# Efficient Notification Ordering
# for Geo-Distributed Pub/Sub Systems

Roberto Baldoni, Silvia Bonomi, Marco Platania, and Leonardo Querzoni

**Abstract**—A distributed event notification service (ENS) is at the core of modern messaging infrastructures providing applications with scalable and robust publish/subscribe communication primitives. Such ENSs can route events toward subscribers using multiple paths with different lengths and latencies. As a consequence, subscribers can receive events out of order. In this paper, we propose a novel solution for ordered notifications on top of an existing distributed topic-based ENS. Our solutions guarantees that each pair of events published in the system will be notified in the same order to all their target subscribers independently from the topics they are published in. It endows a distributed timestamping mechanism based on a multistage sequencer that produces timestamps whose size is dynamically adjusted to accommodate changing subscriptions in the system.

An extensive experimental evaluation based on a prototype implementation shows that the timestamping mechanism is able to scale from several points of view (i.e., number of publisher and subscribers, event rate). Furthermore, it shows how the deployment flexibility of our solution makes it perform better in terms of timestamp size and timestamp generation latency when the system load exhibits geographic topic popularity, that is, matching subscriptions and publications are geographically clustered. This makes our solution particularly well suited to be deployed in geo-distributed infrastructures.

**Index Terms**—Total order, Publish/Subscribe, Geo-Distributed Systems, Logical timestamps, Event based communications, Geographic Topic Popularity.

✦

## 1 INTRODUCTION

Modern large-scale services are usually built on top of asynchronous communication primitives able to mask the unreliability of low-level networks and the dynamism of the application participants by decoupling the interacting parties in space and time. The publish/subscribe paradigm provides communication services where message addressing is implicitly handled by an *Event Notification Service* (ENS), a middleware infrastructure that matches the content of events produced by publishers against interests expressed by subscribers in the form of subscriptions.

Many research efforts in publish/subscribe systems focused on reliability and performance aspects with few contributions in the area of event ordering [1]–[5]. Defining a coherent specification for notification ordering is a fundamental step for a wide range of applications like stock tickers, messaging, command-and-control, or those based on composite event detection [6], where specific event patterns must be concurrently detected by distributed and possibly independent application components. As an example, in distributed online games [7] users in close proximity in a virtual world are supposed

to see events happen in a single consistent order. In Electronic stock tickers [8], investors are required to see coherently ordered stock price evolutions to take their investment decisions.

In this paper we consider the following ordering problem: how to guarantee that two subscribers sharing subscriptions with common interests are notified about events matching those subscriptions in the same order. While the above ordering problem stems from the simple rationale that two participants should always see the notification of two events in the same order, its enforcement in distributed ENSs is far from being trivial. Violations to the ordering property can easily arise due to the fact that two events, possibly published by different publishers, can follow distinct paths through the ENS before reaching the points where they will be notified to the final recipients. Furthermore, non-determinism, in the form of unpredictable network latencies and message losses, can easily exacerbate the problem, especially in geo-distributed application scenarios.

Current solutions either require complex offline set-ups that must be continuously updated when subscribers change their interests [2], or give up some ordering aspects only guaranteeing per-source ordering [9], or are based on synchronization among processes in order to deterministically order conflicting events (i.e., [1], [4], [5]). Such synchronization approaches, either hardware-based [1], or based on total order [4] among all processes receiving an event (using centralized sequencers or distributed consensus primitives like Paxos [10]) or, finally, characterized by the widespread usage of explicit acknowledgments [5], end up into scalability problems with respect to performance when both the

- R. Baldoni, S. Bonomi, and L. Querzoni are with the Department of Computer, Control, and Management Engineering, Sapienza University of Rome, Rome, Italy.
  M. Platania is with the Department of Computer Science, Johns Hopkins University, Baltimore MD, USA
  E-mail: {baldoni|bonomi|querzoni}@dis.uniroma1.it; platania@cs.jhu.edu

event rate and the number of entities involved in the synchronization increase [11] or as soon as WAN links are involved [12].

This paper focuses on studying event ordering for geo-distributed publish/subscribe communication middleware looking for a scalable solution in terms of number of published events and subscriptions, while accommodating subscription changes at run-time. More specifically, the paper introduces a novel timestamping solution designed to be used with topic-based publish/subscribe systems. The purpose of our solution is to generate logical timestamps that can be used, on receiver side, to enforce the following *total notification order* (TNO) property without any explicit synchronization: if two independent subscribers are notified about the same two events, then these two events will be notified to them in the same order[1].

The core of our solution is a new logical timestamping mechanism based on subscribers' interests. The structure and size of each timestamp is automatically calculated at run-time on the basis of current subscription overlapping, in order to keep it to a minimum (as it can vary between a single integer and a vector of integers whose size is the number of topics). Timestamps are built through a multistage sequencer based on a distributed architecture. A key feature of this architecture lies in its deployment flexibility: multiple stages of the sequencer can be co-located on machines deployed in different sites of a geo-distributed infrastructure (e.g., a service provider with multiple data centers). This flexibility allows system designers to exploit specific locality characteristics enjoyed by many large-scale geo-distributed applications, such as YouTube[2], to drastically reduce both timestamp size and generation latency. Geographic topic locality can indeed make most of the overlapping among subscriptions local to a specific geographic area increasing the probability that a timestamp will be entirely generated within a single site of the distributed timestamping architecture, thus avoiding costly inter-site (i.e., WAN) communication.

From an architectural point of view, the timestamping mechanism, encapsulated within a software component that can be deployed on top of existing reliable topic-based publish/subscribe middleware, transparently delivers events notified by the ENS to the application layer guaranteeing the total notification order. When deployed on top of a non-reliable ENS, our mechanism is able to deterministically tag every event whose notification violates the TNO property; this gives application developers the possibility to treat out-of-order events in an appropriate manner.

The performance of our solution have been analyzed through an extensive experimental evaluation. The re-

sults show how it creates timestamps whose size scales with respect to the number of subscribed topics. We developed a prototype implementation of our solution, in order to study its behavior in a realistic geo-distributed setting, mimicking a common architecture employed by cloud-providers and by large companies with several data centers. The experimental evaluation of the prototype takes into account a *geographic topic popularity*, in which subscriptions and publications are geographically strongly clustered [13], and a *spray-and-diffuse* pattern, in which interest clustering is mildly present [13], [14]. Results show that an increasing matching between topic popularity and locality of interest produces a timestamp generation latency of less than 300 ms in the presence of intense publication rates (up to 10000 events/sec).

The rest of this paper is organized as follows: Section 2 introduces the system model and states the problem explaining why its solution includes several difficult aspects; Section 3 describes our algorithm; Section 4 presents some important engineering aspects; Section 5 reports the results of the evaluation of our solution; Section 6 explains how the problem of ordering events has been tackled in the literature and, finally, Section 7 concludes the paper.

## 2 SYSTEM MODEL AND PROBLEM STATEMENT

We consider a system composed by a number of interacting clients that can act as publishers (data producers) or subscribers (data consumers). Clients exchange data in the form of events using a topic-based selection model, thus we assume that they share a common knowledge on a fixed set of available topics. Each piece of data produced by a publisher is published on one of the available topics and takes the form of an event. Each subscriber issues a subscription $S$ containing the set of topics it is interested in. An event $e$ published on a topic $T$ matches a subscription $S$ if and only if $T \in S$; when this happens, the corresponding subscriber must be notified about $e$. Clients do not interact directly: their interactions are mediated by an *Event Notification Service* (ENS) that exposes the fundamental interface of a publish/subscribe system, i.e., the *publish*, *subscribe/unsubscribe* and *notify* primitives. Without loss of generality, here we assume that the ENS is implemented as a distributed middleware.

In addition, in order to simplify the description of our solution, we will initially assume that our system works on top of a reliable communication substrate, that all communication links deliver messages in FIFO order, and that all processes are correct. Section 4 details how some of these assumptions can be removed or relaxed.

The ordering property we want to enforce is defined as follows:

*Property 1:* TOTAL NOTIFICATION ORDER (TNO). Let $e_i$ and $e_j$ be two distinct events notified to a subscriber $s$. If $e_i$ is notified to $s$ before $e_j$, no subscriber will be notified about $e_i$ after being notified about $e_j$.

---

1. The TNO property, also known as *Pairwise Total Order*, is considered in the literature as one of the strongest form of ordering achievable in distributed publish/subscribe middleware [5].

2. Brodersen et al. [13] showed that at least 40% of YouTube videos have 80% of their total views coming from a single country, indicating strong user interest locality.

Note that this definition matches the definition of *Weak Total Order* given in [15] in the context of total order specifications [16]. Differently from those specifications, we do not consider any form of deterministic agreement (uniform or not uniform) because here we are only interested in designing an ordering layer to be transparently plugged on top of a generic ENS which can provide different reliability and agreement properties. The TNO property also matches the *pairwise total order* property defined in [5].
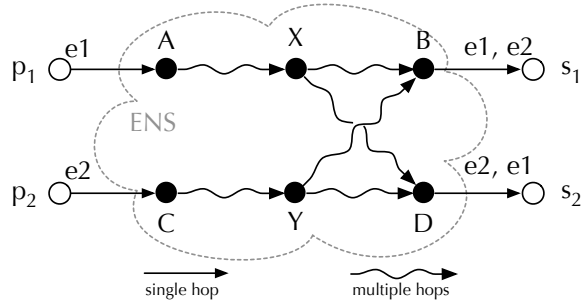


Fig. 1. An example showing how notifications can be performed out of order in a distributed event notification service.

