost also increases for large k . Overall, the average processing time increases slowly w.r.t. k .

Effect of number of subscriptions. We evaluate the scalability of our system in Figure 17, where we vary the number of subscriptions from 1M to 5M. As shown in the figure, our algorithm scales very well with increasing number of subscriptions, thus making it practical to support real-life applications with fast response.

Effect of sliding window size. We turn to evaluate the effect of sliding window size $|\mathcal{W}|$ in this set of experiments. The results are demonstrated in Figure 18, where we vary $|\mathcal{W}|$ from 1M to 5M. It is observed that, when we increase $|\mathcal{W}|$, the *AMP* time decreases, which is due to the fact that a large sliding window usually leads to better top- k results with higher *kScore*. Thus, a new message will affect less subscriptions, resulting in lower *AMP* cost. Regarding the *EMP* time, it fluctuates around a value due to the competitive results of fewer number of re-evaluations and high query cost against the message index as we increase $|\mathcal{W}|$.

Performance over zipf distribution. We evaluate the performance over a different subscription workload where the keywords are sampled from a zipf distribution (Figure 19). It is observed that the zipf distribution workload has similar performance compared to original workload, where keywords are randomly sampled from messages.

Performance over time-based sliding window. We verify the performance over time-based sliding window against *TWEETS* dataset. We set the sliding window size as 4 months and feed the initial sliding window with tweets from January to April 2010. The tweets from recent one month are used to estimate window size and probability discussed in Section 6.3. The *AMP* and *EMP* time w.r.t. *each timestamp* (i.e., 1 sec) are reported by continuously feeding the sliding window with tweets collected in May 2010. The results are shown in Figure 20 where we vary the number of top-*k* results. It is noticed that *Skype* can still achieve an order of magnitude improvement compared to *CIQ-Kmax*, which verifies the efficiency of our extensions.

8.2 Distributed Evaluations

In this section, we verify the performance of our distributed publish/subscribe system, i.e., *DSkype*. All the experiments are conducted on a 10-node (one nimbus and nine supervisors) cluster running Storm 0.10.0¹⁶, with a single node Zookeeper server¹⁷ deployed for coordination. Each node in the cluster is a Debian 6.0.10 server that has 3.4GHz Intel Xeon 8 cores CPU, 16GB memory, and gigabit ethernet interconnect. Each supervisor node is configured to run at most 3 workers at the same time, and each worker can run multiple spouts/bolts concurrently.

We use the same *TWEETS*, *GN* and *YELP* datasets as above and generate subscription workload and message workload accordingly. The default number of subscriptions is 5M, and the default size of sliding window is 1M. The message workload is fed to the system continuously for one hour. The number of subscription spouts, message spouts, distribution bolts, subscription bolts, message bolts and aggregation bolts is set to 1, 1, 1, 32, 3 and 1 respectively by default. It is noticed that the number of subscription bolts is the largest compared to other components, since subscription bolts are the main bottleneck of our system. The parameters in the subscription index are tuned in a similar way to Section 8.1.2. We use the **throughput**, i.e., the average number of messages processed in one

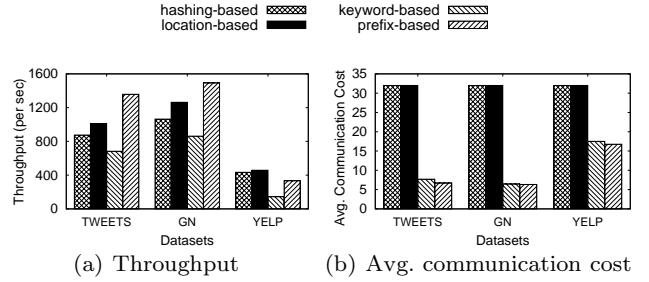


Fig. 21 Compare different distribution mechanisms

Table 5 Replication ratio

Methods	TWEETS	GN	YELP
hashing-based	1	1	1
location-based	1	1	1
keyword-based	3.2	3.3	3.3
prefix-based	1.9	1.9	2.1

second, and the **communication cost**, i.e., the average number of tuples transmitted between distribution bolts and subscription bolts to process one message, as the main measurements. Note that, since the real communication cost heavily depends on the hardware, we use the number of tuples transmitted as the measure of communication cost, which is hardware-independent. Besides, we only consider the communication cost between distribution bolts and subscription bolts because it dominates all the other communication costs. All the measurements are computed after system initialization, which usually takes about 2 minutes.

