

Distributed Ranked Data Dissemination in Social Networks

Middleware Systems Research Group Technical Report
November 2012

Kaiwen Zhang, Mohammad Sadoghi, Vinod Muthusamy, and Hans-Arno
Jacobsen

University of Toronto

Abstract. The amount of content served on social networks can overwhelm users, who must sift through the data for relevant information. To facilitate users, we develop and implement dissemination of ranked data in social networks. Although top-k computation can be performed centrally at the user, the size of the event stream can constitute a significant bottleneck. Our approach distributes the top-k computation to reduce the number of events flowing through. The solution operates on an overlay network model, such as pub/sub. Instead of forwarding all events, nodes closer to the sources can identify and forward top-k events over the collected stream. Nodes closer to the sinks can then merge the top-k streams into the final result. We use a hybrid chunk-based approach which alternates between full forwarding of events and pre-computed top-k streams at upstream nodes. This hybrid approach allows us to maintain correctness by recreating a possible top-k interleaving of the entire event stream. Our results show that our buffered solution exhibits significant traffic reduction in the overlay while maintaining latency comparable to that of the centralized approach. We also offer a sensitivity analysis which demonstrates the impact of various parameters on the performance of the solution. Experiments performed using real Twitter and Facebook datasets with 5K and 30K query subscriptions demonstrate that social workloads exhibit properties that are advantageous for our solution.

1 Introduction

The importance of social networks is prevalent with their ever growing user base, projected to soon reach the billion mark¹. Social networks (SN) at this scale introduce new challenges for both the underlying infrastructure and the perceived user experience.

Examples are the feeds produced and consumed by users and their friends' updates of news feeds, videos, application feeds, and location-based events, offers, and coupons.

¹ As of October 2012, Facebook has over a billion active users;
<http://newsroom.fb.com/content/default.aspx?NewsAreaId=22>.

Other sources of information are rooted in thousands of applications running on SNs which also generate content on behalf of their users (e.g., FarmVille). More recently, the linked data movement is becoming yet another source of publishing information [5, 4] ; for instance, social tags (e.g., keyword annotations, “like” tag) on many websites is a new way of publishing data on SNs. Social tags will (soon) be enabled on billions of web pages, which are visited by millions of users daily. The generated (and delivered) amount of information is overwhelming.

In order to improve the user experience, we aim to enable users to (globally) rank their feeds and be notified only by the top-k results of interest to them. More importantly, many of these feeds naturally have notions of proximity and locality associated to them. For example, in location-based services, people with close proximity and similar interests (eg., students on campus) constitute an ideal basis for joint top-k filtering. Similarly, traditional coupon distribution is also highly localized: People in the same area may receive similar offers, which constitutes a further opportunity to cluster similar interests based on location.

From the infrastructure perspective, the social network must float large volumes of data between its users. From the user experience perspective, large volumes of data must be adequately (pre-)processed to be consumable by users. Thus, it is essential for the user to simplify consumption of the volumes of data, which requires the capabilities to filter and rank information and to selectively pick information. This puts the burden on service providers to offer fine-grained control over the delivered information; yet, it also brings forward new opportunities to optimize the infrastructure to reduce message traffic through the adequate placement of interest queries and information routing operations towards interested users.

The underlying infrastructure that process these sheer volumes of data, which are prevalent in today’s social network domain such as [11, 22, 14, 1], is inevitably running in a distributed computing environment. Indeed, a rising trend is to establish a decentralized architecture in social networks such as Diaspora. The pub/sub model, known for its scalability and decoupled nature, then becomes a logical candidate for the dissemination substrate of distributed social networks, and pub/sub designs for social networks have recently surfaced [27, 30]. Therefore, our goal is to provide a solution for *distributed ranked data dissemination* for pub/sub systems within the context of social networks.

Our work is situated at the internal infrastructure of geographically distributed social sites (e.g., Facebook) in which an effective top-k filtering strategy, especially at the edges of the network, improves both resource usage and user experience.

Thus, in this work, we make the following contributions:

1. Present an approach to aggressively filter at the dissemination network’s edges to substantially reduce the overall network traffic (Sec. 2),
2. Develop novel algorithms that achieve early pruning of messages that are guaranteed not be part of a subscriber’s top-k results (Sec. 3), while maintaining semantical correctness (proofs: Sec. 4.9),

3. Enhancing the algorithm with adaptive features to optimize top-k dissemination from multiple sources and for multiple subscribers,
4. Propose an effective buffering protocol to avoid re-sending messages that fall in the top-k windows of multiple subscribers,
5. Propose a design that supports a rich sliding window semantics including any variations of both count- and time-based sliding windows,
6. Propose a design that makes no assumption about the scoring function used allowing each subscriber to use a scoring function of her choice,
7. Propose a design that imposes no restriction on how subscribers are allocated and assigned to nodes of the dissemination network,
8. Finally, evaluate and compare the performance of our solutions vis-à-vis a baseline algorithm, offer a sensitivity analysis of the various parameters related to top-k semantics and those specific to our solution, and employ datasets from Facebook and Twitter, containing 5K and 30K subscription queries, to project the scalability of our pub/sub architecture by identifying the major properties of social network workloads that impact our solution (Sec. 5)

2 Top-k semantics and correctness criterion

We develop a solution for supporting top-k data dissemination within the context of a pub/sub broker overlay network [8, 15]. This model allows publishers to *publish* data to a *broker* which forwards the data to consumers who are *subscribed* to the data while remaining decoupled. Brokers in the pub/sub system form an overlay and collaborate to disseminate publications from publishers to the appropriate subscribers.

Subscribers define their interest by specifying a conjunction of predicates. Publications contain a set of attribute-value pairs. Those attributes are compared against the predicates of each subscription. If a publication satisfies the predicates of a given subscriber, the publication is considered *matched* and must be delivered to the subscriber.

We extend the functionalities of pub/sub systems to allow top-k filtering on a publication stream. Subscribers include a scoring function with their subscription which assigns a numerical score to an input publication. Windows of publications are derived from the stream of matching publications according to parameters provided by the subscriber. For each window, the top-k publications with the highest score are determined and delivered to the subscriber, while the rest are discarded. Windows can be either count- or time-based. A count-based window contains a fixed number of events, while a time-based window contains publications within a certain time range.

Although the top-k stream can be computed a posteriori, our approach distributes the top-k computation to the upstream brokers as shown in Fig. 1(a). Top-k filtering is applied on partial streams collected upstream and disseminated towards a target subscriber. These top-k streams are smaller than the original streams; thus, reducing traffic within the system. These are then merged downstream to form the final top-k results.

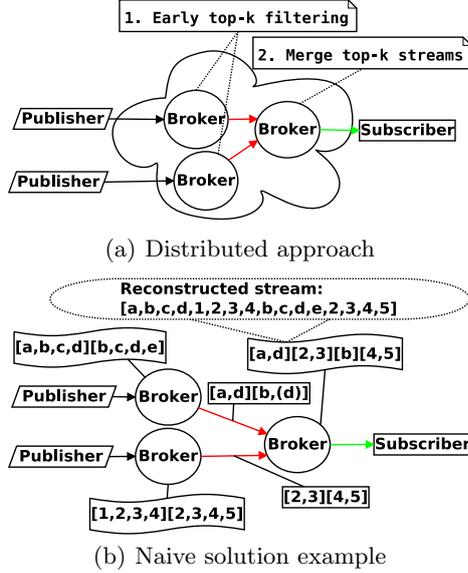


Fig. 1. Top-k matching overview

With each subscription, the subscriber specifies the following top-k parameters:

- W : The window size. For count-based windows, W is the number of elements, and for time-based windows, W is the duration.
- k : The number of highest ranking elements within a window that should be delivered.
- δ : The window shift size. For count-based windows, $\delta = 1$ results in a sliding window, $\delta = W$ is a tumbling window, $1 < \delta < W$ is hopping window, and $\delta > W$ is a “sampling” window. For time-based windows, δ is specified in units of time, and corresponding window semantics are possible.

The semantics are as follows. Let E be an ordered stream of events that match subscription s . Let \mathbb{W} be the set of windows over E according to the window semantics. For each event $e \in E$, if there is a window $w \in \mathbb{W}$, where e is one of the top k ranked events in w , then e should be delivered to the subscriber.

We suggest a correctness criterion for top-k algorithms called *stream reconstructability*: Given a set of top-k events T generated by a solution for any finite stream of events E using any window semantics WS , there exists at least one permutation (interleaving) E' of events in E possible under the specified pub/sub model where applying top-k semantics over E' using WS will produce T .

The correctness criterion states that a correct solution selects the same set of top-k events as the centralized solution, where events are all delivered to the subscriber before top-k is applied. Due to the asynchronous nature of pub/sub,

several interleavings of publications exist, depending on the reliability guarantees provided by the system. In count-based windows, a different permutation of publications may produce different results, while in time-based windows, propagation delays can result in publications being considered as part of a different time window. For the remainder of the paper, we will focus on a model where messages in the network can be arbitrarily delayed but never lost or reordered. Per-publisher order is therefore enforced: Multiple publications from a single publisher are delivered to all interested subscribers in FIFO order. According to our correctness criterion, there exists one interleaving of publications that can be processed by the centralized solution to obtain a set of top-k results which is identical to our distributed solution. Because this correctness criterion is agnostic to the scoring function, correctness is only maintained if a distributed solution considers every window from a possible interleaving of all publications. Note that the criterion does not guarantee a specific order for the correct set of top-k events.

We show a naive distributed solution for count-based windows which is incorrect: Source brokers simply filter local publications and select the top publications out of each window to disseminate downstream. Non-source brokers then simply forward any received publications to the subscriber. Fig. 1(b) is a sample execution of the naive distributed algorithm for sliding ($\delta = 1$) count-based windows. Each source broker independently selects and forwards the top-k publications out of their respective windows. For instance, the top broker forwards publications $[a, d]$ for the first window and only $[b]$ for the second window (publication d is selected again but since it has already been disseminated, it is omitted). The top-k results from each source broker are delivered in a certain interleaving to the subscriber broker.