Guaranteeing TNO in a distributed setting is a complex task. As an example, consider the toy system depicted in Figure 1: the six black dots represent processes constituting the ENS, the white dots on the left ($p_1$ and $p_2$) are two publishers and those on the right ($s_1$ and $s_2$) are two subscribers. A common solution for ordering events published on a specific topic is based on the usage of sequencers: a single node in the ENS is elected as a "sequencer" for all the events published in that topic. In our example $X$ acts as the sequencer node for topic $T1$, receiving all the events published in $T1$ (i.e., event $e_1$ published by $p_1$), adding a sequence number to them, and then routing the events toward the intended destinations (i.e., $s_1$ and $s_2$ notified by nodes $B$ and $D$). Similarly, $Y$ is the node in charge of sequencing events published in topic $T2$ (event $e_2$ in the example).

This simple approach, however, is not useful when subscriptions intersect in multiple topics. For example, assume that both $s_1$ and $s_2$ are subscribed to $T1$ and $T2$. In this case the sequence numbers attached by $X$ and $Y$ would be completely uncorrelated and useless to check for a correct notification order on the subscribers' side. Centralized solutions, i.e., using a single sequencer for all the topics, have important scalability drawbacks. Similarly, distributed consensus algorithms impose stringent latency requirements to provide synchronization in a timely manner [11], [12]. Therefore, they cannot be realistically considered in scenarios where large scale or large loads are expected.

## 3 THE EVENT ORDERING ALGORITHM

In this Section we first introduce an abstraction to illustrate how the event ordering problem in publish/subscribe systems can be theoretically addressed, and the design principles underlying a possible distributed implementation. Then, we detail the algorithm that implements the proposed solution. Due to space constraints the algorithm pseudocode and the related correctness proofs are reported in the supplemental material.
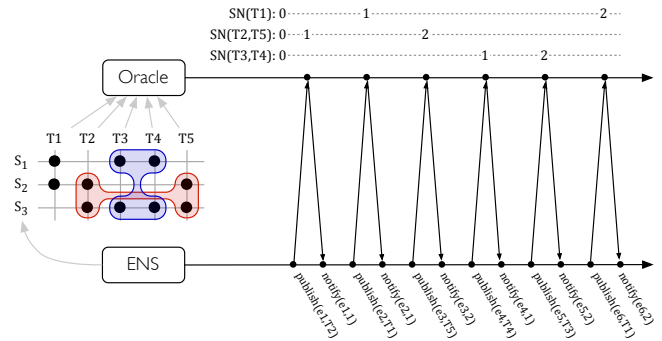


Fig. 2. The oracle assigns sequence numbers to events published in different topics on the basis of intersections among subscriptions.

### 3.1 Event Ordering Abstraction

From a theoretical standpoint we need an event ordering abstraction (*oracle* in the following) able to produce sequence numbers. Such oracle would provide, for each topic, the sequence number representing the timestamp of the next event published in that topic. Sequence numbers would need to be generated such that two events published on topics that are both subscribed by at least two subscribers have different comparable timestamps. To enforce this condition, and thus guarantee TNO, the oracle must check all subscription intersections and thus requires complete knowledge of subscriptions. Figure 2 shows an example where the oracle accesses knowledge of subscriptions from the ENS to identify intersections among $s_2$ and $s_3$ on topics $T2$ and $T5$, and among $s_1$ and $s_3$ on topics $T3$ and $T4$. As a consequence the oracle will maintain a sequence number for $(T2, T5)$ and one for $(T3, T4)$. A third sequence number will be maintained to timestamp events published in $T1$ only as events for this topic must not be ordered with respect to events published in any other topic. Thanks to these sequence numbers, events $e1$ and $e3$, published in $T2$ and $T5$ respectively, will be notified by all subscribers ($s_2$ and $s_3$) in the order defined by their timestamps (1 for $e1$ and 2 for $e3$), independently from possible reordering happening in the ENS during the event diffusion phase. Event $e_2$ published in topic $T1$ can be notified from $s_2$ independently from the former events (i.e., without a precise order), as the only other subscriber in $T1$, i.e., $s_1$, will not be notified about $e1$, nor $e3$.

In the following sections we introduce a distributed implementation of this oracle and show how to use it to guarantee TNO on the subscriber side. However, before delving in these design details, it is useful to clearly outline the design principles underlying our solution.

## 3.2 Design principles

A system implementing the event ordering abstraction should match the following design principles:

**Full decoupling** - It must retain the typical full decoupling characteristics of publish/subscribe, i.e., no direct interactions should happen between publishers and subscribers.

**Support for large subscription loads** - It must gracefully scale in scenarios with large number of subscriptions and multiple possible interest intersections. This must be achieved by both adopting short timestamps to reduce the overhead on the ENS and a clever internal design. The system should be designed in a distributed fashion with the aim of sharing the load imposed by timestamp generation on multiple machines (i.e., no machine should maintain full knowledge of the system state) and thus avoid possible bottlenecks.

**Asynchronous one-way message flows** - It must avoid any kind of explicit synchronization among its internal processes by adopting a one-way message flow strategy (i.e., no ACKs are required) for the most common operations, like event timestamping or subscription management; the lack of explicit synchronization is a crucial principle needed to support intense event publication rates [11].

**Full independence from the ENS** - It must be designed to be transparently pluggable on top of existing topic-based publish/subscribe middleware platforms. This principle facilitates the deployments within existing infrastructures.

**Flexible deployment** - It must allow the deployment of its internal components in a flexible manner so to efficiently support different application scenarios, from simple centralized setups to distributed deployments with geographic topic popularity, while maximizing the available resource usage.

## 3.3 Architectural aspects

Our solution assumes that all participants to the system (publishers and subscribers) are equipped with an *Ordering module* that implements the algorithm described in the next Section (see Figure 3). This module mediates the interactions between application level software components, that act as information producers (publisher applications) or consumers (subscriber applications), and a standard ENS.

We assume that the ENS implements a standard topic-based publish/subscribe interface (here represented by the *ENSpublish*, *ENSsubscribe/ENSunsubscribe* and *ENSnotify* operations). The same interface is offered by
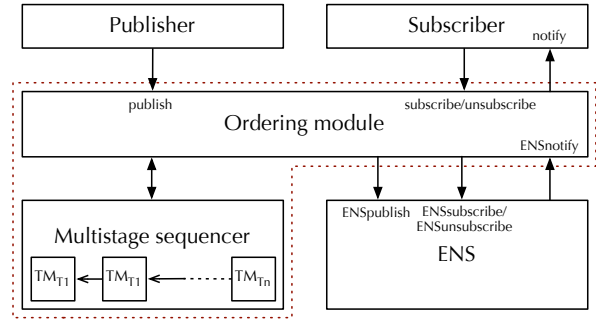


Fig. 3. Architectural view that shows how the Ordering module acts as a mediating software layer between the applications and an existing event notification service. The dashed line represents the whole timestamping mechanism that also includes a multistage sequencer used to produce timestamps.

the ordering module to the application level, therefore neither the applications, nor the ENS must be changed in order to work with our solution. In the following of this section we will consider a deployment where our solution is coupled with a reliable ENS. In particular, we assume that when the ENS returns from an invocation to *ENSsubscribe*, all events published on the subscribed topic after that point in time will be eventually *ENSnotified* to the subscriber. Section 4 will provide further details on how the solution behaviour changes when it is coupled with a non-reliable ENS. Note that, thanks to the full adherence of the ordering module with the standard publish/subscribe interface, it could also be used in conjunction with a mapping layer that efficiently allows the adoption of a content-based interface on top of a topic-based publish/subscribe system [17].

The ordering module also needs to access a point-to-point communication primitive that can be offered by the operating system or by other solutions like an overlay network. We also assume that the set of available topics is fixed and a precedence relationship $\rightarrow$ holds among topic identifiers inducing a total order on them: if $T \rightarrow T'$ we say that $T$ has a higher rank in the relationship than $T'$.