Compare different distribution mechanisms. In this set of experiments, we compare the performance of our proposed distribution mechanisms over *TWEETS*, *GN* and *YELP* datasets. We denote the four distribution mechanisms as **hashing-based**, **location-based**, **keyword-based** and **prefix-based** respectively. The results are depicted in Figure 21. Regarding *TWEETS* dataset in Figure 21(a), we observe that **location-based** method achieves about 15% throughput improvement than **hashing-based** method. This is mainly due to the better spatial similarity bound computed in **location-based** method. It is also noticed that **keyword-based** method performs the worst among all the methods because of the extra overhead incurred by duplicate subscription allocation, which is shown in Table 5. However, **prefix-based** method has the best performance, about 56% faster than **hashing-based**, at slightly extra cost of replication ratio (i.e., 1.9). Similar trend is also observed on *GN* dataset. We remark that the distributed system is about 26 ~ 49 times faster than its centralized version, which demonstrates the superiority of distributed processing. On the other hand, the communication cost of both **keyword-based** and **prefix-based** methods in *TWEETS* dataset (Figure 21(b)) achieve about 80% reduction compared to **hashing-based** and

¹⁶ <http://storm.apache.org/2015/11/05/storm0100-released.html>

¹⁷ <http://zookeeper.apache.org/doc/r3.4.8/>

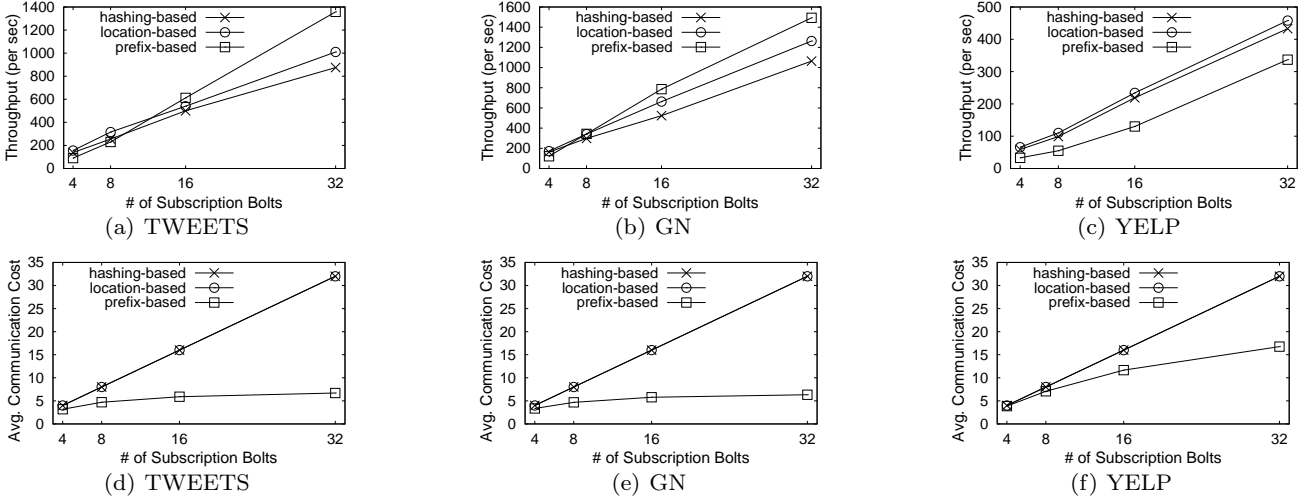


Fig. 22 Effect of number of subscription bolts

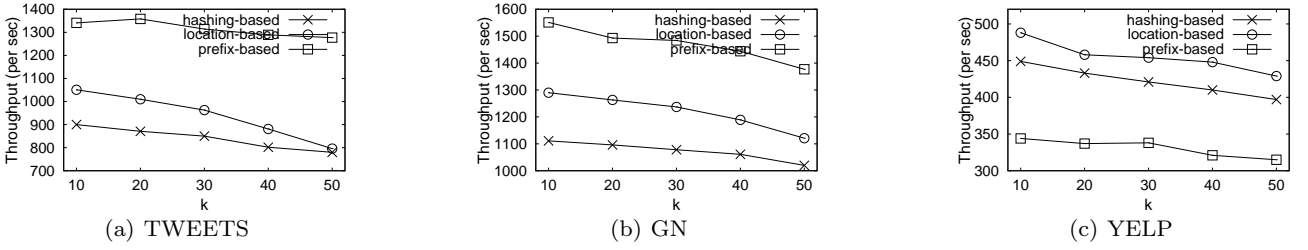


Fig. 23 Effect of number of top-k results

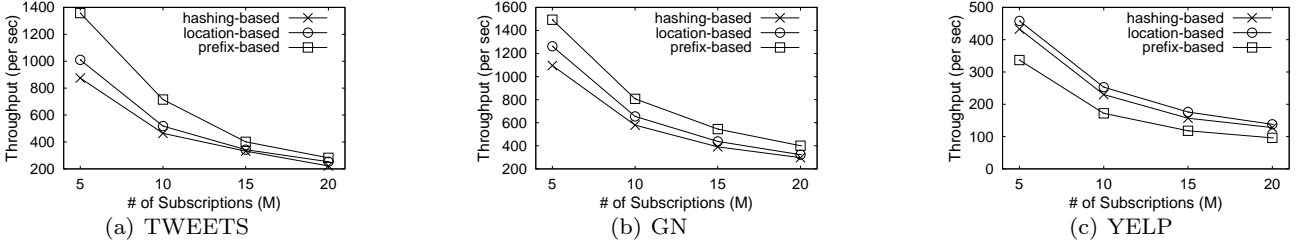


Fig. 24 Effect of number of subscriptions