Suppose the order is $[a, d][2, 3][b][4, 5]$. Such interleaving of top-k results can be derived from a reconstructed stream of original publications as $[a, b, c, d, 1, 2, 3, 4, b, c, d, e, 2, 3, 4, 5]$. However, applying our sliding window semantics to this stream would also require results for windows such as $[b, c, d, 1]$ and $[2, 3, 4, b]$. In a naive distributed solution, such windows are not considered since each source broker considers only windows of publications originating from their own publishers. Because it is impossible to construct an interleaving of publications where there does not exist at least one sliding count-based window which includes a publication from both brokers, this solution fails the correctness criterion.

3 Distributed chunking top-k solution

Our algorithm for efficient top-k pub/sub *dissemination* consists of three key ideas:

1. Maintain correctness using *chunking*.
2. Avoid sending duplicate messages using *buffering*.
3. Defer (adaptive) covering to further points in the network in order to exploit both distributed top-k computation and subscription covering gains.

Note that efficient (centralized) top-k *computation* by a single broker is complementary to the focus of this paper and has been thoroughly studied elsewhere [9, 26]. We are concerned with the *distribution* of the top-k computation, not the actual computation itself.

3.1 Chunking

As shown in the previous section, the naive distributed solution is incorrect because it fails to consider windows which cut across publisher boundaries. In other words, when publishers are located at different brokers, the naive distributed solution will omit certain results because certain combinations of publications coming from multiple publishers are not considered.

The intuition behind our solution is to disseminate additional publications downstream such that the broker connected to the subscriber can process those missing windows. This means that every publication found in those windows must be sent downstream. Our solution is therefore hybrid in nature: Publisher edge brokers switch between full forwarding of publications necessary to maintain correctness, and dissemination of top-k results only. Additionally, the missing windows are ones which contain publications from multiple sources. We thus need a mechanism to pick an interleaving of publications which reduces the occurrence of such windows.

We introduce the notion of *chunks* and *guards* to maintain correctness. Each publisher edge broker stream is divided into contiguous chunks. Every broker or subscriber can independently choose the chunk size (a parameter referred to as C), and can decide to change C at any time. Chunks can be count- or time-based independent of the properties of the top-k window semantics.

For each chunk C , the first and last W events must be included in the chunk and serve as *guards* (left guard and right guard, respectively). For all windows of length W within the chunk, the top-k events are computed and forwarded. Publications which are not part of any top-k or guards are discarded.

The broker may forward events in a chunk as they arrive or as a whole. In either case, the downstream broker must process the events in the chunk as a whole. In particular, it must not interleave other events within those in a chunk.

Upon receipt of a chunk, an edge broker can compute the final top-k results for its subscribers through a process called *dechunking*. Publications located within a chunk (outside of the guards) are already part of processed windows upstream and are simply delivered to the subscribers. The guards, which are complete windows of publications, are processed to compute intra-chunk windows (ie., windows which contain publications from two different chunks). Only publications which have been identified as top-k within those guards will ultimately be delivered. The example of Fig. 2 is dechunked as follows: Events from LG_1 marked as top-k events will be delivered. Then, windows w_1 to w_n will be delivered as is. Finally, top-k events will be computed and delivered for all windows with publications between RG_1 to LG_2 . The process is repeated for chunk 2 and any subsequent chunks.

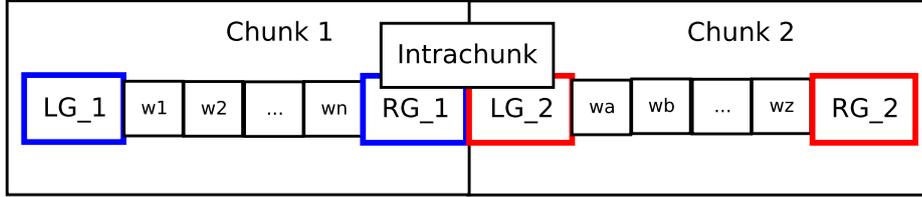


Fig. 2. Chunk structure and intra-chunk overlaps

The correctness of our solution depends on Lemma 1 and Lemma 4 (see Sec. 4.9). These two lemmas establish the stream reconstructability property: They demonstrate that our chunking algorithm creates valid interleavings of events, and that the guards provide enough publications to compute all the necessary windows.

In addition, our solution provides support for adaptive rechunking techniques which allow chunk size to change dynamically, both per broker and per subscription. Because chunks are processed atomically, a large chunk size can increase the latency since a chunk from another source may be complete but not processable until the currently selected chunk is completed. Another issue arises when a publisher stops publishing, leaving a chunk incomplete. In both cases, it is beneficial to have a mechanism to allow the chunk size to end prematurely or to be merged with another chunk (from another source) to create a larger chunk. Finally, we will also see that rechunking is useful for synchronizing chunks for different subscriptions to maximize the traffic reduction gains offered by our solution. More discussions on this aspect are available in Sec. 5.

3.2 Buffering

Due to the heterogeneous top-k specifications of the subscriptions, it is possible for the same publication to be matched at different times for different subscribers. For instance, it is possible for a publication to be a guard event for some subscription and later be matched as the top-k of a window for another subscription. This could also occur if different subscribers have different resolutions for their time-based window.

In such cases, the same publication may be forwarded down a different path through the overlay network. However, it is possible that the paths to the different subscriptions share common brokers. Each broker is equipped with standard knowledge to understand which outgoing links a publication can *potentially* be forwarded to. We require only storage of the standard subscription set and subscription last hop. This information already exists in the pub/sub broker through the operation of pub/sub routing (e.g., [15]).

3.3 Covering

A subscription S that covers another subscription S' is marked as a covering subscription. We make sure brokers do not perform top-k filtering on a covering subscription, but simply forward all matching events. To be more useful, we also need to be more flexible in whether a covered subscription should be forwarded.

A broker can determine to forward a covered subscription S' ; the result is that the upstream brokers will no longer treat S as a covering subscription and will perform top-k processing on it. By forwarding S' , the current broker effectively delegates some of its top-k processing to the upstream broker, reducing both the processing load and event traffic at the current broker.

3.4 Window semantics

In this work, we focus on two different types of window semantics: count- and window-based windows.

For the count-based windows, the following high-level procedure is employed: (1) the publisher edge broker *matches* events against subscriptions and performs *chunking*; (2) any broker may optionally perform merging of chunks: we call this *rechunking*. Optionally, any pair of neighbouring brokers may eliminate duplicates. We say that the source broker in this pair *deduplicates* the event stream, and the target broker *rehydrates* the stream. The broker pipeline is shown in Fig. 3.

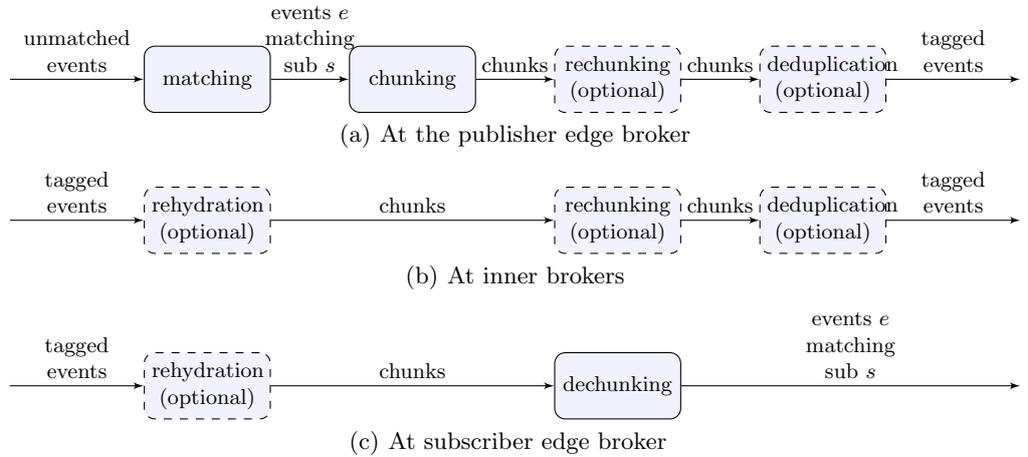


Fig. 3. Broker pipeline for count-based windows

For time-based windows, each broker maintains a local window of events and forwards the top-k within each window. This works assuming all brokers can compute time *durations* accurately. Note that the clocks do not have to be synchronized.

4 Algorithms

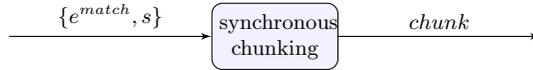
We now present the algorithm for each component in the broker processing pipeline. The chunking and rechunking components have two variants: a synchronous one that queues all events in a chunk and outputs one chunk at a time, and an incremental one that attempts to output the events with less queuing.

4.1 Matching

For each incoming event e , the matching component simply outputs a set of event-subscription pairs for each subscription s that matches e . As well, the output event is tagged as a matching event. Event tags are used in the incremental chunking algorithm.



4.2 Synchronous chunking



The synchronous chunking component, used by a publisher edge broker, outputs a chunk for a set of incoming events (Alg. 1). A chunk c has two parts: $c.sub$ refers to the subscription associated with the chunk, and $c.evts$ is the sequence of events in the chunk. This component internally maintains a data structure $s.evts$ to store a sequence of events associated with subscription s .

The $getFirstWindow(evts)$ and $getLastWindow(evts)$ functions return the set of events in the first and last windows in the sequence of events. $getAllWindows(evts)$ returns the set of windows, each containing W events, in the sequence of events. These functions encapsulate the window semantics. Note that the algorithm computes top-k of windows within a chunk, even those with events already included in a guard. If a publication is part of a top-k it will be tagged as top-k to ensure it will be delivered to the subscriber.