Ordering modules communicate with a multistage sequencer, constituted by a set of processes, one per topic, called *topic managers* ($TM$s). Finally, we assume there is a method to univocally map a topic $T$ to its topic manager $TM_T$. This problem can be solved in several different ways, i.e., through a static mapping provided as a configuration parameter, using a DNS, or resorting to a distributed hash table as in rendez-vous based publish/subscribe systems [18]. In the following, whenever there is no ambiguity, we will use the terms *publisher* and *subscriber* to refer the parts of our ordering module located respectively at the publisher and at the subscriber.
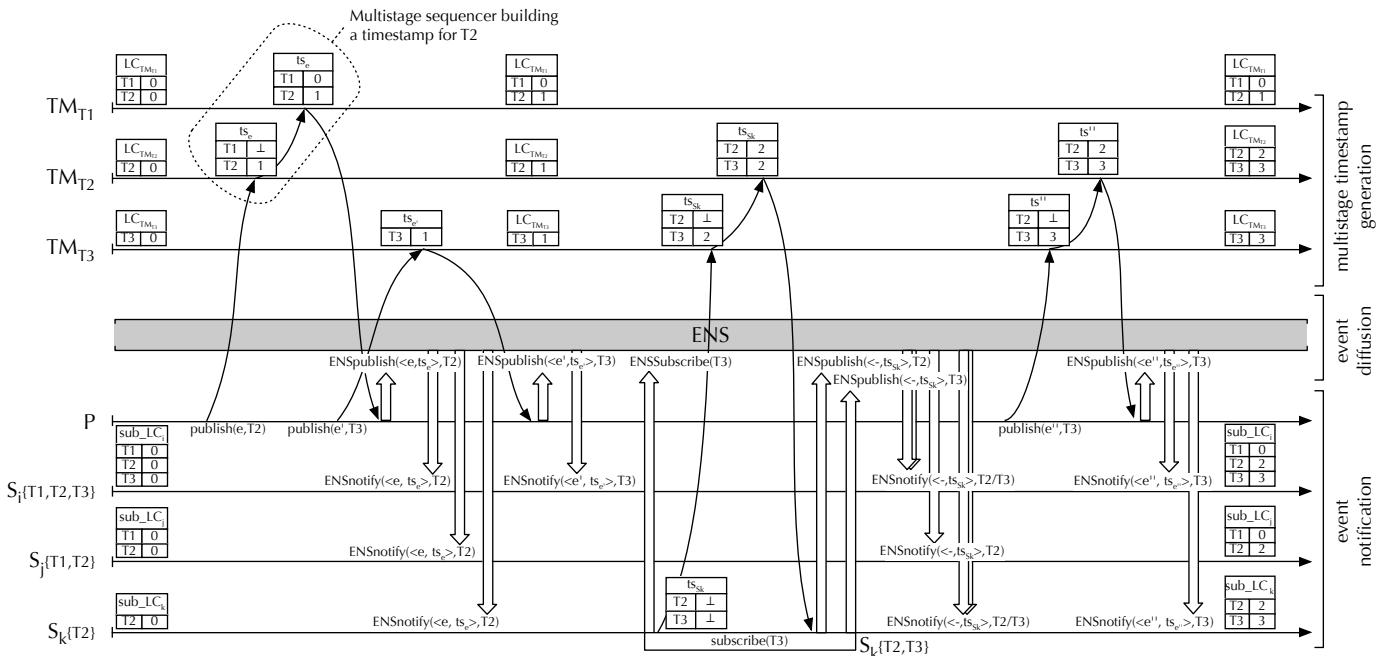
Fig. 4. Example of a system run with three subscribers, $S_i$, $S_j$ and $S_k$, and a publisher $P$.

## 3.4 Algorithm description

The basic idea behind the algorithm is to assign a logical timestamp *ts* to each event. By looking at a timestamp, a subscriber must be able to decide if the event can be immediately notified to the application level or there is another event that precedes it in the TNO order. In this latter case the current event is locally buffered while the subscriber waits for the missing event to be notified by the ENS.

The algorithm to be executed when an event is published is split in three phases (Figure 4): (i) *multistage timestamp generation*, where a timestamp is generated for the event; (ii) *event diffusion*, where the ENS delivers the event and its timestamp to all the intended subscribers; and (iii) *event notification*, where subscribers, by looking at the timestamp content, decide if the event is the next in the TNO order to be notified. The algorithm uses only local information maintained by each process. A topic manager $TM_T$ stores all subscriptions containing topic $T$, a sequence number $LC_T$ that counts the number of events published in $T$ and a (possibly empty) set of sequence numbers $LLC_T$ storing identifiers and sequence numbers of topics $T', T'', \cdots$ with lower ranks than $T$. Each subscriber stores its subscription $S$ and a set $sub\_LC$ containing the sequence number of the last event notified on $T$, for each topic $T \in S$ (i.e., it maintains a local subscription clock).

In the following, before describing the multistage timestamp generation procedure, we first provide a few formal definitions we will use later (i.e. *sequencing group*, *timestamp*, and *order relation between two timestamps*).

Informally, a *sequencing group* for a topic $T$ contains the (ordered) set of all the other topics whose events must be ordered with respect to events published in $T$ to guarantee TNO.

*Definition 1:* A *sequencing group* of a topic $T$ ($\mathcal{SG}_T$) is a set of topics including $T$ and all $T'$ such that there are at least two subscriptions including both $T$ and $T'$.

Therefore, $\mathcal{SG}_T$ is a one-way sequence of topics, whose direction is determined by the precedence relationship $\rightarrow$, and whose content may only change with subscription updates. A topic manager $TM_T$ can calculate $\mathcal{SG}_T$ by looking at the list of subscriptions containing $T$ it holds. Specifically, $TM_T$ (i) calculates the union of all topics in the subscriptions it knows and (ii) removes all topics that appear in a single subscription (except $T$). The resulting set is $\mathcal{SG}_\mathcal{T}$.

Given the sequencing group of a topic $T$ we can now define how a timestamp for events published in $T$ must be structured. Informally, the timestamp contains a set of entries, one for each topic contained in the sequencing group of $T$. Each entry represent a sequence number for a specific topic $T' \in \mathcal{SG}_\mathcal{T}$.

*Definition 2:* Let $e$ be an event published in a topic $T$, $T_1 \rightarrow T_2 \rightarrow \cdots \rightarrow T_n$ the topic ordering according to the precedence relation $\rightarrow$, and $\mathcal{SG}_T$ the sequencing group of $T$. A timestamp $ts_e$ for $e$ is a set of pairs $< T_i, sn_i >$ ordered according to $\rightarrow$, where $T_i \in \mathcal{SG}_T$ is a topic identifier and $sn_i$ is the sequence number for $T_i$.

Now that we know how a timestamp is internally structured, we can define when and how two timestamps can be compared.

*Definition 3:* Let $ts_e$ and $ts_{e'}$ be two timestamps associated with two different events $e$ and $e'$. We say that $ts_e$ and $ts_{e'}$ are *comparable* if there exists at least one topic identifier present both in $ts_e$ and $ts_{e'}$ (i.e., $\exists\, t_{id}, i, j \mid (< t_{id}, i > \in ts_e) \wedge (< t_{id}, j > \in ts_{e'}))$.

From the three definitions above, it is easy to see that given two events $e$ and $e'$ published respectively in topics $T$ and $T'$, the corresponding timestamps $ts_e$ and $ts_{e'}$ are comparable if and only if $\mathcal{SG}_T \cap \mathcal{SG}_{T'} \neq \emptyset$. The fact that their sequencing groups intersect means that a total order must be defined for events published within them.

*Definition 4:* Let $ts_e$ and $ts_{e'}$ be two timestamps associated with two different events $e$ and $e'$. We say that $ts_e$ is *smaller* than $ts_{e'}$ (i.e., $ts_e < ts_{e'}$) if

1) $ts_e$ and $ts_{e'}$ are comparable, and
2) $\forall < t_{id}, sn > \in ts_e \mid \exists < t_{id}, sn' > \in ts_{e'},\ sn \leq sn'$, and
3) $\exists < t_{id}, sn > \in ts_e \mid \exists < t_{id}, sn' > \in ts_{e'},\ sn < sn'$

As an example, in Figure 4 we show the timestamps for two published events $e$ and $e'$. Considering the timestamp $ts$ associated with $e$ and the timestamp $ts'$ associated to $e'$ we have that they are not comparable.

**Multistage timestamp generation:** When the publication of an event $e$ in a topic $T$ occurs, the Ordering module on the publisher side contacts the topic manager $TM_T$ to obtain a timestamp for $e$. $TM_T$ starts a collaborative multistage timestamp generation procedure that involves the subset of $TM$s associated with topics belonging to the *sequencing group* of $T$. Specifically, $TM_T$:

1) Creates a timestamp structure with an entry for each topic $T'$ such that $T' \in \mathcal{SG}_T$;
2) Increases its local sequence number $LC_T$;
3) Inserts this value in $T$'s entry in the timestamp;
4) Inserts in the timestamp the values related to all the topics $\overline{T} \in \mathcal{SG}_T$ such that $T \to \overline{T}$;
5) Forwards the timestamp to the $TM$ associated to the first topic in $\mathcal{SG}_T$ that precedes $T$ according to the precedence relation $\to$.

Note that, given a specific order $T_1 \to T_2 \to \cdots \to T_n$ among topics, the multistage timestamp generation flow proceeds in the opposite direction (i.e., given a topic $T_i$, $TM_{T_i}$ will fill in the timestamp and forward it to some $TM_{T_j}$ such that $T_j \to T_i$). In addition, based on the definition of *sequencing group*, the multistage timestamp generation is a one-way flow of messages, so to avoid loops among $TM$s that may prevent a correct timestamp construction [2] or create instability in large scale high-throughput systems [19]. The receiving $TM$ also inserts in the timestamp the sequence number of the topic it manages, without increasing it. Then it forwards the timestamp to the next $TM$ in $\mathcal{SG}_T$. When the last $TM$ completes the timestamp, it is returned to the publisher that will publish the event in the ENS together with the

timestamp. During this process each $TM_{T'}$ receiving a partially filled timestamp, uses its content to update the local clocks $LC_{\overline{T}}$ for all topics $\overline{T} \in \mathcal{SG}_{\mathcal{T}}$ such that $T' \to \overline{T}$. Note that, in order to increase the scalability of the multistage sequencer, only event $id$s travel within requests, while event payloads are buffered on the publisher side while it waits for the multistage timestamp generation procedure to complete.

