

Universal Cross-Cloud Communication

Chamikara Jayalath, Julian Stephen, and Patrick Eugster

Abstract—Integration of applications, data-centers, and programming abstractions in the cloud-of-clouds poses many challenges to system engineers. Different cloud providers offer different communication abstractions, and applications exhibit different communication patterns. By abstracting from hardware addresses and lower-level communication, the *publish/subscribe* paradigm seems like an adequate abstraction for supporting communication across clouds, as it supports many-to-many communication between publishers and subscribers, of which one-to-one or one-to-many can be viewed as special cases. In particular, *content-based publish/subscribe* (CPS) systems provide an expressive abstraction that matches well with the key-value pair model of many established cloud storage and computing systems, and decentralized overlay-based CPS implementations scale up well. However, CPS systems perform poorly at small scale, e.g., one-to-one or one-to-many communication. This holds especially for multi-send scenarios which we refer to as *entourages* that may range from a channel between a publisher and a single subscriber to a broadcast between a publisher and a handful of subscribers. These scenarios are common in cloud computing, where cheap hardware is exploited for parallelism (efficiency) and redundancy (fault-tolerance). With CPS, multi-send messages go over several hops before their destinations are even identified via predicate matching, resulting in increased latency, especially when destinations are located in different data-centers or zones. *Topic-based publish/subscribe* (TPS) systems support communication at small scale more efficiently, but still route messages over multiple hops and inversely lack the flexibility of CPS systems. In this paper, we propose CPS protocols for cloud-of-clouds communication that can dynamically identify entourages of publishers and corresponding subscribers. Our CPS protocols dynamically connect the publishers with their entourages through *überlays*. These überlays can transmit messages from a publisher to its corresponding subscribers with low latency. Our experiments show that our protocols make CPS abstraction viable and beneficial for many applications. We introduce a CPS system named Atmosphere that leverages out CPS protocols and illustrate how Atmosphere has allowed us to implement, with little effort, versions of the popular HDFS and ZooKeeper systems which operate efficiently across data-centers.

Index Terms—Publish/subscribe, content-based, unicast, multicast, multi-send

1 INTRODUCTION

THE advent of *cloud brokers* [19] for mediating between different cloud providers, and the *cloud-of-clouds* [9] paradigm denoting the integration of different clouds—including clouds offering different abstractions (e.g., infrastructure as a service versus platform as a service)—pose new challenges to software developers. Indeed, while support for developing specific types of applications to run in different individual cloud infrastructures is slowly becoming established, there is little support for programming applications that run across several clouds or types of clouds.

1.1 Cross-Cloud Communication

One particular aspect of such integration, addressed herein, is communication. The cloud-of-cloud paradigm postulates integration of multiple data-centers, cloud abstractions, and applications. Can a single communication middleware fulfill all different arising needs for communication? In particular, a candidate middleware system must be able to

- R1. support a variety of *communication patterns* (e.g., communication rate, number of interacting entities)

effectively. Given the variety of target applications (e.g., social networking, web servers), the system must be able to cope with one-to-one communication as well as different forms of multicast (one-to-many, many-to-many). In particular, the system must be able to scale up *and* down (“elasticity”) based on current needs [27] such as number of communicating endpoints;

- R2. run on standard “low-level” network layers and abstractions without relying on any specific protocols such as IP Multicast [17] that may be deployed in certain clouds but not in others or across clouds [36];
- R3. provide an interface which hides cloud-specific hardware addresses and integrates well with abstractions of widespread cloud storage and computing systems in order to support a wide variety of applications;
- R4. operate efficiently despite varying network latencies within/across datacenters.

- The authors are with the Department of Computer Science, Purdue University, West Lafayette 47907, IN.
E-mail: {cjayalat, stephe22, peugster}@cs.purdue.edu.