Note that the algorithm works even when a chunk has fewer than W events.

The algorithm does not specify when a chunk ends. This is a configuration parameter that can be used to tradeoff latency versus efficiency. One possible choice is to end chunks periodically, that is, chunk windows are time-based tumbling windows.

Algorithm 1: Synchronous chunking

```

1 on receive event  $e$  matching subscription  $s$  do
2    $s.evts.add(e)$  ;
3 on finish chunk for subscription  $s$  do
4   /* This is a higher level version of the below */
5    $c \leftarrow$  init chunk for  $s$  ;
6   add first  $W$  events in  $s.evts$  to  $c.evts$  in order ;
7   foreach window  $w$  in  $s.evts$  do
8      $\lfloor$  add top  $k$  events in  $w$  to  $c.evts$  in order ;
9   add last  $W$  events in  $s.evts$  to  $chunk.evts$  in order send chunk  $c$  to last hop
10  of  $s$  ;
11 on finish chunk for subscription  $s$  do
12   $leftguard \leftarrow \{e | s.evts.seqOf(e) \leq W\}$  ;
13   $rightguard \leftarrow \{e | s.evts.seqOf(e) > |s.evts| - W\}$  ;
14   $topk \leftarrow \emptyset$  ;
15   $\mathbb{W} \leftarrow$  getAllWindows( $s.evts$ ) ;
16  foreach window  $w$  in  $\mathbb{W}$  do
17     $\lfloor topk \leftarrow topk \cup$  getTopK( $w, k$ ) ;
18   $c.sub \leftarrow s$  ;
19   $c.evts \leftarrow leftguard \cup topk \cup rightguard$  ;
20  sort( $c.evts$ ) ;
21  forward( $c, s.lasthop$ ) ;

```

**4.3 Synchronous rechunking**

The synchronous rechunking component merges a set of chunks (Alg. 2). This component internally maintains a data structure $s.chunks$ to store a set of chunks associated with subscription s . After concatenating a set of chunks, the algorithm computes the top- k events across the chunk boundaries, namely the windows that span the guard events.

The algorithm does not specify which chunks should be rechunked, the order of these chunks, or when rechunking should occur. These are configuration parameters that can be used to prioritize certain streams, and tradeoff latency for efficiency.

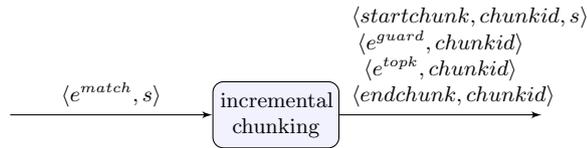
4.4 Incremental chunking

The incremental chunking component performs chunking but does not wait until the chunk is complete before forwarding events (Alg. 3). The component indicates which chunk c each event belongs to, and also tags outgoing events so that

Algorithm 2: Synchronous rechunking

```

1 on receive chunk  $c$  for subscription  $s$  do
2    $s.chunks.add(c)$  ;
3 on perform rechunk of chunks  $\mathbb{C}$  for subscription  $s$  do
4    $s.chunks \leftarrow s.chunks \setminus \mathbb{C}$  ;
5    $allevents \leftarrow \text{getEvents}(\mathbb{C})$  ;
6    $leftguard \leftarrow \{e \mid allevents.seqOf(e) \leq W\}$  ;
7    $rightguard \leftarrow \{e \mid allevents.seqOf(e) > |allevents.events| - W\}$  ;
8    $topk \leftarrow \emptyset$  ;
9   for  $i \leftarrow 1$  to  $|\mathbb{C}| - 1$  do
10     $guardevents \leftarrow \text{getRightGuard}(\mathbb{C}_i) \cup \text{getLeftGuard}(\mathbb{C}_{i+1})$  ;
11     $\mathbb{W} \leftarrow \text{getAllWindows}(guardevents)$  ;
12    foreach window  $w \in \mathbb{W}$  do
13       $topk \leftarrow topk \cup \text{getTopK}(w, k)$  ;
14   $newchunk.sub \leftarrow s$  ;
15   $newchunk.events \leftarrow leftguard \cup topk \cup rightguard$  ;
16   $\text{forward}(newchunk, s.lasthop)$  ;
    
```



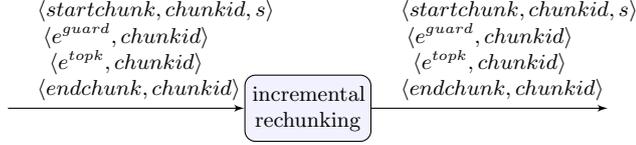
downstream brokers can distinguish between the guards and top- k events. As well, it sends messages to indicate when each chunk begins and ends.

The algorithm forwards the first W events as the left guard, and when the chunk ends, the last W events are forwarded as the right guard. As each event arrives, it computes a new window (according to the window semantics) and forwards the top k events that have not already been forwarded.

Internally, the component stores in the $s.events$ data structure the most recent two *disjoint* windows of incoming events that match subscription s . We need two windows to ensure there is one full window of events in addition to one window of right guard events in case the chunk ends soon. As well, $s.sentevents$ stores the events the components has already forwarded, and $s.chunkid$ is the unique ID of the current chunk for subscription s .

4.5 Incremental rechunking

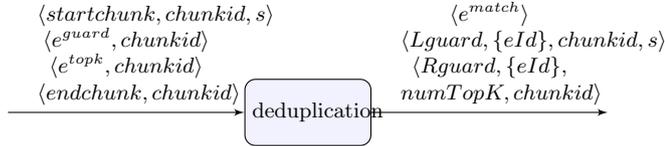
Since we are not allowed to interleave events across chunks, rechunking needs to wait for complete chunks. Therefore, the incremental rechunking algorithm is essentially the same as the synchronous variant but the outgoing events are individually tagged as guard or top- k events.



4.6 Dechunking:

This component, used by a broker for subscribers directly connected to it, computes the final top-k stream by concatenating chunks, computing and forwarding top-k results over guards and forwarding pre-computed top-k results. It functions similarly to rechunking.

4.7 Deduplication



The deduplication and rehydration components operate as a pair to avoid sending duplicate messages between brokers. The deduplication component operates at the tail end of an upstream broker’s pipeline and essentially “compresses” the outgoing chunked event stream. The rehydration component works at the head of a downstream broker’s pipeline and reconstructs the original chunks. To clarify, rehydration restores deduplicated chunks and may be performed at any broker while dechunking refers to the reconstruction of the final top-k results at the subscriber edge.

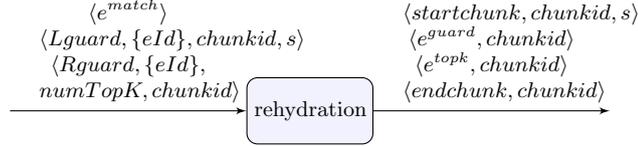
The deduplication component maintains a history of the events forwarded to each neighbour and never forwards duplicate messages over the same overlay link (Alg. 4). To ensure the downstream broker can reconstruct the chunk stream, the deduplication component forwards the IDs of the guard events; this is relatively lightweight information but must be done for each chunk. The top-k events that fall between the guards, however, can be inferred by the rehydration component at the downstream broker. The details of this algorithm are described below in the discussion about the rehydration component.

Internally, the deduplication component maintains information about the incoming chunk stream for each subscription s . In particular, $s.chunkid$ is the current chunk ID, $s.guardids$ are the event IDs of the guard events, and $s.numTopK$ is the number of top-k events in the chunk.

The algorithm below assumes the previous component in the pipeline is the incremental rechunking component, but it is easy to adapt it to work with the chunking or synchronous rechunking component.

Note that deduplication always propagates events immediately, but the corresponding rehydration phase does queue events and increases event delivery latency.

4.8 Rehydration



The rehydration component reconstructs the chunk streams (Alg. 5). For each incoming event, it first finds the matching subscriptions and records the event in a buffer associated with each subscription. This buffer is the sequence of *potential* events within each subscription’s chunk. The guards within each chunk are explicitly specified by the upstream broker’s deduplication component. To determine the top-k events within each chunk, it waits for the count n of top-k events within the chunk (sent within the right guard message *Rguard* from the deduplication phase), and selects the n highest ranked events in the buffer within the guards.

Note that rehydration is the only component that requires matching events against subscriptions after the initial matching at the publisher edge broker. Therefore, while deduplication and rehydration reduce message traffic, they increase the processing overhead at the brokers.

4.9 Proofs

In this section, we prove the correctness of the chunking technique. In particular, we show that the chunking, rechunking, and dechunking components do not violate the top-k semantics defined in Section 2.

Recall that the semantics are defined on a computation of top-k windows over the stream of events received at a subscriber. Since the network only guarantees per-publisher ordering of messages, and messages may be arbitrarily delayed, there are many *possible interleavings* of events at the subscriber.

With chunking, top-k computations are performed on a partial stream of events before they arrive at the subscriber, in fact before it is even known what the complete event stream at the subscriber will look like. The proof below will argue that performing top-k computations on these partial streams is equivalent to a top-k computation over a stream that represents a *possible interleaving* of events at the subscriber.

The proofs exploit the requirement that chunks must have $n\delta$ events, where δ is the window shift size, and $n > 0$. Each chunk then has n windows. Windows

whose events all appear within the chunk are referred to as intra-chunk windows, and windows whose right boundaries extend beyond the chunk are called inter-chunk windows. Note that a chunk can have at most $\lceil W/\delta - 1 \rceil$ inter-chunk windows, which must all appear as a contiguous set at the end of the chunk. In all the notation below, items are numbered from one: the first event is a chunk or window is e_1 , and the first window in a chunk is W_1 .

Lemma 1. *The interleaving of chunks will reconstruct a possible interleaving of events.*

Proof. From our assumptions, the network preserves per-publisher message order, but may arbitrarily delay messages, so it remains to show that the interleaving of chunks preserves per-publisher ordering.