Figure 4 shows a run of the algorithm in a system with three subscriptions $S_i : \{T_1, T_2, T_3\}$, $S_j : \{T_1, T_2\}$, and $S_k : \{T_2\}$ and the precedence relation as $T_1 \to T_2 \to T_3$. The intersection of $S_i$ and $S_j$, and the precedence relation determine the following *sequencing groups*: $\mathcal{SG}_{T1} = \mathcal{SG}_{T2} = \{T1, T2\}$, $\mathcal{SG}_{T3} = \{T3\}$. The publisher $P$ publishes an event $e$ on topic $T2$ and asks $TM_{T2}$ to create the timestamp. $TM_{T2}$ creates the structure of the timestamp with entries for topics $T2$ and $T1$, puts its sequence number in the timestamp and forwards it to $TM_{T1}$ that, in turn, will complete the timestamp and return it to $P$. Finally $P$ publishes both $e$ and its timestamp on the ENS. Local clocks $LC_{TM_{T1}}$ and $LC_{TM_{T1}}$ are updated accordingly.

In the event notification phase, once an event $e$ and its timestamp are notified by the ENS, the subscriber checks if the timestamp attached to the event is coherent with the event order maintained through the local subscription clock $sub\_LC$. This check is performed by looking at the sequence numbers included in the timestamp: if the values for all the topics are equal to the corresponding ones stored locally in $sub\_LC_i$, except for the topic where the event has been published, that must have a value greater than the local one by one unit, then no event with a smaller timestamp exists that must be still received by the subscriber, and the received event can thus be notified. Conversely, if a gap in the timestamp sequence is detected $e$ is locally buffered. Buffered events are delivered as soon as all the events preceding them in the TNO order have been delivered. The assumption that the underlying ENS is reliable guarantees that buffered events will be eventually notified to the application level. Figure 5 shows an example where network lags induce a mis-ordering between two events $e1$ and $e2$ at $S_k$. A gap in the correct notification sequence is locally detected at the subscriber that buffers $e2$ waiting for the notification of $e1$ from the ENS. As soon as $e1$ is notified by the ENS, both events, in the correct order, are notified at the application level.

**Timestamp properties:** At a first glance, timestamps provided by the multistage sequencer resemble vector clocks, but they have very different structures. Vector clocks have a well defined and fixed structure that depends on the number of processes in the computation. On the contrary, our timestamp structure is related to topics rather than processes and its size depends on the current set of subscriptions. In particular, the size of a timestamp associated with events published in a topic $T$ is equal to the dimension of $\mathcal{SG}_T$. In the best case $\forall T, |\mathcal{SG}_T| = 1$. This is the case, for example, in which
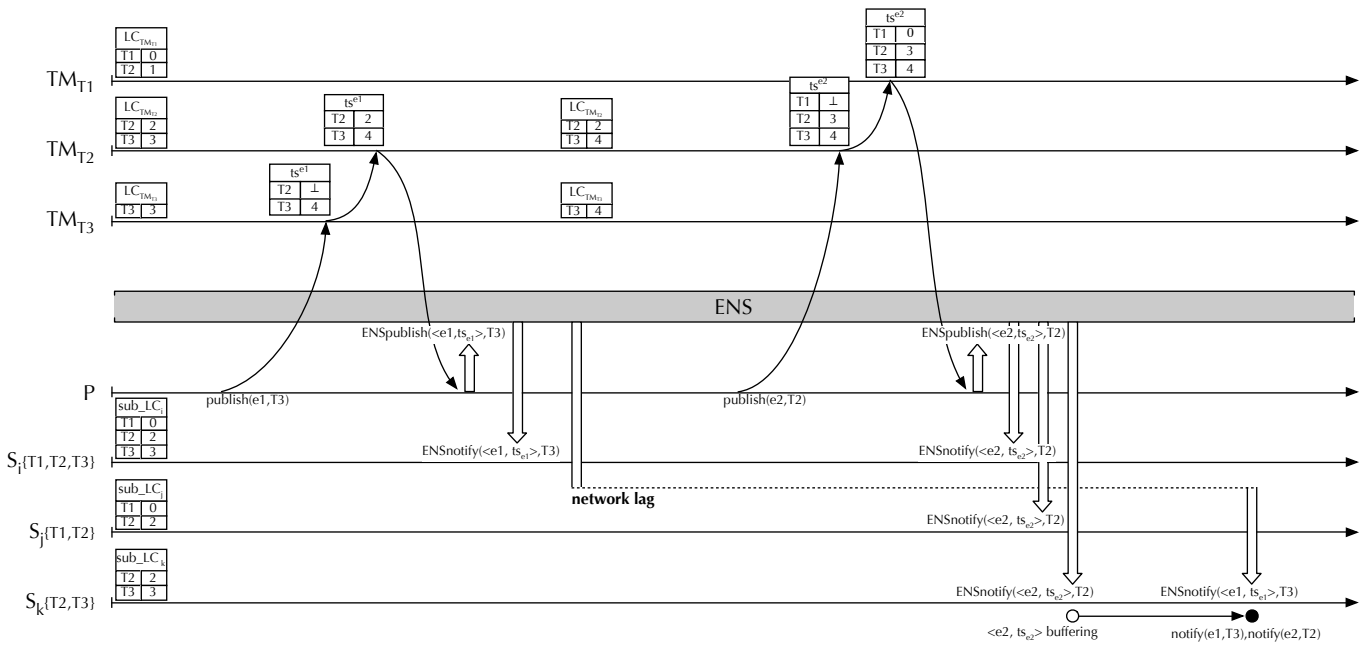
Fig. 5. Example of reordering for out of order notifications.

subscribers only subscribe a single topic each, thus avoiding intersections among subscriptions. In the worst case, instead, the size of $\mathcal{SG}_T$ equals the number of topics in the system, and this happens when all subscribers subscribe all topics. However, typically subscribers are interested just in a subset of topics. Hence, the timestamps size is expected to be much smaller than the topic cardinality. Typical applications are characterized by a total number of topics that is much lower than the number of participating processes. To this respect, our timestamps provide a more scalable solution to order events.

In addition, differently from vector clocks, where the ordering relationship holding among two events can be detected just comparing (entry by entry) the two vector clocks associated with the two events, with our timestamps the ordering relation can be detected looking firstly at the timestamp structure (i.e., the events have to be comparable according to Definition 3), and secondly, by examining the values contained in common entries of the timestamps. Let us finally remark that in our timestamping technique, the timestamp associated with an event does not bring any information about the producer of the event, thus preserving the space decoupling principle of the publish/subscribe paradigm.

**Subscription management:** Every new subscription (unsubscription) to (from) topic $T$ induces a modification of $\mathcal{SG}_T$ and, then, of the set of $TM$s used in the multistage timestamp generation phase. Therefore, when a topic $T$ is added to a subscription $S$, each topic manager $TM_{T'}$ associated with a topic $T'$ in $S$ must be advertised in order to let it recalculate $\mathcal{SG}_{T'}$ and thus avoid possible TNO violations. Furthermore, as soon as a subscriber subscribes a new topic, he will

start being notified about events published in that topic. Due to the lack of synchronization between the time an event is timestamped and published and the time it is notified by the ENS to the intended recipients, the subscriber could both receive events timestamped before or after its subscription. Such events could possibly bring different timestamps (because sequencing have changed with its subscription). Those produced after its subscription are always comparable with other events published in different topics the subscriber is notified about. However, timestamps of events generated before its subscription could be non-comparable, making it impossible for the subscriber to infer the correct notification order. It is thus necessary for the subscriber to gather enough information before completing the subscription to deterministically discern events produced before it from those produced after.

To this aim, the subscriber first subscribes to $T$ and, as soon as the invocation returns from the ENS, it starts buffering incoming events notified for $T$. Buffering at this stage is fundamental because only a subset of the buffered events, those published after the subscription, shall be notified in the right order at the application level. The subscriber then creates an empty subscription timestamp, containing one entry for each topic in the subscription. The timestamp and the new subscription including $T$ are forwarded to the $TM$ associated with the lowest ranked topic among the subscribed ones. Similarly to event timestamps, each entry of the subscription timestamp is filled in by the corresponding topic manager, which, however, also increases its local sequence number. In addition, each $TM$ that receives a request, updates the list of subscriptions it holds. When the timestamp is complete, it is sent back to the

subscriber, which uses it to update values contained in its local clock. Thanks to this subscription timestamp the subscriber is now able to clearly distinguish events produced before its subscription from those produced later, and can thus analyze the content of the buffer to discard old events and notify received events for $T$ with timestamps greater than the one associated to its subscription.

Due to the fact that the subscription procedure increases the value of timestamps for all topics included in $S$ and possibly changes multiple sequencing groups beside $\mathcal{SG}_T$, subscribers of these topics must be informed of this increase, otherwise they will start to endlessly buffer new event notifications as they suppose that the gap in the sequence number sequence is due to a missing notification. In order to avoid this undesirable side effect, the subscriber, after subscribing to the new topic also publishes a "dummy" event in all topics included in $S$ attaching to it the subscription timestamp it just received. This dummy event has no application meaning, and thus it will never be notified to the application level, but its accompanying timestamp is used by all notified subscribers to correctly update their local clocks.