location-based methods, which demonstrates the benefit of former two methods, especially when the communication cost becomes system bottleneck. Similar trends can also be observed from *GN* dataset. However, as to *YELP* dataset, we notice that even though prefix-based method is much faster than keyword-based method, it is slower than hashing-based and location-based methods. This is mainly because the number of keywords in *YELP* message is much larger than *TWEETS* and *GN*, and thus the reduction in communication cost (Figure 21(b)) is not significant enough to pay off the cost contributed by duplicate allocation. In the following experiments, we omit keyword-based method since prefix-based method has much better overall performance in terms of both throughput and replication ratio.

Effect of number of subscription bolts. We evaluate the effect of number of subscription bolts N_{sb}

in Figure 22 by varying N_{sb} from 4 to 32. In terms of throughput, all the three algorithms can scale very well with increasing number of subscription bolts. We observe that when N_{sb} is small (i.e., 4 or 8), prefix-based method performs worse than hashing-based and location-based methods due to the duplicate subscriptions in *TWEETS* and *GN* (Figure 22(a) and Figure 22(b)). However, when N_{sb} reaches 16 or 32, the benefit of communication cost reduction has become significant (Figure 22(d) and Figure 22(e)), thus contributing to the high throughput of prefix-based method. For *YELP* dataset, prefix-based method processes less messages per second compared to hashing-based and location-based methods while managed to reduce communication by about 50%.

Effect of number of top-k results. We turn to evaluate the effect of number of top-k results in Figure 23.