Consider the concatenation of two chunks C and C' , which generates a reconstructed stream of events that contains two events e and e' from the same publisher, such that e was published before e' .

Now, if e and e' belong to the same chunk, then e must appear before e' , since event ordering is preserved within a chunk. On the other hand, suppose without loss of generality that e belongs to chunk C , and e' to C' . Since e was published before e' , chunk C must have been generated before C' . For a broker to concatenate C' before C , either it received C' and C from different brokers, or it received C' before C from the same broker. The former case is impossible since e and e' originated at a single publisher, and by the acyclicity of the network, there must be a single path from this publisher to the current broker. The latter case is also impossible: since e was published before e' , chunk C must have been generated before C' . Given that C and C' traversed the same path to the current broker, and brokers along the path process chunks for a given subscription from the same upstream broker in FIFO order, a broker must receive C before C' .

Therefore, e must appear before e' in the reconstructed stream, preserving per-publisher ordering.

Lemma 2. *The windows considered over the events within a chunk correspond to windows on the reconstructed stream of events.*

Proof. Let chunk C contain a stream of $n\delta$ events $e_1, \dots, e_{n\delta}$ and chunk C' contain a stream of $m\delta$ events $e'_1, \dots, e'_{m\delta}$, where δ is the window shift size, and n and m are positive integers. By concatenating chunks C and C' , we get the reconstructed stream of events $e_1, \dots, e_{n\delta}, e_{n\delta+1}, \dots, e_{n\delta+m\delta+1}$, where for convenience event $e_{n\delta+i}$ is an alias for e'_i for $i > 0$.

There are n windows on chunk C , where the first event in window W_i is event $e_{(i-1)\delta+1}$, for $0 < i \leq n$. Likewise, there are m windows on chunk C' , where the first event in window W'_j is event $e'_{(j-1)\delta+1}$, for $0 < j \leq m$. The reconstructed stream, having $(n+m)\delta$ events, has mn windows, where the first event in window W''_k is event $e_{(k-1)\delta+1}$, for $0 < k \leq mn$.

We now show that each window W''_k on the reconstructed stream corresponds to a window in chunks C or C' , that is, W_i or W'_j , respectively.

In chunk C , window W_i starts at event $e_{(i-1)\delta+1}$, and hence corresponds to window W_i'' which also starts at the same event according to the definitions of W_i and W_k'' above. By a similar argument, for chunk C' , window W_j' starts at event $e'_{(j-1)\delta+1} = e_{(n+j-1)\delta+1}$, and hence corresponds to window W_{n+j}'' .

Since, each window in each chunk corresponds to a window in the reconstructed stream, and the number of windows in the reconstructed stream is the sum of the number of windows in the two chunks, each window in the reconstructed stream corresponds to a window in one of the two chunks.

The above reasoning can be extended to an arbitrary number of chunks.

Lemma 3. *The top- k events computed over intra-chunk windows are the top- k events in the reconstructed stream. Similarly, events that are not in the top- k in intra-chunk windows are not among the top- k events in the reconstructed stream.*

Proof. From Lemma 2, we know that intra-chunk windows in a chunk correspond to windows in the reconstructed stream. Moreover, since intra-chunk windows always contain W events, the top- k events found over these windows exactly match those in the reconstructed stream.

Lemma 4. *All events contained in an inter-chunk windows are guard events.*

Proof. In a chunk with p events, for window W_i to be an inter-chunk window, the first event in the window must be event e_j where $p - W + 1 < j \leq p$. Otherwise, if $j < p - W + 1$, the last event in the window will be e_{j+W-1} which appears in the chunk, and hence W_i will not be an inter-chunk window.

We require that the first and last $W - 1$ events in a chunk be guard events, so all events in W_i will be guard events.

Theorem 1. *The distributed chunking algorithms compute top- k results that are possible were top- k computed locally over the stream of events received at the subscriber.*

Proof. From Lemmas 1 and 2, we know that the windows computed over the chunks are aligned with windows over a possible stream of events that a subscriber may observe. From Lemma 3, we know that the top- k results computed over intra-chunk windows are correct. Finally, by Lemma 4, we know that the remaining (inter-chunk) windows consist of guard events that are propagated to the subscriber, where the inter-chunk window top- k results can be computed.

5 Evaluation

We experimentally evaluate our various count-based algorithm implementations. Our experiments are divided into two parts. The first part employs a synthetic workload and contains a performance analysis of the following algorithms:

1. Baseline where all computations are performed at the subscriber's edge broker (EDGE),

2. Incremental chunking (CHUNK),
3. Incremental chunking with the buffering deduplication/rehydration option (B-CHUNK).

We also perform a sensitivity analysis of the B-CHUNK solution.

The second part evaluates the B-CHUNK algorithm within the context of the online social network use case and is performed using real datasets.

5.1 Performance and sensitivity analysis

Setup: The algorithms are implemented in Java for the PADRES pub/sub prototype². Experiments are conducted on the SciNet testbed in the General Purpose Cluster (GPC) using 24 machines³.

The workload used for the baseline comparison and the sensitivity analysis is synthetic: Publishers send data every 4 seconds. To maximize the top-k overhead at that scale (worst case scenario), the subscription predicates are identical and match every publication sent. However, the scoring function associated with each subscription can differ, as described below. The workload for the online social networks use case is described separately in Sec. 5.4.

The overlay topology consists of several core brokers connected in a chain. In addition, these core brokers are each connected to 5 additional edge brokers. Publishers are located on edge brokers called source brokers, while subscribers are uniformly distributed on the rest. This setup models a network of data centers connected through gateways and focuses on measuring the impact of our algorithms on delivery paths with multiple broker hops.

A score is assigned to each publication for each subscriber, generated using a certain distribution function. The top-k consists of the k publications with the highest scores within the current window. Since we are focused on the dissemination aspect, we minimize the overhead of the scoring function. Thus, the overhead of the scoring function is not a bottleneck in our experiments, but could be a factor if a real function was used. We also employ a deterministic scoring function (DET) which assigns strictly decreasing scores over time. This deterministic scoring function is used when we want different subscribers to assign the same score to the same publications.

Parameters: The various window related parameters, such as W , k and δ vary during the experiments. In particular, we are interested in a sliding window ($\delta = 1$) and a tumbling window ($\delta = W$). We also vary the number of subscribers, publishers, and the top-k distribution. In a uniform distribution, each publication is equally likely to be part of the top-k of subscription, whereas in a Zipfian distribution, the top-k of each subscription are more overlapping. Finally, we also vary the publication delay. In a *mixed* setup, publications are injected as subscribers join the system. In a *2-phase* setup (2P), subscriptions are first submitted before publications start to flow.

² <http://padres.msrg.toronto.edu/>

³ <http://www.scinet.utoronto.ca/>

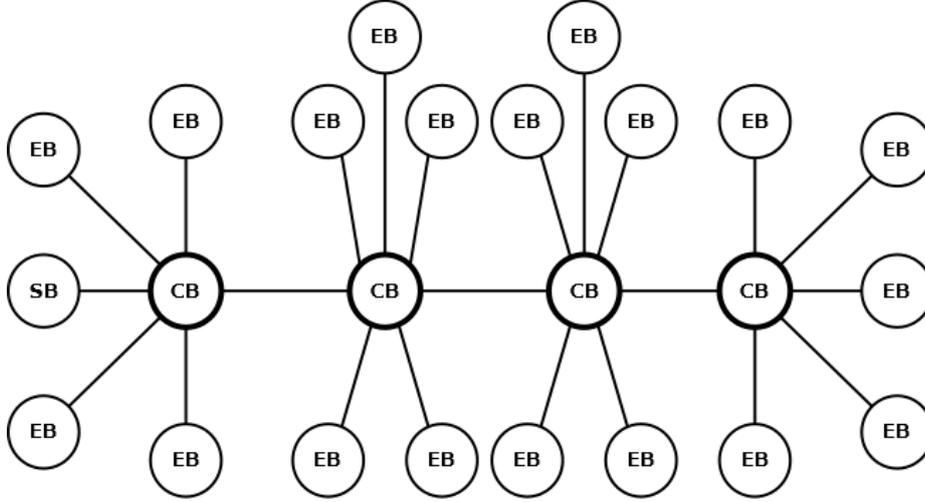


Fig. 4. Broker overlay topology

Metrics: *Outgoing traffic* - The main advantage of early top-k matching is to allow unselected publications to be discarded early, which reduces the amount of traffic in the system. We are therefore interested in the traffic reduction of our distributed solutions.

End-to-end latency - Although the matching overhead associated to top-k is negligible in our evaluation, publications can be deferred until they are selected as part of the top-k. This delay is dependent upon the window parameters used. We measure the end-to-end latency at the subscribers which are the furthest away from the publishers. In our topology, publishers are 6 broker hops away from the measured subscribers.

5.2 Performance Comparison

We measure the performance of the various solutions presented over an increasing number of subscribers. Traffic throughput have been measured at different points in the network, namely at the source (publisher) broker, core brokers and edge (subscriber) brokers. We employ a single publisher and each subscriber is interested in the top-k of every publications published using a deterministic scoring function. We use $k = 1$, $W = \delta = 20$, and $C = 100$ (chunk size). These top-k parameters are applied to every subscription.