Figure 4 shows an example where subscriber $k$ starts with subscription $S_k = \{T2\}$ and updates it at runtime by subscribing $T3$. Event $e'$ published in topic $T3$ before $k$'s subscription is timestamped with a single serial number added by $TM_{T3}$ as $\mathcal{SG}_{T3} = \{T3\}$. When $k$ starts its subscription, it first subscribes to $T3$ with the ENS ($ENSsubscribe(T3)$) and then sends a request to $TM_{T3}$ containing a subscriptions timestamp with an empty entry for every subscribed topic, including also an entry for $T3$ (i.e., $\{T2, T3\}$ in the example). By receiving this request, $TM_{T3}$ updates its local list of subscriptions, increases $LC_{T3}$ and fills the corresponding entry in the timestamp with its value before forwarding it to $TM_{T2}$, which, in turn, will repeat the same operation and finally return the timestamp to $k$. By receiving the completed timestamp, $k$ will use its content to update its local clock (that from now on will also contain an entry for $T3$) and will then flush the buffer selecting received events that must be notified at the application level (not shown in the example).

Finally, $k$ publishes two dummy events in $T2$ and $T3$ that are notified to the corresponding subscribers ($i, j, k$ and $i, k$ respectively). The subscribers update their local clocks with the information contained in the event timestamp. The figure also reports the example of a third event $e''$ published in topic $T3$ showing how the topic timestamp has been dynamically adapted by the system to take into account new ordering opportunities with events produced in $T2$.

A similar approach is used to unsubscribe a topic $T$. However, during an unsubscription operation a subscriber can simply start to ignore further events received for $T$ since the moment $unsubscribe(T)$ is invoked, thus avoiding any possible TNO violation. In order to maintain the system efficient, the subscriber must also inform all the relevant $TM$s about its subscription change. Without this step, in fact, sequencing groups would remain unchanged despite the possible removal of subscription intersections, and this would cause timestamps to be possibly larger than what is required by our solution to guarantee TNO. To this aim the subscriber sends to the $TM$ associated to each of its subscribed topics a message containing its updated subscription, and to $TM_T$ an empty subscription (because it is no longer subscribed to $T$, $TM_T$ does no longer need to maintain its subscription). Upon receiving such messages $TM$s simply updates their local subscription lists, without updating any local clock.

## 4 RELIABILITY AND ORDERING ASPECTS

**Working with unreliable ENSs** - The algorithm introduced in Section 3.4 assumes that the ordering module is deployed in conjunction with a reliable ENS. While this is a desirable setup, sometimes adopting a non-reliable publish/subscribe middleware, i.e. a middleware that does not guarantee the notification of all events to all the intended recipients, may be preferable. As an example, non-reliable publish/subscribe middleware often sport better performance with lower end-to-end notification latencies, and are able to scale to larger sizes. This setup creates a simple but fundamental problem to our solution: whenever a gap in the sequence of events is spotted by a subscriber during the notification phase, the subscriber cannot decide if the missing event has been lost (because the ENS is unreliable) or its notification is just late. The strategy our solution adopts with respect to this issue is optimistic: every event notified in order by the ENS is notified to the application level, while late notifications are tagged as *out-of-order* and immediately notified to the application level. This leaves to the application developer the choice to discard these events or treat them in an appropriate way.

The main source of out-of-order notifications lies in the fact that two events, possibly published by different publishers, can follow distinct paths through the ENS, before reaching the point where they will be notified to the final recipients. To reduce the number of out-of-order notifications we can use again a buffering strategy on the subscriber side, but with some important differences in its management with respect to the solution shown in Section 3.4.

Every time the ENSnotify() primitive returns a new event $e$, the algorithm checks through the attached timestamp whether some other event may exist with a smaller timestamp. If there is a possibility that an event with a smaller timestamp exists but has not been delivered to the subscriber so far, then the event $e$ is enqueued to a buffer able to host a maximum of $b$ events and a timer for $e$ is started ($TTL_e$). The event $e$ is delivered through the notify() primitive when one of the following conditions holds: (i) all the events with smaller timestamps have been notified, (ii) $TTL_e$ expires or (iii) the buffer is full,

a new event must be buffered and $e$ is at the head of the queue. By carefully tuning the buffer size and the timer length, system integrators can compensate for possible network delay fluctuations by paying some further end-to-end notification latency induced by event buffering.

**Reliability** - Making the algorithm presented so far work reliably in an environment where messages can be lost requires some minor changes. The loss of a message during the multistage timestamp generation phase, for instance, could lead a publisher to wait forever before publishing an event in the ENS. This problem can be solved with a simple retransmission approach: the publisher periodically re-initiates the procedure for building the timestamp until it receives a correct timestamp for the event. The same solution can be applied for subscription/unsubscription timestamp requests as well. Finally, the internal state of $TM$s should be preserved despite possible process failures in order to avoid possible TNO violations. This can be obtained by adopting standard replication techniques [20], [21].

**Ordering Features** - A total ordering imposed on event notifications is important to provide distinct subscribers with the same view of the evolution of an application. However, this kind of ordering does not capture the natural *cause-effect* relationship that may relate events. This kind of relationship may be extremely important in applications where subscribers are expected to observe a coherent evolution of events with respect to a given application-level semantics. Specifically, here we are interested in analyzing if our solution may be used to also guarantee the causal ordering of events produced by publishers [22]. This order is particularly relevant because it allows to recognize *cause-effect* relationships among published events [23].

Our solution is clearly able to guarantee causally ordered notifications within each single topic. If a publisher publishes two events in the same topic, in fact, it will sequentially request two timestamps for the same topic whose content, by construction, will guarantee the FIFO-order as defined at the publisher side. If a process is notified about an event and then publishes a new event in the same topic, this event's timestamp will be greater than the one attached to the notified event.

If events are published in different topics, their timestamps are generally not comparable, thus causal ordering cannot be enforced. Ensuring this ordering imposes some slight modifications to the ordering algorithm introduced in Section 3.4. In particular, the definition of $\mathcal{SG}_\mathcal{T}$ should be revised as follows: the *sequencing group* of a topic $T$ ($\mathcal{SG}_T$) is a set of topics including $T$ and all $T'$ such that there is at least a subscription including both $T$ and $T'$. This update changes the behaviour of the multistage sequencer and produces timestamps that, given an event $e$ published on topic $T$, totally order it with respect to all events published in topics subscribed by subscribers that will be notified about $e$. Thanks to this change, we can guarantee causal order also in the case there is a subscriber subscribed to $T$ and $T'$, and

a second subscriber of $T$ that is notified of an event $e$ in that topic, and afterwards publishes an event $e'$ in topic $T'$. The new structure of the timestamps in this case guarantees that $\mathcal{SG}_\mathcal{T} \cap \mathcal{SG}_{\mathcal{T'}} \supseteq \{\mathcal{T}, \mathcal{T'}\}$, and thus that event $e'$ published in $T'$ will get a timestamp larger than the one associated to $e$.

# 5 PERFORMANCE EVALUATION

In this section we evaluate the performance of the proposed solution. First, we describe the system deployment, the metrics of interest for the evaluation, and the load used to stress the system. Then, we analyze how the characteristics of our timestamps (i.e., their sizes) vary depending on different application loads. Finally, we evaluate the overall performance of a prototype implementing our solution.
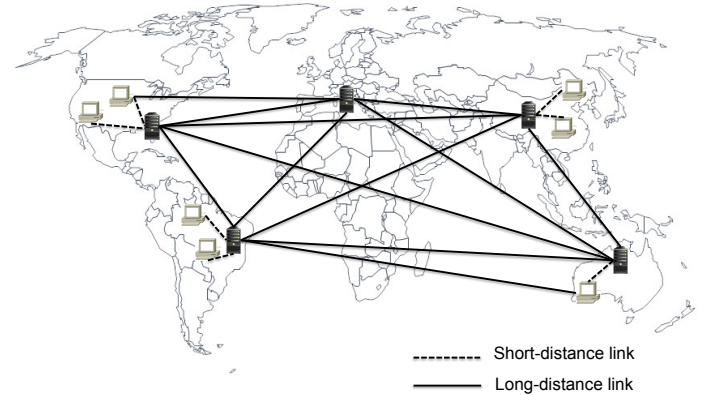


Fig. 6. A geo-distributed system interconnected by long-distance links. Clients connect to servers by means of two kinds of link: short- and long-distance.

## 5.1 System deployment

We consider the scenario depicted in Figure 6: a geo-distributed infrastructure using an asynchronous messaging layer to convey client updates that have to be applied in a consistent order. This order is provided by our ordering solution. Servers are located in dispersed geographic *sites* and connected to each other through high-latency links, while clients connect to servers through either low- or high-latency links. These links represent short- or long-distance WAN connections, respectively, and refer to the geographic distance between a client and a server or between two servers. We assume that the geo-distributed application running on top of this infrastructure implements an event-based communication pattern leveraging an underlying reliable topic-based publish/subscribe middleware. Our ordering solution is deployed on top of this middleware such that each server in the ENS hosts the $TM$s of the topics that are more popular in terms of publications and subscriptions in that geographic location. This represents a typical scenario, in which the topic popularity follows a locality

principle, as the one shown in [13]. As an example, consider a geo-distributed system that disseminates news, results, and statistics about soccer. Servers in Italy, US, Brazil host the $TMs$ of topics that are more popular in these countries, such as *teams, results, players, records, ...* of the Italian, American, and Brazilian leagues, respectively. A client in one of those countries issues with higher probability publications/subscriptions to topics related to the local league. In the next subsection we will discuss the publication/subscription model that captures this scenario.