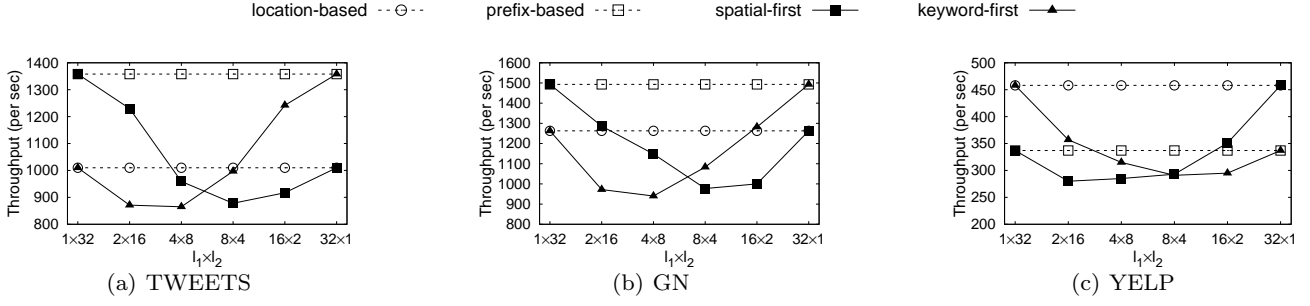


Fig. 25 Throughput of hybrid methods

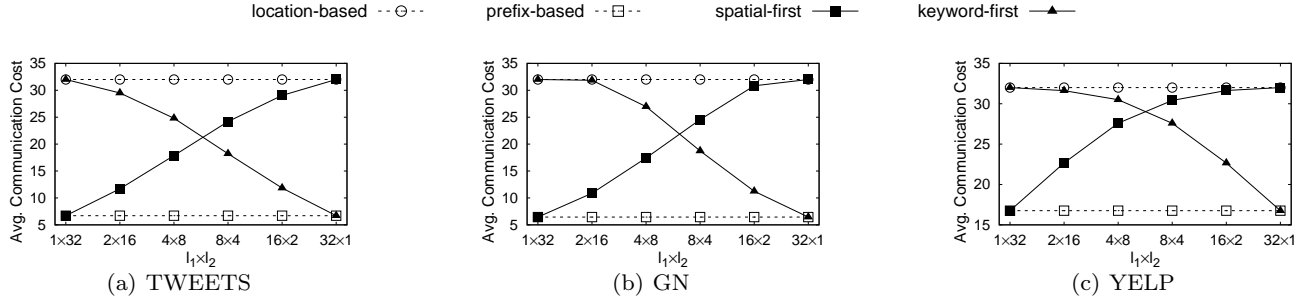


Fig. 26 Communication cost of hybrid methods

We do not show the communication cost because it is irrelevant to the number of top- k results. As shown in Figure 23, the throughputs in all the datasets decrease when we vary k from 10 to 50. This is because large k usually yields high processing cost in subscription index. However, the influence of k is not very significant as the throughputs decrease slowly and linearly.

Effect of number of subscriptions. In Figure 24, we vary the number of subscriptions from 5M to 20M. It is obvious that the throughputs drop in all the datasets when we increase the number of subscriptions. However, the decreasing trend indicates that the average processing time of an incoming message is still linearly proportional to the number of subscriptions, considering the factor that the average processing time is inversely proportional to the throughput.

Hybrid methods. In this set of experiments, we compare our methods with two possible hybrid methods: spatial-first method and keyword-first method. Specifically, spatial-first method employs two-level partition scheme. On the first level, it employs location-based method while on the second level it employs prefix-based method. keyword-first method is similar to spatial-first method except that it employs prefix-based method first and then location-based method. We compare the throughput and communication cost of these two hybrid methods while changing the number of partitions on each level. The results of location-based and prefix-based methods are also shown in dotted line for comparison purpose. Regarding throughput in Figure 25 where l_1 is the number of partitions on the first level

and l_2 is the number of partitions on the second level, all the datasets exhibit similar trends. Specifically, for spatial-first method, when l_1 increases and l_2 decreases, the benefit of keyword partition becomes less while the benefit of spatial partition is still not significant, which leads to decreasing throughput. As we further increase l_1 (e.g., 16×2), the spatial partition becomes dominant and thus improves the throughput which finally reaches the same throughput as location-based method at 32×1 . The trends of keyword-first method can be explained similarly. The communication costs of both spatial-first and keyword-first methods in Figure 26 are between those of location-based and prefix-based due to its hybrid nature. Overall, the hybrid methods do not exhibit large improvement compared to location-based and prefix-based methods.

9 Conclusion

The popularity of streaming geo-textual data offers great opportunity for applications such as information dissemination and location-based campaigns. In this paper, we study a novel problem of continuous top- k spatial-keyword publish/subscribe over sliding window. To maintain top- k results for a large number of subscriptions over a fast stream simultaneously and continuously, we propose a novel indexing structure, which employs both individual pruning technique and group pruning technique, to process a new message instantly on its arrival. In addition, to handle the re-evaluations incurred by expired messages from the slid-

ing window, we develop a novel cost-based k -skyband model with theoretical analysis to judiciously maintain a partial k -skyband for each subscription. Furthermore, a distributed stream processing system called DSkype is developed, which extends Skype to Storm to exploit the benefit of parallel processing. The experiments demonstrate that both Skype and its distributed version DSkype can achieve high throughput performance over geo-textual stream.