Fig. 5(a) shows the relative throughput at the publisher broker of each solution, normalized relative to the EDGE solution. This solution always forwards every publication received since no top-k processing occurs at the source. The chunking solution without buffering performs top-k filtering at the source; however, the top-k publications are sent separately for each subscription. Thus, a subscription could be sent multiple times over the same link. This duplication

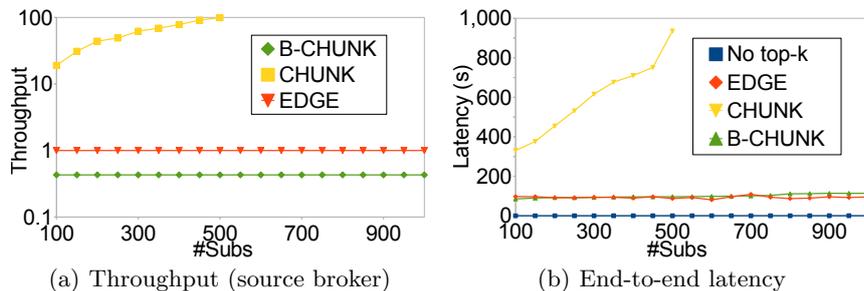


Fig. 5. Performance comparison

means the traffic increases linearly to the number of subscribers. In our experiments, the throughput of CHUNK starts at an order of magnitude higher in traffic to that of EDGE. In fact, the throughput is saturated at around 500 subscribers. With buffering enabled, the chunking solution (B-CHUNK) scales much better than the centralized solution. B-CHUNK shows a constant 57% traffic reduction with varying subscriber loads. In our setup, a large portion of the remaining traffic is used to forward guard publications.

At the subscriber brokers, the throughput results are straightforward: Each solution reduces the traffic by 95%. Since the selectivity of the top-k filter is 1 out of 20, only 5% of the publications are ultimately sent to the clients. No further processing is required by the clients themselves: The processed stream sent by the edge broker to a client corresponds exactly to the final result. Therefore, all three solutions are able to produce that stream at the edge broker.

Compared to EDGE, core brokers will receive more inbound traffic for CHUNK, since the source broker generates more traffic. Each publication received is then forwarded down a single link to the corresponding subscriber. Thus, the outgoing traffic is exactly equal to the inbound traffic. Since core brokers have 6 outgoing edges, each publication is forwarded to only 16.67% of the links. However, the duplication of publications at the source broker causes core brokers to forward a publication down a specific link multiple times to satisfy different subscribers. Therefore, CHUNK generates more traffic at the core brokers than EDGE, even if individual inbound publications are forwarded down only a single link. For B-CHUNK, every publication received is forwarded down every link, since every subscriber uses the same scoring function. The reduction in traffic at core brokers comes only from the fact that there is less inbound traffic from the source broker.

We compare the end-to-end latency between the various solutions (see Fig. 5(b)). Because top-k is a stream processing operation, publications can sit in a queue waiting for the window to fill up. In a count-based solution, the delay experienced due to queuing can vary widely between individual publications. In our experiments, we control the variation by using a tumbling window with deterministic scoring. A single publication is sent for every 20 publications received. The sent

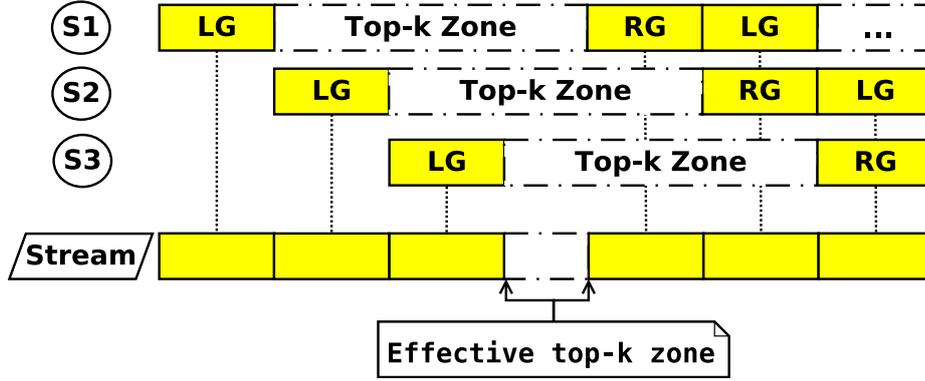


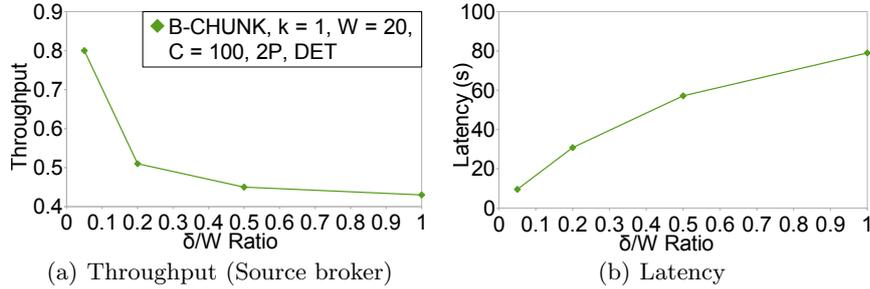
Fig. 6. Chunk guards offset

publication is the one which entered the window first, since publications have monotonously decreasing scores.

For CHUNK, the latency is dependent on the number of publications and performs much worse than other solutions. This is because that traffic is duplicated per each subscriber. For B-CHUNK and EDGE, the latency is dominated by the queuing time (as seen by the difference when no top-k is used). The queuing time is mostly dependent upon three factors: The scoring function used, the window size, and the publication rate. All three of these factors affect both solutions equally, hence, there is no advantage in choosing one solution over another with respect to latency. Thus, the overall latency is comparable between B-CHUNK and EDGE. The difference is that B-CHUNK queues publications mostly at the publisher brokers, while EDGE queues at the subscribers' edges.

As explained before, the scoring function used in the experiment incurs minimal overhead. Therefore, it does not factor in our latency measurement. However, the matching overhead is distributed differently depending on the solution used. In B-CHUNK, the source brokers must process a single publisher stream for multiple subscribers, while in EDGE, the opposite is true. Thus, a workload which contains more publishers than subscribers will distribute the matching overhead more evenly using B-CHUNK.

Summary: The CHUNK solution is simply not scalable due to the publication duplication issue and always performs worse than the EDGE solution in both throughput and latency even for small number of subscribers. B-CHUNK demonstrates lower traffic than EDGE due to the source filtering of top-k publications. Latency for both solutions is comparable since it is dominated by the queuing time, which is the same.

Fig. 7. δ/W ratio analysis

5.3 Sensitivity analysis

The parameters used in Sec. 5.2 can be considered optimal in terms of traffic reduction for the B-CHUNK algorithm. We now evaluate the impact of each parameter and justify its effect within the context of our solution.

Shift: The shift parameter has a significant impact on the throughput at the core brokers. With a tumbling window, the throughput flows at a constant rate to the subscriber: For every W publications, k publications will be sent to the subscriber. For a sliding window, the rate can vary depending on the actual score of the publications. In the best case, a publication can dominate up to $W - 1$ publications. For instance, for $k = 1$, a publication p which has a higher score than the next $W - 1$ publications following it will only trigger one message (used to forward p). It is only when that publication falls off the window (at the W^{th} publication following p) that a new top- k publication will be selected. In the worst case, a publication can be sent out for every incoming publication if the top- k is always at the head of the window. At every shift, one of the current top- k publications fall off the window and is replaced by another which must be forwarded. In this case, once the window is initially filled, maximal throughput is achieved: Every incoming publication generates an outgoing publication.

We evaluate the impact of the shift in the worst case scenario, where the scoring function used generates progressively decreasing scores. Fig. 7(a) shows the impact of the shift parameter on the throughput. With a window size of 20, the first data point corresponds to a sliding window shift of 1 (ratio of 0.05). Except for the first window of each chunk, the throughput is maximized and there is no traffic reduction in those subsequent windows. At the other extreme, the tumbling window (ratio of 1.0) provides the best traffic reduction at 43% of the baseline traffic at the source broker.

The throughput is inversely proportional to latency (see Fig. 7(b)). A smaller shift will create a more frequent publication rate, which minimizes the queuing time. A longer shift empties the window more quickly, which delays the next top- k computation until the buffer is filled. We observe that this particular tradeoff is not specific to our solution, but rather comes from the window semantics

themselves. A workload with decreasing publication scores is sensitive to the throughput issue in sliding windows, but will result in shorter publication delays.

Publication delay: We evaluate B-CHUNK with varying publication delay. Publication delay refers to the time elapsed before publications are released in the system. In a *mixed* setup, publishers start their publications while subscribers join the system. This means that earlier subscribers will receive more publications than those who join later. In a *2-phase* setup, we first disseminate all the subscriptions before sending publications. Subscribers thus receive the same stream of publications when using the 2-phase setup in conjunction with the deterministic scoring function.

Publication delay has a major impact on the throughput of our B-CHUNK solution. This is due to the fact that publication forwarding is a binary decision: A broker **must** forward a publication down a certain link if at least one subscriber requires this publication at that path. For the source brokers, this means a publication can be dropped only if **all subscribers** do not require that publication. This occurs when the publication is not selected as part of any top-k window of all subscribers **and** the publication is not part of any guards of any chunks. In a mixed setup, the subscriptions will start their first chunk at different times, which correspond to the first publication each subscription matches. This means that chunks belonging to different subscriptions can offset one another such that guards for one subscription will overlap into the top-k zone of another subscription. When this happens, publications in that overlap must be forwarded by the source broker to satisfy the guard requirements, even if they are not part of any top-k window computed by other subscriptions. This is illustrated in Fig. 6. The source broker can only drop publications which are located only in the *effective top-k zone*, where no subscription is in a guard. In the worst case, a publication stream can be completely covered by guard zones, in which case every publication must be forwarded regardless of the actual top-k computations being performed. In the best case, guard zones are completely disjoint with each others' top-k zones, which maximizes the effectiveness of the top-k computations.

In our mixed setup, subscriptions are submitted every second (publication rate is once per 4 seconds). At the source broker, the throughput is already degraded after 50 subscribers and is saturated at 100 subscribers and more (see Fig. 8(a)). At the core brokers (Fig. 8(b)), there is stable throughput reduction of 20% at the core brokers. Since core brokers have multiple outgoing edges, it is less likely for subscriptions down a certain link to offset each other. In our solution, there are 50 subscribers per outgoing edge, which is not enough to saturate the throughput at the core brokers, even in a mixed setup.