We implement a prototype of the ordering protocol and deploy it on 5 physical servers. Each server is a Dell PowerEdge R210 II, with an Intel Xeon E3 1270v2 3.50 GHz processor and 16GB of memory. We consider a publish/subscribe system with 1000 topics, with the precedence relation $\rightarrow$ defined by the topic identifier. $TMs$ are equally partitioned onto servers (i.e., each server hosts 200 $TMs$). We assume that $TMs$ of topics with id 1-200 are on server 1, $TMs$ of topics with id 201-400 are on server 2, and so on. We deploy 1000 clients into 5 physical machines. We use the *netem* network emulator to emulate short- and long-distance WAN links. A short-distance link has average delay of 10 ms and standard deviation of 2 ms, while a long-distance link has average delay of 100 ms and standard deviation of 10 ms.

## 5.2   Settings and metrics

The following metrics have been considered for our evaluation:

**Timestamp size** Number of entries in a timestamp.
**End-to-end latency** Time taken for the construction of a timestamp. This time includes the timestamp request by a publisher, the actual timestamp generation by $TMs$, and the response issued by some $TM$.
**Throughput** The number of timestamps generated by $TMs$ in a time unit.
**Percentage of outgoing bandwidth** The fraction of bandwidth that a physical server that hosts one or more $TMs$ reserves to forward timestamps to other servers during their generations. This fraction is computed as the number of bytes of timestamps forwarded to other servers divided the total number of bytes of all timestamps handled in that server during an experiment.

The tests were performed by varying two basic parameters: event rate (number of timestamp requests per time unit) and total number of subscribed topics (global sum of the size of all subscriptions).

The publication and subscription models *at each site* of the geo-distributed system were varied as follows:

**Publication model** We model publications as a probability distribution over the set of topics. We consider a *power-law* distribution with shape $0.901$ or $0.349$. The former refers to the $0.5\%$ of topics having a probability of $80\%$ to be selected for a new publication. The latter, instead, refers to the $40\%$ of topics having a probability of $80\%$ to be selected for a new publication;
**Subscription model** As for publications, subscriptions are modeled as a probability distribution over the set of topics. Again, we consider *power-law* distributions with shapes $0.901$ and $0.349$.

The distribution with shape $0.901$ represents a model in which a few topic per each site have high popularity within a geographic region, and are subscribed with very high probability by clients of the same region. Hence, just a small percentage of the clients of a geographic region subscribe to topics that are outside their region. In other words, most of the content generated in a certain region is mostly consumed in the same region. This captures the *geographic topic popularity* pattern described in [13] in the context of YouTube video popularity. Considering the prototype deployment, this means that most of the computation for timestamp generation is local to each server. This property is confirmed in the performance study of Section 5.4.

The distribution with shape $0.349$ represents a model in which a few topic per each site have high popularity within a geographic region. However, differently from the previous distribution, these topics are subscribed by clients that can be of any region as topics might diffuse to other geographical regions along time. This phenomenon captures the *spray-and-diffuse* pattern described in [13] and [14] in the context of YouTube videos and twitter hastags respectively. In this case, considering the prototype deployment, several geo-distributed servers could be involved in timestamp generation.

## 5.3   Timestamp size evaluation

We first analyze the scalability of the proposed timestamping technique over the number of subscribed topics by calculating the average size of timestamps. Note that this size only depends on issued subscriptions and is completely independent from the specific deployment scenario. For this reason, and only for the evaluation of this aspect, we simplified the setting described in section 5.1 by considering a single server hosting all the $TMs$. The evaluation is conducted in a system with 10000 clients and 1000 available topics, with the total number of subscribed topics that varies in the range [10k - 2000k]. Topic ranking was defined to match the topic popularity as defined by the publication and subscription models. Note that in this setting, geographic topic popularity and spray-and-diffuse boil down to popularity topic distribution with two different shapes.

Figure 7 shows the mean timestamp size by varying the number of subscribed topics. The leftmost point of both curves represents the case where each of the 10000 available clients, acting as subscribers, subscribe 10 topics. In both cases the average timestamp size is in the same order of magnitude of the number of subscribed
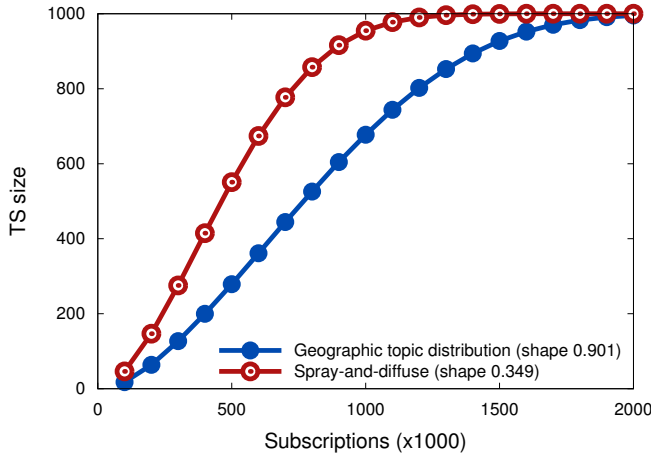
Fig. 7.  Mean timestamp size varying the number of subscribed topics. This experiment considers all $TM_i$ hosted in a single site.



Fig. 8.  Standard deviation of the mean timestamp size varying the number of subscribed topics. This experiment considers all $TM_i$ hosted in a single site.

topics per clients (46 and 17 topics for the curves with shape 0.349 and 0.901 respectively). These size would both decrease to 0 by further reducing the amount of subscribed topics to 1 for subscriber as there would be no intersections among subscriptions. On the right side of the picture, as the number of topics subscribed grows approaching the maximum of 1000 per subscriber, both curves asymptotically approach their maximum. It is worth noticing that (i) more skewed popularity distributions (i.e., shape 0.901) produce smaller timestamps as less popular topics are rarely subscribed by more than one subscriber, and (ii) in both cases the size of the timestamps remains fairly small in meaningful scenarios where the number of topics subscribed per subscriber is not huge[3]. Fig 8 reports the standard deviation for the same experiments.

### 5.4   Prototype performance evaluation

Figure 9 shows the mean end-to-end latency by varying the event rate, i.e., the number of timestamp requests per second. The subscription size for each client is 10 topics (10k topics subscribed in total). Figure 9 evidences that our protocol imposes a small delay for timestamp generation when subscriptions include with very high probability topics whose $TMs$ reside on the same physical host (power-law distribution with shape 0.901). In this case most of the computation for a timestamp generation is local as discussed in Section 5.2. In addition, according to the publication model, the timestamp requests are more likely to come from publishers that are geographically closer. Hence, low-distance links are used most of the time. On the contrary, when publication and subscription models follow the *spray-and-diffuse* pattern (power-law distribution with shape 0.349), more physical machines are involved in the timestamp generation and the mean

3. Typically, scenarios where most subscriber subscribe a vast majority of the available topics are better served by broadcast primitives rather than publish/subscribe ones.
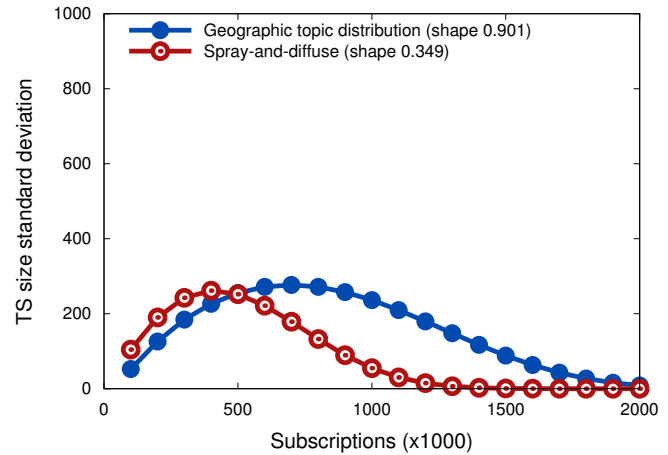
end-to-end latency is affected by communication on long-distance links.
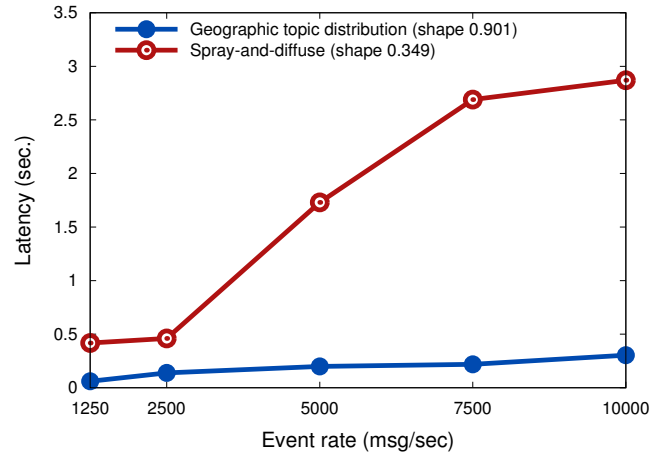


Fig. 9.  Mean end-to-end latency for different publication and subscription models by varying the event rate.

Similar results are obtained when measuring the throughput, i.e., the number of timestamps generated in a second. Our protocol sports its peak performance in the proposed deployment when fed with 5000 timestamp requests per second for both probability distributions. When these requests follow the *geographic topic popularity* pattern the system shows its best performance, reaching a peak of more that 4k timestamps generated per second.