References

1. Aji, A., Wang, F., Vo, H., Lee, R., Liu, Q., Zhang, X., Saltz, J.H.: Hadoop-gis: A high performance spatial data warehousing system over mapreduce. *PVLDB* **6**(11), 1009–1020 (2013)
2. Aly, A.M., Mahmood, A.R., Hassan, M.S., Aref, W.G., Ouzzani, M., Elmeleegy, H., Qadah, T.: AQWA: adaptive query-workload-aware partitioning of big spatial data. *PVLDB* **8**(13), 2062–2073 (2015)
3. Avriel, M.: Nonlinear programming: analysis and methods. Courier Corporation (2003)
4. Babcock, B., Babu, S., Datar, M., Motwani, R., Widom, J.: Models and issues in data stream systems. In: *PODS* (2002)
5. Bayardo, R.J., Ma, Y., Srikant, R.: Scaling up all pairs similarity search. In: *WWW*, pp. 131–140 (2007)
6. Bentley, J.L.: Multidimensional binary search trees used for associative searching. *Commun. ACM* **18**(9), 509–517 (1975)
7. Böhm, C., Ooi, B.C., Plant, C., Yan, Y.: Efficiently processing continuous k-nn queries on data streams. In: *ICDE*, pp. 156–165 (2007)
8. Broder, A.Z., Carmel, D., Herscovici, M., Soffer, A., Zien, J.Y.: Efficient query evaluation using a two-level retrieval process. In: *CIKM*, pp. 426–434 (2003)
9. Buckley, C., Lewit, A.F.: Optimization of inverted vector searches. In: *SIGIR*, pp. 97–110 (1985)
10. Chaudhuri, S., Ganti, V., Kaushik, R.: A primitive operator for similarity joins in data cleaning. In: *ICDE*, p. 5 (2006)
11. Chen, L., Cong, G., Cao, X.: An efficient query indexing mechanism for filtering geo-textual data. In: *SIGMOD* (2013)
12. Chen, L., Cong, G., Cao, X., Tan, K.: Temporal spatial-keyword top-k publish/subscribe. In: *ICDE* (2015)
13. Chen, L., Cong, G., Jensen, C.S., Wu, D.: Spatial keyword query processing: An experimental evaluation. *PVLDB* (2013)
14. Christoforaki, M., He, J., Dimopoulos, C., Markowetz, A., Suel, T.: Text vs. space: efficient geo-search query processing. In: *CIKM*, pp. 423–432 (2011)
15. Cong, G., Jensen, C.S., Wu, D.: Efficient retrieval of the top-k most relevant spatial web objects. *PVLDB* **2**(1) (2009)
16. Dean, J., Ghemawat, S.: Mapreduce: simplified data processing on large clusters. *Commun. ACM* **51**(1), 107–113 (2008)
17. Ding, S., Suel, T.: Faster top-k document retrieval using block-max indexes. In: *SIGIR*, pp. 993–1002 (2011)
18. Eldawy, A., Mokbel, M.F.: Spatialhadoop: A mapreduce framework for spatial data. In: *ICDE*, pp. 1352–1363 (2015)
19. Felipe, I.D., Hristidis, V., Rishe, N.: Keyword search on spatial databases. In: *ICDE*, pp. 656–665 (2008)
20. Feller, W.: An introduction to probability theory and its applications, vol. 2. John Wiley & Sons (2008)
21. Guo, L., Zhang, D., Li, G., Tan, K., Bao, Z.: Location-aware pub/sub system: When continuous moving queries meet dynamic event streams. In: *SIGMOD*, pp. 843–857 (2015)
22. Guo, T., Cao, X., Cong, G.: Efficient algorithms for answering the m-closest keywords query. In: *SIGMOD* (2015)
23. Hariharan, R., Hore, B., Li, C., Mehrotra, S.: Processing spatial-keyword (SK) queries in geographic information retrieval (GIR) systems. In: *SSDBM*, p. 16 (2007)