Nevertheless, our solution remains sensitive to offsets in subscription chunks due to the presence of guard zones which can nullify top-k filtering. It is therefore particularly useful to employ rechunking techniques to synchronize chunks such that their guards overlap each other, maximizing the effective top-k zones.

Scoring function: We evaluate the impact of the scoring function used. In our random selection tests, subscribers select different top-k streams according to

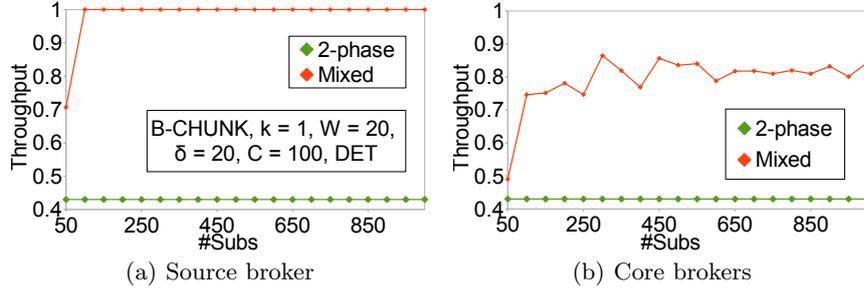


Fig. 8. Throughput vs. publication delay

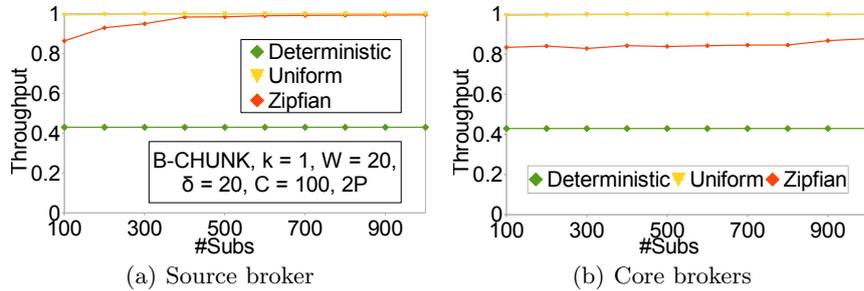


Fig. 9. Throughput vs. scoring function

the score they individually assign to each publication. Again due to the binary nature of publication matching, the throughput of our solution is sensitive to the scoring function used. Scoring functions that produce disjoint sets of top- k publications result in a higher throughput since publications matching at least one subscriber must be forwarded. In contrast, similar scoring functions can be leveraged to reduce traffic by forwarding publications which will match multiple top- k windows.

We simulate subscriptions with similar interests using a Zipfian scoring function. A random global score is assigned to each publication and each subscription derives its personal score for the subscription by modifying the global score with a random value chosen from a Zipfian distribution. Thus, subscriptions will have similar scores for each publication. In contrast, the uniformly random scoring function assigns an arbitrary value for each publication per subscription.

Fig. 9(a) shows the relative throughput of the various scoring functions at the source broker. The deterministic scoring function is the most efficient since all the subscribers' top- k publications are the same. Zipfian scoring provides some clustering of top- k publications, which provides a 17% traffic reduction (with respect to EDGE) at 100 subscriptions. However, the throughput saturates beyond 400 subscriptions. The uniform scoring function does not provide any benefit at the source broker, even with a small number of subscriptions.

At the core brokers, we observe that the throughput stays constant with respect to the number of subscribers in all cases (see Fig. 9(b)). This is due to the fact that each outgoing link for a core broker has at most 50 subscribers. Therefore, a core broker can safely eschew forwarding of a publication down a specific link if none of those 50 subscribers require it. This is advantageous for the Zipfian scoring function, which provides a sizable reduction in traffic when the downstream link has a small number of subscribers (e.g., 50).

In general, subscriptions with similar scoring functions should be clustered together. While the source brokers will not benefit from such clustering, core brokers will as they skip forwarding down unmatched links.

Number of publishers: We increase the number of publishers and test various deployment strategies. In a *cluster* strategy, we increase the number of publishers at a single source broker, whereas in an *uniform* strategy, the publishers are uniformly distributed around the topology. We keep the rate of publications constant (15 pubs/min), regardless of the number of publishers.

We first observe that the various algorithms are not sensitive to the number of publishers with regards to throughput. Using the same setup as the performance evaluation, but varying the number of publishers, the same pattern is followed: B-CHUNK uses 43% of the original traffic at the source broker, while CHUNK performs increasingly worse as the number of subscribers increase. This pattern holds true no matter how the publishers are distributed.

This result is justified by the distributed and decoupled nature of our B-CHUNK algorithm. When multiple publishers are connected to the same source broker, all publishers' event streams are treated together as a single merged event stream for that source broker. Because the scoring function is deterministic and is based on the order of publications, the top-k computation produces the same results. When multiple publishers are connected to different brokers, each broker locally performs top-k computations independently. Thus, they will behave and perform identically.

For end-to-end delay, we find that the placement strategy used, rather than the actual number of publishers, is a major factor. Fig. 10 shows that in a clustered strategy, the latency remains constant over an increasing number of publishers. This is due to the decoupled property mentioned previously: The output events from the different publishers are merged into a single input stream for the source broker, which is no different from a single publisher publishing at their aggregated rate. When publishers are located on different source brokers (uniform strategy), the latency grows exponentially in the number of publishers. Because the publication rate is kept constant, a greater number of publishers also imply that each publisher will publish at a reduced rate. Since each publisher is publishing at a different broker, each window fills up at a slower rate, thus increasing the queuing time for publications. Furthermore, the edge broker at a subscriber's end will process one chunk at a time, originating from a single source broker. The edge broker will therefore select one chunk to process and buffer the other chunks. Since each source broker is filling its chunk concurrently, as soon as the selected chunk is completed, the edge broker will quickly process

the completely buffered chunks before selecting a new chunk among the batch of newly formed chunks. This means the queuing time of publications in buffered chunks also includes the time required to finish the currently selected chunk. This phenomenon motivates the use of publisher clustering techniques.

Summary: The shift parameter, the publication delay and the scoring function are three major performance factors for B-CHUNK. Selecting a small shift value can lead to high throughput in certain workloads. However, one must be aware of the throughput-latency tradeoff surrounding the shift parameter. Offset chunk guards can dominate top-k filtering and saturate the throughput. Without chunk synchronization, a mixed workload will not scale as the number of subscriptions increases and the likelihood of guards offset to occur. Synchronizing chunks is necessary in order to correct the offsets. The scoring function can vary the amount of overlap between the top-k sets of the subscriptions. As the number of subscribers increase, the publications to be forwarded is likely to increase. Eventually, every publication must be forwarded by the publishers. However, clustering subscriptions with similar scoring functions in the same parts of the overlay will yield benefits for core brokers with multiple links.

With regards to latency, The chunking solutions are sensitive to the number of source brokers and not the number of publishers. As the number of source brokers increases, the number of concurrent chunks also increases, with subscribers only able to process one chunk at a time. This provides an incentive to cluster publishers to reduce the number of source brokers, as well as reducing the chunk size to allow the subscribers to switch between chunks faster.

5.4 Online social networks use case

We evaluated our solution in the context of online social networks. We employ workloads extracted from public datasets from Facebook [29] and Twitter [19].

Subscriptions: The datasets contain the social relationships between the different users in the network. In Facebook, the *friendship* relation is modeled as a pair of subscriptions which subscribes one user to another and vice-versa. In Twitter, the *follow* relation is modeled as a subscription between a *follower* to a *followee*. All these subscriptions are topic-based: Subscribers are interested in the totality of publications posted by a user. From the original datasets, we extracted a sample set contained subscriptions for 1000 users, using the technique employed in [27]. The samples retain the original properties of the original datasets, specifically with regards to the relative popularity of the topics. Fig. 11 shows that the popularity of both datasets exhibit a long tail, with the Twitter tail being broader than Facebook. The Facebook and Twitter samples contain 5K and 30K subscriptions, respectively. For each user, we combine all subscriptions together to form one subscription spanning the various topics this user is interested in. This allows our pub/sub system to perform top-k over the entire stream of matching publications for that given user, rather than compute the top-k publications over each topic separately. The intended purpose of our solution is to provide subscribers with a top-k selection of the most relevant

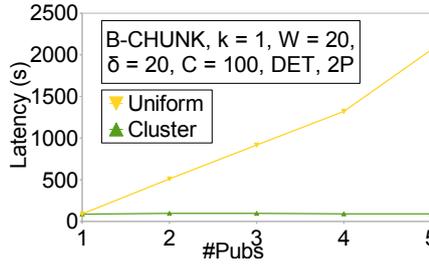


Fig. 10. End-to-end latency vs. publisher placement

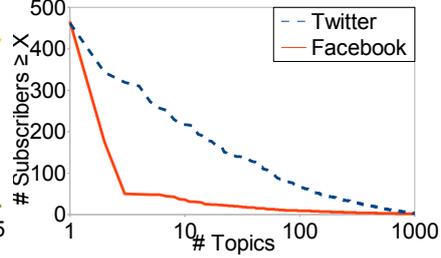


Fig. 11. Social popularity tail distribution

publications among all the topics this user is subscribed to. Thus, our sample actually contains only 1K subscriptions, one for each user.

Publications: Publications are generated synthetically at a fixed rate of one per second. The publications represent *tweets* (Twitter) and *status updates* (Facebook) posted by users ie., topics, which must be disseminated to all subscribers of the user eg., followers and friends. The publications are assigned a topic selected separately for each publication. According to [17], there is a correlation between user popularity and activity. We take this into account by skewing topic selection towards popular users: The probability of a topic being selected for a publication is linearly dependent on the topic popularity.

Scoring function: Publications are scored according to the popularity of its source (namely, its topic). A publication coming from a popular topic will gather a higher score than an unpopular topic and thus be more likely to be part of the top-k scores for a window. We argue that since our sample subscriptions come from a cluster of linked users, this popularity-based scoring function is akin to collaborative filtering, which is employed in social networks.