The rationale behind the results showed above is confirmed by Figures 11 and 12, which depict the percentage of outgoing bandwidth for each physical server in the presence of *geographic topic popularity* and *spray-and-diffuse* patterns, respectively. When submissions show *geographic topic popularity*, the $TMs$ involved in the generation of a timestamp are likely to be hosted by the same physical server. Similarly, under that pattern, the vast majority of timestamp requests are issued on low-
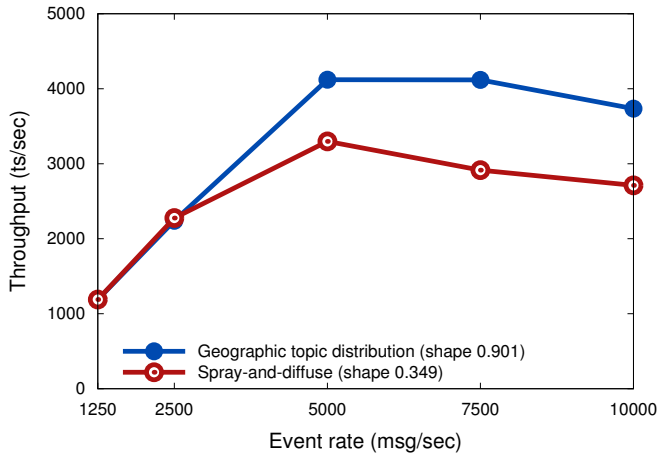
Fig. 10. Mean throughput for different publication and subscription models by varying the event rate.
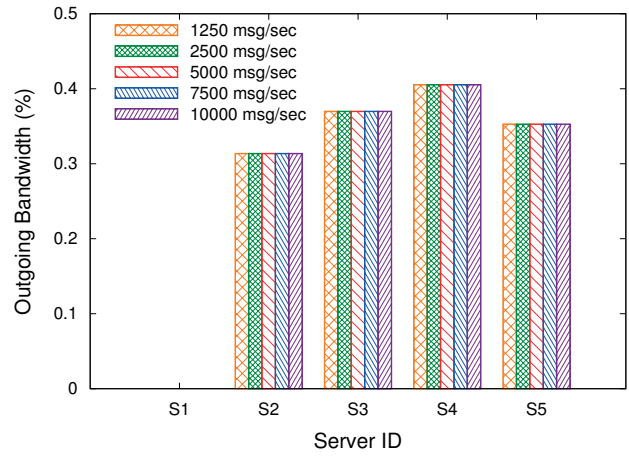


Fig. 12. Percentage of outgoing bandwidth in case of *spray-and-diffuse* pattern (power-law distribution with shape 0.349).

distance WAN links. Figure 11 evidences how only $TMs$ on servers 4 and 5 forward timestamp requests to $TMs$ on different servers. Indeed, $TMs$ on servers 4 and 5 are the $TMs$ of topics with lower precedence according to their identifiers. The percentage of outgoing bandwidth is negligible for server 3 and null for servers 1 and 2.

On the contrary, when subscriptions and publications follow the *spray-and-diffuse* pattern, the percentage of outgoing bandwidth for each server is considerably higher (except for server 1).
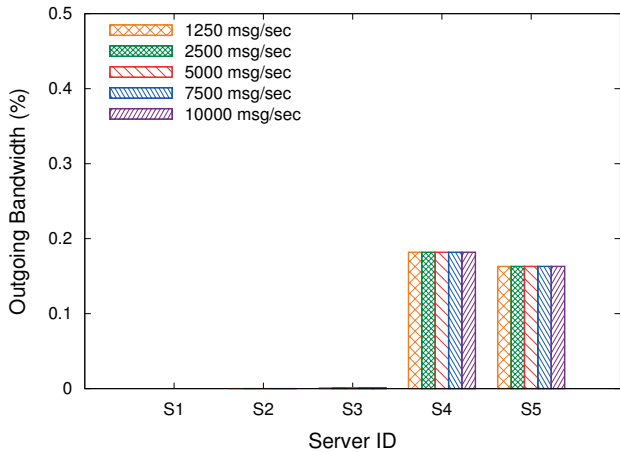
col when the subscription model follows the *geographic topic popularity* pattern. A larger number of subscriptions results just in some additional local computations. On the contrary, when subscriptions follow the *spray-and-diffuse* pattern, a larger number of subscriptions impose computations needed to generate timestamps which can involve multiple sites.
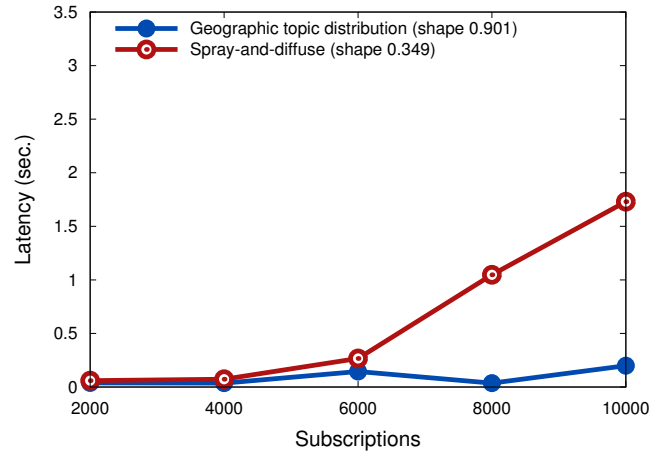


Fig. 13. Mean end-to-end latency for different publication and subscription models by varying the number of subscriptions.



Fig. 11. Percentage of outgoing bandwidth in case of *geographic topic popularity* pattern (power-law distribution with shape 0.901).

Figures 13 and 14 show the mean end-to-end latency and throughput by varying the total number of subscriptions in the system and fixing the event rate to 5000 timestamp requests per second. As the previous experiments, even in this case the publication and subscription models affect the performance of the system. The interesting fact, however, is that the number of subscriptions does not impact the performance of the proto-

## 6 RELATED WORK

Although many real world applications require support for QoS such as ordering, the majority of current publish/subscribe solutions mainly operates on a best-effort basis [24]. Among the QoS-enabled event dissemination services, JEDI [25] is a publish/subscribe middleware that satisfies a causal order delivery of events. It is obtained by means of a *return value*, i.e., a response
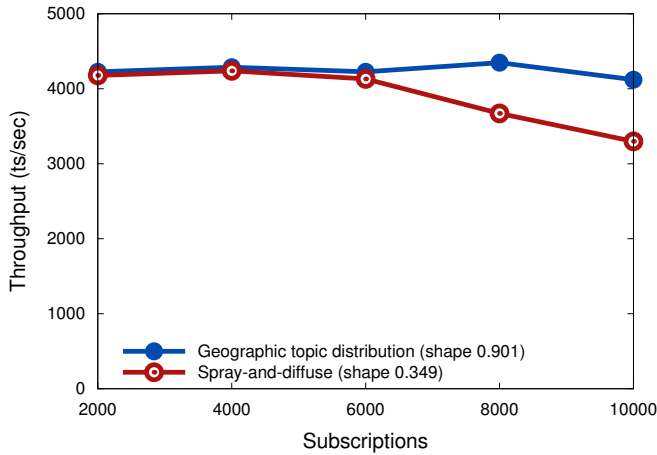
Fig. 14. Mean throughput for different publication and subscription models by varying the number of subscriptions.

message that a receiver uses to notify an event delivery to the producer. This mechanism is clearly not scalable in the presence of a high number of nodes and high event rate. A total order algorithms for publish/subscribe systems is presented in [2]. Similar to our work, authors use a *sequencing network* to order events across multiple groups of subscribers. However, their solution is not able to handle subscription dynamics. A new subscription/unsubscription can create loops in the sequencing network (*circular dependency problem*), which, thus, must be rebuilt from scratch. On the contrary, our solution solves these problems by defining a total order relation among topics that determines a one-way sequence of topic managers that establish an order for events.

Two interesting solutions for ordering in content-based publish/subscribe middleware recently appeared in [5] and [26]. The paper in [26] presents a partition-tolerant content-based publish/subscribe algorithm that can tolerate concurrent failures of brokers or communication links up to a value $\delta$. However, this solution only guarantees a per-source FIFO order. The authors of [5] introduce a *Pairwise Total Order* specification that, in its *weak* variant, fully matches our TNO property. Both solutions follow the same path: discard liveness (delivery reliability) in favor of safety (ordering). However, the system presented in [5] supports the content-based event selection model by deeply integrating within the PADRES publish/subscribe system. Its applicability to other ENSs is thus limited. Conversely, our solution is completely decoupled from the underlying messaging layer and it can be easily adapted to several different contexts (i.e., ENS or group toolkit). In [5], the correct order of events is reconstructed by brokers, by using advertisements and subscriptions to detect *conflicts* and buffers to maintain them during the resolution phase. Conflicts are resolved through an acknowledgment mechanism. Note that this solution increases the load on brokers due to the widespread usage of ACK messages that enforce a tight synchronization among multiple brokers,

hampering the ability of the system to support strong loads [11]. To this respect, our solution has been natively designed by enforcing a *one-way message flow* design within the multistage sequencer that, by removing explicit synchronization among $TMs$, helps supporting growing loads and thus scales gracefully even if the order of events must be reconstructed on subscribers' side. Recently, Yahoo started adopting HedWig [27] as a topic-based publish/subscribe solution for event dissemination within the PNUTS system. HedWig provides strong data delivery guarantees but, contrarily to our solution, it only enforces per-topic ordering within a single datacenter. Therefore, topics are independent in HedWig, as there is no way to order notifications related to different topics.