24. Hu, H., Liu, Y., Li, G., Feng, J., Tan, K.: A location-aware publish/subscribe framework for parameterized spatio-textual subscriptions. In: *ICDE*, pp. 711–722 (2015)
25. Li, G., Wang, Y., Wang, T., Feng, J.: Location-aware publish/subscribe. In: *SIGKDD*, pp. 802–810 (2013)
26. Lu, J., Lu, Y., Cong, G.: Reverse spatial and textual k nearest neighbor search. In: *SIGMOD*, pp. 349–360 (2011)
27. Mahmood, A.R., Aly, A.M., Qadah, T., Rezig, E.K., Daghistani, A., Madkour, A., Abdelhamid, A.S., Hassan, M.S., Aref, W.G., Basalamah, S.M.: Tornado: A distributed spatio-textual stream processing system. *PVLDB* **8**(12), 2020–2031 (2015)
28. Manning, C.D., Raghavan, P., Schütze, H., et al.: Introduction to information retrieval, vol. 1. Cambridge Press (2008)
29. Mouratidis, K., Bakiras, S., Papadias, D.: Continuous monitoring of top-k queries over sliding windows. In: *SIGMOD*, pp. 635–646 (2006)
30. Nishimura, S., Das, S., Agrawal, D., El Abbadi, A.: Mdbase: A scalable multi-dimensional data infrastructure for location aware services. In: *MDM* 2011, pp. 7–16 (2011)
31. Pripuzic, K., Zarko, I., Aberer, K.: Time and space-efficient sliding window top-k query processing. *TODS* (2015)
32. Ranjan, R.: Streaming big data processing in datacenter clouds. *IEEE Cloud Computing* **1**(1), 78–83 (2014)
33. Rocha-Junior, J.B., Gkorgkas, O., Jonassen, S., Nørnvåg, K.: Efficient processing of top-k spatial keyword queries. In: *SSTD*, pp. 205–222 (2011)
34. Sadoghi, M., Jacobsen, H.: Be-tree: an index structure to efficiently match boolean expressions over high-dimensional discrete space. In: *SIGMOD*, pp. 637–648 (2011)
35. Shraer, A., Gurevich, M., Fontoura, M., Josifovski, V.: Top-k publish-subscribe for social annotation of news. *PVLDB* (2013)
36. Shvachko, K., Kuang, H., Radia, S., Chansler, R.: The hadoop distributed file system. pp. 1–10 (2010)
37. Wang, X., Zhang, Y., Zhang, W., Lin, X., Wang, W.: Ap-tree: Efficiently support continuous spatial-keyword queries over stream. In: *ICDE*, pp. 1107–1118 (2015)
38. Whang, S., Brower, C., Shanmugasundaram, J., Vassilvitskii, S., Vee, E., Yerneni, R., Garcia-Molina, H.: Indexing boolean expressions. *PVLDB* **2**(1), 37–48 (2009)
39. Xiao, C., Wang, W., Lin, X., Yu, J.X., Wang, G.: Efficient similarity joins for near-duplicate detection. *TODS* (2011)
40. Xie, D., Li, F., Yao, B., Li, G., Zhou, L., Guo, M.: Simba: Efficient in-memory spatial analytics. *SIGMOD* (2016)
41. Yi, K., Yu, H., Yang, J., Xia, G., Chen, Y.: Efficient maintenance of materialized top-k views. In: *ICDE* (2003)
42. Zhang, D., Chan, C., Tan, K.: An efficient publish/subscribe index for ecommerce databases. *PVLDB* **7**(8), 613–624 (2014)

-
43. Zhang, D., Chan, C., Tan, K.: Processing spatial keyword query as a top-k aggregation query. In: SIGIR, pp. 355–364 (2014)
 44. Zhang, Y., Lin, X., Yuan, Y., Kitsuregawa, M., Zhou, X., Yu, J.X.: Duplicate-insensitive order statistics computation over data streams. TKDE **22**(4), 493–507 (2010)
 45. Zhou, Y., Xie, X., Wang, C., Gong, Y., Ma, W.: Hybrid index structures for location-based web search. In: CIKM, pp. 155–162 (2005)