Topology: We use a topology similar to the previous experiments, with edge brokers connected to a chain of core brokers. Each edge broker is connected to a single publisher. We argue that in a real social network cluster, state updates are handled by an external service which is logically represented as a single source of data for our pub/sub system. Each edge broker is connected to a single subscriber, who is in charge of submitting subscriptions for several users. The subscriptions are allocated randomly to each subscriber. Each pub/sub subscriber logically represents a collection of users that reside on the same machine. Furthermore, even though the subscriptions are assigned randomly, the subscriptions in the sample data are closely related and follow the aforementioned long tail distribution: Subscriptions residing on the same physical machine have similar interests, which adheres to the clustering techniques prescribed by [24]. Clustering the interest of multiple users to a single subscriber is an effective technique for scaling up a pub/sub-based social network dissemination middleware to the required specifications.

Top-k parameters: We use $k = 2$, $W = 20$, $\delta = 20$ (tumbling window), and $C = 100$ (chunk size), with a 2-phase setup. The algorithm employed is B-CHUNK (chunking with deduplication/rehydration).

Performance evaluation: Fig. 12(a) shows the normalized throughput of the B-CHUNK solution over the centralized baseline for varying number of subscriptions at the publisher edge broker. For Facebook, the reduction gain starts at 27.6% for 100 subscribers and deteriorates to 6.7% for 1000 subscribers, while in Twitter the reduction starts at 16.9% and drops to 0.5%. Two main properties of the workload justify these results.

First, the subscriptions have varying degrees of overlap in the publication space. This means that the windows of different subscriptions progress at different rates since they are not filled by every published event (as in our sensitivity analysis). This causes offsets in the chunks and reduces the size of the effective top-k zone, as described earlier. Due to this phenomenon, the results are worse than during the performance evaluation using synthetic data. Synchronizing guards is therefore a concern as it has practical benefits in real workloads.

Second, the skew towards popular subscriptions is beneficial to our solution. The publication workload contains a higher frequency of popular publications, which are matching a large number of subscriptions. Therefore, it is likely that for any given window of any subscriber, it contains some publications of popular topics. Due to our scoring functions, these publications are then able to “overshadow” any lower ranking publications and be selected as part of top-k. Thus, any subscription that overlap in those popular topics will likely obtain the same top-k results. This allows the publisher broker to safely filter out a large volume of publications very quickly and only keep the higher ranking ones which appear in the top-k list of a large number of subscriptions. It also follows that having a long tail, as exhibited by social network workloads, is a benefit as it allows the solution to “cut” the tail at the publisher broker and keep only a smaller core of publications. This is demonstrated in the results, where the Facebook experiments perform better than Twitter due to their narrower tails. Due to this positive effect, our results are better than the measurements obtained during the publication delay analysis, even though we are still under the chunk offset issue.

The positive gain of the popularity skew is more apparent in the core brokers performance (see Fig. 12(b)). Because the core brokers have multiple outgoing hops, the chunk offset occurs only between subscription from the same hop. The problem is therefore lessened, and we see the traffic reduction jumps 48% and 38% at 1000 subscribers for Facebook and Twitter. This means our popularity-based social network experiment outperforms the Zipfian scoring function used in the scoring sensitivity analysis and is closer to the uniform scoring performance.

Scalability tests: We also evaluate the impact of the number of brokers on traffic reduction. We extend the topology to 96 and 960 brokers while maintaining the degree of each node constant. For instance, a 960 brokers topology consist of 160 core brokers, each connected to 5 edge brokers. For the publisher edge broker, we find no impact for the number of brokers in the system. The traffic reduction is found to be the same: 0.5% for Twitter and 6.7% for Facebook

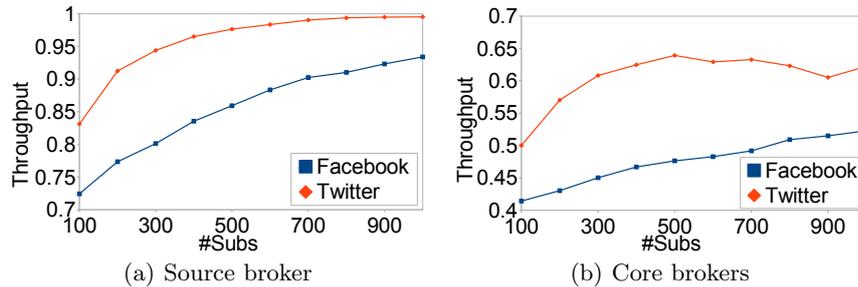


Fig. 12. Throughput vs. number of subscriptions

at 1000 subscribers. This is because the algorithm is lightweight: The publication edge broker computes the top-k of each subscriber locally and forwards the resulting publications downstream, irrespective of the rest of the topology.

For the core brokers, we find the size of the topology to have no impact on the performance as well, other than the fact that the same number of subscriptions is now spread out on a larger number of edge brokers. Thus, the first core broker directly attached to the publisher edge broker is in charge of disseminating for all 1000 subscriptions, and will therefore exhibit traffic reduction in the range of 48% and 38% for Facebook and Twitter, as per Fig. 12(b). Core brokers towards the end of the chain have a smaller number of subscriptions to disseminate to, and reduce traffic at a higher rate.

Summary: The popularity-based scoring function and publication workload, coupled with the heavy tail distribution of subscriptions in social networks, creates large overlaps in the top-k results of the subscribers which increases the traffic reduction capabilities of our solution, especially apparent at the core brokers. On the other hand, the uneven overlaps between the subscriptions create offsets in the chunks which minimize the effective top-k zone, which indicates the solution can benefit from chunk synchronization optimizations (eg., rechunking).

The scale of the topology has no direct impact on the performance of the solution. The only benefit gained by increasing the number of brokers is to distribute the subscription load onto a larger number of subscriptions, which is the main factor for our solution. We envision our pub/sub model to scale adequately through partitioning of the social networks according to interests.

6 Related Work

Generally speaking, the problems related to retrieving the most relevant answers have been studied in different contexts including database (distributed) top-k querying ([13, 18, 3, 7, 21]) and publish/subscribe (pub/sub) matching techniques [31, 2, 12, 6, 28, 16, 25].

The most widely adopted database top-k processing model [13, 18] differs with our proposed top-k model in an important respect: Our top-k model solves

the reverse problem. In the database context, top-k querying means finding the most relevant tuples (events) for a given query (subscription). But in our pub/sub abstraction, matching means finding the relevant subscriptions (queries) for a given event (tuple). Furthermore, the database literature is best suited for lower dimensional data [18], while in pub/sub context, a dimensionality in the order of hundreds is commonplace [31, 2, 12, 6, 28, 16, 25, 26], which is orders of magnitude larger than capabilities of existing database techniques.

Broadly speaking, two classes of matching algorithms have been proposed for pub/sub: Counting-based [31, 12, 28] and tree-based [2, 6, 25, 26] approaches. A fresh look at enhancing pub/sub matching algorithms is to leverage top-k processing techniques to improve matching performance. An early top-k model is presented in [20]; however, this model leverages a fixed and predetermined scoring function, i.e., the score for each expression is computed independent of the incoming event. In addition, this approach is an extension of *R-Tree*, the interval tree, or the segment tree structure; as a result, it is ideal for data with few dimensions [20]. In contrast, a scalable top-k model that supports up to thousands of dimensions while incorporating a generic scoring function, i.e., takes the event into consideration, is introduced in [28], which relies on a static and single-layered pruning structure [28]. To alleviate these challenges, a new dynamic and multi-layered pruning top-k structure is developed in [26]. However, our proposed top-k model attempts to solve a different problem, namely, distributed top-k processing; therefore, our model can leverage any of the existing top-k work as building blocks (e.g. [20, 28, 26]).

Another important aspect of pub/sub top-k matching is to explore and identify a plausible top-k semantics. Unlike in the database context, formalizing top-k semantics in pub/sub is more involved and not limited to a single interpretation [23, 9, 10]. The most widely used pub/sub top-k semantics is defined with respect to subscribers, i.e., consume-centric semantics, in which the subscription language is extended (with a scoring function) in order to rank each incoming publication (over a time- or count-based sliding window); thus, delivering only the top-k matched publications to each subscriber [23, 9, 10].

Alternatively, the top-k semantics can be defined with respect to a publisher, i.e., produce-centric semantics, which extends the publication language for ranking subscribers and delivering publications only to the top-k matched subscribers [20, 28]. Produce-centric semantics is suitable for targeted advertisement (e.g., targeting a specific demographic group) and diversified advertisement (e.g., reaching out to most eligible members within each demographic group). Finally, hybrid semantics can be foreseen such that both subscribers and publishers have control on how data is received and disseminated, respectively. In our current work, we focus primarily on the well-adopted consume-centric top-k semantics, although some of our techniques could be leveraged for both semantics.

7 Conclusions

We have developed a distributed ranked data dissemination algorithm. The solution performs aggressive top-k filtering closer to the sources within an overlay network. The resulting top-k streams are propagated and recombined by downstream nodes. Early filtering discards low ranking data, thus saving traffic over a centralized solution at the end-user. The solution uses a chunking algorithm which switches between full forwarding of a sequence of events and selective forwarding of top-k events. This chunking algorithm ensures correctness by providing enough data to reconstruct a possible interleaving of the original event stream.

We have implemented our solution within the context of pub/sub overlay networks. Brokers closer to publishers are responsible for the top-k computation while brokers closer to subscribers recombine the chunks. Our evaluation shows that we obtain significant throughput reduction in the system compared to the centralized solution, while maintaining comparable end-to-end latency. A sensitivity analysis reveals the importance of rechunking and clustering to maximize the performance of our solution. We also show how the social network use case exhibits favorable properties for our solution in terms of scalability and efficiency.