## 7 CONCLUSIONS

Totally ordering notifications produced by a publish/subscribe system deployed on top of a distributed infrastructure is a complex problem due to the inherent dynamism of pub/sub interactions, and to the fact that the order has to be achieved in a fully decoupled way, without any explicit synchronization among components of the systems itself.

This paper presented an ordering solution that can be used to order notifications delivered by a reliable event notification service to satisfy a *total notification order* specification. The main contribution is a timestamping mechanism that considers topics and overlapping subscriptions to automatically build appropriately sized timestamps that are attached to published events. Subscribers can use timestamps to infer the correct sequence of notifications to be enforced locally. The ordering mechanism uses a multistage sequencer to build timestamps. A precedence relationship among topics is used to determine the order in which stages are executed, i.e., to determine a one-way sequencing path (no loops or feedbacks that may create instability in the network) for the timestamp generation. Thanks to its internal structure the sequencer can be deployed in multiple flexible ways: several stages can be co-located in the same server of a distributed system thus making local to that server the entire generation of several timestamps.

The performance of the proposed solution has been evaluated through a prototype implementation. In particular, the results show that our solution is able to deliver very good performance for geo-distributed applications characterized by a load following a *geographic topic distribution* pattern where user interests are geographically clustered.

like to thank the anonymous reviewers for all their comments that helped to improve both presentation and content of the paper.

# REFERENCES

[1] C. Liebig, M. Cilia, and A. P. Buchmann, "Event composition in time-dependent distributed systems," in *Proceedings of the 4th International Conference on Cooperative Information Systems (IFCIS), Edinburgh, Scotland, September 2-4.* IEEE, 1999, pp. 70–78.

[2] C. Lumezanu, N. Spring, and B. Bhattacharjee, "Decentralized message ordering for publish/subscribe systems," in *Proceedings of the 7th ACM/IFIP/USENIX International Middleware Conference (Middleware), Melbourne, Australia, November 27 - December 1*, 2006, pp. 162–179.

[3] A. Malekpour, A. Carzaniga, G. T. Carughi, and F. Pedone, "Probabilistic FIFO ordering in publish/subscribe networks," Faculty of Informatics, University of Lugano, Technical Report USI-INF-TR-2011-2, Tech. Rep., 2011.

[4] H. Garcia-Molina and A. Spauster, "Ordered and reliable multicast communication," *ACM Transaction on Computer Systems*, vol. 9, no. 3, pp. 242–271, 1991.

[5] K. Zhang, V. Muthusamy, and H.-A. Jacobsen, "Total order in content-based publish/subscribe systems," in *Proceedings of the 32nd International Conference on Distributed Computing Systems (ICDCS), Macau, China, June 18-21.* IEEE, 2012.

[6] P. R. Pietzuch, B. Shand, and J. Bacon, "Composite event detection as a generic middleware extension," *IEEE Network*, vol. 18, no. 1, pp. 44–55, 2004.

[7] A. R. Bharambe, S. G. Rao, and S. Seshan, "Mercury: a scalable publish-subscribe system for internet games," in *Proceedings of the 1st Workshop on Network and System Support for Games (NETGAMES), Braunschweig, Germany, April 16-17.* ACM, 2002, pp. 3–9.

[8] Y. Tock, N. Naaman, A. Harpaz, and G. Gershinsky, "Hierarchical clustering of message flows in a multicast data dissemination system." in *IASTED PDCS*, 2005, pp. 320–326.

[9] RTI, "Real time messaging and integration middleware," http://www.rti.com/.

[10] L. Lamport, "Paxos made simple," *ACM SIGACT News*, vol. 32, no. 4, pp. 18–25, 2001.

[11] K. Birman, G. Chockler, and R. van Renesse, "Toward a cloud computing research agenda," *SIGACT News*, vol. 40, no. 2, pp. 68–80, 2009.

[12] R. Barrett, "Transactions across datacenters," Google I/O developer conference http://www.google.com/events/io/2009/sessions/TransactionsAcrossDatacenters.html, 2009.

[13] A. Brodersen, S. Scellato, and M. Wattenhofer, "YouTube around the world: geographic popularity of videos," in *Proceedings of the 21st international conference on World Wide Web.* ACM, 2012, pp. 241–250.

[14] K. Y. Kamath, J. Caverlee, K. Lee, and Z. Cheng, "Spatio-temporal dynamics of online memes: a study of geo-tagged tweets," in *Proceedings of the 22nd international conference on World Wide Web*, 2013, pp. 667–678.

[15] R. Baldoni, S. Cimmino, and C. Marchetti, "A classification of total order specifications and its application to fixed sequencer-based implementations," *Journal of Parallel and Distributed Computing*, vol. 66, no. 1, pp. 108–127, 2006.

[16] X. Defago, A. Schiper, and P. Urban, "Total order broadcast and multicast algorithms: Taxonomy and survey," *ACM Computer Survey*, vol. 36, no. 4, pp. 372–421, 2004.

[17] R. Baldoni, C. Marchetti, A. Virgillito, and R. Vitenberg, "Content-based publish-subscribe over structured overlay networks," in *Proceedings of the 25th International Conference on Distributed Computing Systems (ICDCS), Columbus, OH, USA, June 6-10.* IEEE, 2005, pp. 437–446.

[18] M. Castro, P. Druschel, A.-M. Kermarrec, and A. I. Rowstron, "Scribe: A large-scale and decentralized application-level multicast infrastructure," *IEEE Journal on Selected Areas in Communications*, vol. 20, no. 8, pp. 1489–1499, 2002.

[19] K. Birman, "Rethinking multicast for massive-scale platforms," in *Proceedings of the 29th IEEE International Conference on Distributed Computing Systems (ICDCS), Montreal, Canada, June 22-26.* IEEE Computer Society, 2009.

[20] N. Budhiraja, K. Marzullo, F. B. Schneider, and S. Toueg, "Distributed systems (2nd edition)," S. Mullender, Ed. New York, NY, USA: ACM Press/Addison-Wesley Publishing Co., 1993, ch. The primary-backup approach, pp. 199–216.

[21] F. B. Schneider, "Distributed systems (2nd edition)," S. Mullender, Ed. New York, NY, USA: ACM Press/Addison-Wesley Publishing Co., 1993, ch. Replication management using the state-machine approach, pp. 169–198.

[22] K. P. Birman and T. A. Joseph, "Reliable communication in the presence of failures," *ACM Transaction on Computer Systems*, vol. 5, no. 1, pp. 47–76, 1987.

[23] L. Lamport, "Time, clocks, and the ordering of events in a distributed system," *Commununication of the ACM*, vol. 21, no. 7, pp. 558–565, 1978.

[24] G. Muehl, L. Fiege, and P. R. Pietzuch, *Distributed event-based systems.* Springer, 2006.

[25] G. Cugola, E. D. Nitto, and A. Fuggetta, "The JEDI event-based infrastructure and its application to the development of the OPSS WFMS," *IEEE Transaction on Software Engineering*, vol. 27, no. 9, pp. 827–850, 2001.

[26] R. S. Kazemzadeh and H.-A. Jacobsen, "Partition-tolerant distributed publish/subscribe systems," in *Proceedings of the 30th IEEE Symposium on Reliable Distributed Systems (SRDS), Madrid, Spain, October 4-7.* IEEE, 2011, pp. 101–110.

[27] HedWig, http://wiki.apache.org/hadoop/HedWig.

**Roberto Baldoni** is the director of the Sapienza Research Center for Cyber Intelligence and Information Security and, at National level, he is director of the Cyber Security National Laboratory. In 2013-2014, Roberto has been Chair of the IEEE Technical committee on Dependable Computing and Fault Tolerance and Chair of the steering committee of IEEE Dependable Systems and Networks Conference. Roberto conducts research (from theory to practice) in the fields of distributed, pervasive and P2P computing, middleware platforms and information systems infrastructure with a specific angle on dependability and security aspects.

**Silvia Bonomi** is a PhD in Computer Science at Sapienza University of Rome. She is doing research on various computer science fields including dynamic distributed systems and event-based systems. In these research fields, she published several papers in peer reviewed scientific forums. As a part of the MIDLAB research group, she has been involved in an EU-funded project dealing with energy saving in private and public buildings (GreenerBuildings project) and she worked on dependable distributed systems (ReSIST network of excellence) and on the definition of new semantic tools for e-Government (SemanticGov).

**Marco Platania** is a Postdoctoral Fellow at the Department of Computer Science, Johns Hopkins University (Baltimore, USA). He got the PhD in Computer Science Engineering in 2012 at Sapienza University of Rome. His main research interests are resilient clouds, publish/subscribe architecture, critical infrastructure protection, and P2P systems. He regularly serves as a reviewer in several leading journals and conferences in the field of Distributed Systems.

**Leonardo Querzoni** is an assistant professor at Sapienza University of Rome. He obtained a PhD in in computer engineering with a work on efficient data dissemination though the publish/subscribe communication paradigm. His research activities are focused on several computer science fields including distributed stream processing, parallel computing, large scale and dynamic distributed systems, dependability in distributed systems, publish/subscribe middleware services.