References

1. Tumblr architecture - 15 billion page views a month and harder to scale than twitter. <http://highscalability.com/blog/2012/2/13/tumblr-architecture-15-billion-page-views-a-month-and-harder.html> (2012)
2. Aguilera, M.K., Strom, R.E., Sturman, D.C., Astley, M., Chandra, T.D.: Matching events in a content-based subscription system. In: PODC'99
3. Babcock, B., Olston, C.: Distributed top-k monitoring. In: SIGMOD'03
4. Bizer, C., Heath, T., Berners-Lee, T.: Linked data - the story so far. *Int. J. Semantic Web Inf. Syst.*'09
5. Bizer, C., Heath, T., Idehen, K., Berners-Lee, T.: Linked data on the web. In: WWW'08
6. Campailla, A., Chaki, S., Clarke, E., Jha, S., Veith, H.: Efficient filtering in publish-subscribe systems using binary decision diagrams. In: ICSE'01
7. Cao, P., Wang, Z.: Efficient top-k query calculation in distributed networks. In: PODC'04
8. Carzaniga, A., Rosenblum, D.S., Wolf, A.L.: Design and evaluation of a wide-area event notification service. *TOCS*'01
9. Drosou, M., Pitoura, E., Stefanidis, K.: Preferential publish/subscribe. In: PersDB'08
10. Drosou, M., Stefanidis, K., Pitoura, E.: Preference-aware publish/subscribe delivery with diversity. In: DEBS'09
11. Eyers, D., Freudenreich, T., Margara, A., Frischbier, S., Pietzuch, P., Eugster, P.: Living in the present: on-the-fly information processing in scalable web architectures. In: CloudCP'12
12. Fabret, F., Jacobsen, H.A., Llirbat, F., Pereira, J., Ross, K.A., Shasha, D.: Filtering algorithms and implementation for fast pub/sub systems. *SIGMOD*'01

13. Fagin, R., Lotem, A., Naor, M.: Optimal aggregation algorithms for middleware. In: PODS'01
14. Fan, M., Chen, R.: Real-time analytics streaming system at Zynga. XLDB'12
15. Fidler, E., Jacobsen, H.A., Li, G., Mankovski, S.: The PADRES distributed publish/subscribe system. In: ICFI'05
16. Fontoura, M., Sadanandan, S., Shanmugasundaram, J., Vassilvitski, S., Vee, E., Venkatesan, S., Zien, J.: Efficiently evaluating complex Boolean expressions. In: SIGMOD'10
17. Huberman, B.A., Romero, D.M., Wu, F.: Social networks that matter: Twitter under the microscope. First Monday'09
18. Ilyas, I.F., Beskales, G., Soliman, M.A.: A survey of top-k query processing techniques in relational database systems. Comput. Surv.'08
19. Kwak, H., Lee, C., Park, H., Moon, S.: What is twitter, a social network or a news media? In: WWW'10
20. Machanavajjhala, A., Vee, E., Garofalakis, M., Shanmugasundaram, J.: Scalable ranked publish/subscribe. VLDB'08
21. Michel, S., Triantafillou, P., Weikum, G.: Klee: a framework for distributed top-k query algorithms. In: VLDB'05
22. Neumeier, L., Robbins, B., Nair, A., Kesari, A.: S4: Distributed stream computing platform. In: ICDM Workshops'10
23. Pripuzić, K., Žarko, I.P., Aberer, K.: Top-k/w publish/subscribe: finding k most relevant publications in sliding time window w. In: DEBS'08
24. Pujol, J.M., Erramilli, V., Siganos, G., Yang, X., Laoutaris, N., Chhabra, P., Rodriguez, P.: The little engine(s) that could: scaling online social networks. In: SIGCOMM'10
25. Sadoghi, M., Jacobsen, H.A.: BE-Tree: An index structure to efficiently match Boolean expressions over high-dimensional discrete space. In: SIGMOD'11
26. Sadoghi, M., Jacobsen, H.A.: Relevance matters: Capitalizing on less (top-k matching in publish/subscribe). In: ICDE'12
27. Setty, V., Voulgaris, S., Vitenberg, R., van Steen, M.: Poldercast: Fast, robust, and scalable architecture for p2p pub/sub. In: Middleware'12
28. Whang, S., Brower, C., Shanmugasundaram, J., Vassilvitskii, S., Vee, E., Yerneni, R., Garcia-Molina, H.: Indexing boolean expressions. In: VLDB'09
29. Wilson, C., Boe, B., Sala, A., Puttaswamy, K.P., Zhao, B.Y.: User interactions in social networks and their implications. In: EuroSys'09
30. Wong, B., Guha, S.: Quasar: a probabilistic publish-subscribe system for social networks. In: IPTPS'08
31. Yan, T., Garcia-molina, H.: Index structures for selective dissemination of information under the boolean model. TODS'94

Algorithm 3: Incremental chunking

```

1 on receive event  $e$  matching subscription  $s$  do
2    $s.evts.add(e)$ ;
3   if  $|s.sentevts| = 0$  then
4     // Start a new chunk.
5      $s.chunkid \leftarrow generateUniqueId()$ ;
6     send( $(startchunk, s.chunkid, s), s.lasthop$ );
7   if  $|s.sentevts| < W$  then
8     //  $e$  is a left guard.
9      $s.sentevts.add(e)$ ;
10    send( $(e^{guard}, s.chunkid), s.lasthop$ );
11  else if  $getNumDisjointWindows(s.evts) < 2$  then
12    // Wait for more events.
13  else
14    // Forward top-k events.
15     $w \leftarrow getFirstWindow(s.evts)$ ;
16     $topk \leftarrow getTopK(w, k)$ ;
17    foreach event  $e' : e' \in topk \wedge e' \notin s.sentevts$  do
18       $s.sentevts.add(e')$ ;
19      send( $(e'^{topk}, s.chunkid), s.lasthop$ );
20    // Expire events.
21     $w' \leftarrow getLastWindow(s.evts)$ ;
22    foreach event  $e' : e' \in w \wedge e' \notin w'$  do
23       $s.evts.remove(e')$ ;
24       $s.sentevts.remove(e')$ ;
25 on finish chunk for subscription  $s$  do
26    $w \leftarrow getLastWindow(s.evts)$ ;
27   foreach event  $e : e \in w$  do
28     send( $(e^{guard}, s.chunkid), s.lasthop$ );
29   send( $(endchunk, s.chunkid), s.lasthop$ );
30    $s.chunks \leftarrow s.chunks \setminus \mathbb{C}$ ;
31    $s.evts \leftarrow \emptyset$ ;
32    $s.sentevts \leftarrow \emptyset$ ;

```

Algorithm 4: Deduplication

```

1 on receive  $\langle \text{startchunk}, c, s \rangle$  do
2    $s.\text{chunkStartTime} \leftarrow \text{NOW}$  ;
3    $s.\text{chunkid} \leftarrow c.\text{id}$  ;
4    $s.\text{guardids} \leftarrow \emptyset$  ;
5    $s.\text{numTopK} \leftarrow 0$  ;
6 on receive event  $e^{\text{guard}}$  for chunk for subscription  $s$  do
7    $\text{sendOnce}(e^{\text{match}}, s.\text{lasthop})$  ;
8    $s.\text{guardids.add}(e.\text{id})$  ;
9 on receive event  $e^{\text{topk}}$  for chunk for subscription  $s$  do
10   $\text{sendOnce}(e^{\text{match}}, s.\text{lasthop})$  ;
11  if  $s.\text{numTopK} = 0$  then
12    // Left guard is done.
13     $\text{send}(\langle \text{Lguard}, s.\text{guardids}, s.\text{chunkid}, s \rangle, s.\text{lasthop})$  ;
14     $s.\text{guardids} \leftarrow \emptyset$  ;
15   $s.\text{numTopK} \leftarrow s.\text{numTopK} + 1$  ;
16 on call  $\text{sendOnce}(e, \text{nexthop})$  do
17    $\text{send}(\langle \text{Rguard}, s.\text{guardids}, s.\text{numTopK}, s.\text{chunkid}, s \rangle, s.\text{lasthop})$  ;
18   if  $e \notin \text{nexthop.evts}$  then
19      $\text{send}(e^{\text{match}}, \text{nexthop})$  ;
20   // Expire old events.
21    $\text{oldestChunkStartTime} \leftarrow \min_s(s.\text{chunkStartTime})$  ;
22   foreach neighbour  $n$  do
23      $n.\text{evts.removeOlderThan}(\text{oldestChunkStartTime})$  ;

```

Algorithm 5: Rehydration

```

1 on receive event  $\langle e^{match} \rangle$  do
2   foreach subscription  $s$  that matches  $e^{match}$  do
3      $s.evts.add(e)$  ;
4 on receive  $\langle Lguard, \{eId\}, c, s \rangle$  do
5   send( $\langle startchunk, c, s \rangle, s.lasthop$ ) ;
6   foreach  $id \in \{eId\}$  do
7      $e \leftarrow s.evts.get(id)$  ;
8     send( $\langle e^{guard}, c \rangle, s.lasthop$ ) ;
9      $s.evts.removeOlderThan(e)$  ;
10 on receive  $\langle Rguard, \{eId\}, n, c \rangle$  do
11    $guardevts \leftarrow \emptyset$  ;
12   foreach  $id \in \{eId\}$  do
13      $guardevts \leftarrow guardevts \cup s.evts.get(id)$  ;
14     // Compute and send top-k events.
15      $potentialtopkevts \leftarrow s.evts.getOlderThan(guardevts)$  ;
16      $topkevts \leftarrow potentialtopkevts.getTopRanked(n)$  ;
17     foreach  $e \in topkevts$  do
18       send( $\langle e^{topk}, c \rangle, s.lasthop$ ) ;
19     // Send guard events.
20     foreach  $e \in guardevts$  do
21       send( $\langle e^{guard}, c \rangle, s.lasthop$ ) ;
22        $s.evts.removeOlderThan(e)$  ;

```