Manuscript received 15 Sept. 2013; revised 11 Feb. 2014; accepted 5 Mar. 2014. Date of publication 10 Apr. 2014; date of current version 30 July 2014. Recommended for acceptance by I. Bojanova, R.C.H. Hua, O. Rana, and M. Parashar.

For information on obtaining reprints of this article, please send e-mail to: reprints@ieee.org, and reference the Digital Object Identifier below.
Digital Object Identifier no. 10.1109/TCC.2014.2316813

1.2 Publish/Subscribe for the Cloud

One candidate abstraction is *publish/subscribe*. Components act as *publishers* of messages, and dually as *subscribers* by delineating messages of interest. Examples of publish/subscribe services designed for and/or deployed in the cloud include Amazon’s simple notification service (SNS) [2], LinkedIn’s Kafka [26], Hadoop Hedwig [20], or Blue Dove [27]. Intuitively, publish/subscribe is an

adequate abstraction for cross-cloud communication because it supports generic many-to-many interaction, shields applications from specific lower-level communication—in particular hardware addresses—thus supporting application interoperability and portability. In particular, *content-based publish/subscribe* (CPS) [11], [31], [18], [8], [33] promotes an addressing model based on message *properties* and corresponding values (with subscribers delineating values of interest for relevant properties) which matches well the *key-value* pair abstractions used by many cloud storage (e.g., [16], [14]) and computing (e.g., [15]) frameworks.

1.3 Limitations

However, existing publish/subscribe systems for the cloud are not designed to operate beyond single datacenters, and CPS systems focus on scaling *up* to large numbers of subscribers: to “mediate” between published messages and subscriptions, CPS systems typically employ an overlay network of brokers, with filtering happening downstream from publishers to subscribers based on upstream aggregation of subscriptions. When messages from a publisher are only of interest to one or few subscribers, such overlay-based multi-hop routing (and filtering) will impose increased latency compared to a direct *multi-send* via UDP or TCP from the publisher to its subscribers. Yet such scenarios are particularly wide-spread in third-party computing models, where many cheap resources are exploited for parallelism (efficiency) or redundancy (fault-tolerance). A particular example are distributed file-systems, which store data in a redundant manner to deal with crash failures [3], thus leading to frequent communication between an updating component and (typically 3) replicas. Another example for multi-sends are (group) chat sessions in social networks.

Existing approaches to adapting interaction and communication between participants based on *actual* communication patterns (e.g., [37], [27], [35]) are agnostic to deployment constraints such as network topology. *Topic-based publish/subscribe* (TPS) [11], [6]—where messages are published to *topics* and delivered to consumers based on topics they subscribed to—is typically implemented by assigning topics to nodes. This limits communication hops in multi-send scenarios, but also the number of subscribers.

In short, CPS is an appealing, generic, communication abstraction (R2, R3), but existing implementations are not efficient at small scale (R1), or, when adapting to application characteristics, do not consider deployment constraints in the network (R4); inversely, TPS is less expressive than CPS, and existing systems do not scale up as well.

1.4 Atmosphere

This paper describes Atmosphere, a middleware solution that aims at supporting the expressive CPS abstraction across datacenters and clouds in a way which is effective for a wide range of communication patterns. Specifically, our goal is to support the extreme cases of communication between individual publisher-subscriber pairs (unicast) and large scale CPS, and to elastically scale both up *and*

down between these cases, whilst providing performance which is comparable to more specialized solutions for individual communication patterns. This allows applications to focus on the logical content of communication rather than on peer addresses even in the unicast case: application components need not contain hardcoded addresses or use corresponding deployment parameters, as the middleware automatically infers associations between publishers and subscribers from advertisements and subscriptions.

Our approach relies on a CPS-like peer-based overlay network which is used primarily for “membership” purposes, i.e., to keep participants in an application connected, and as a fallback for content-based message routing. The system dynamically identifies clusters of publishers and their corresponding subscribers, termed *entourages* while taking network topology into account. Members of such entourages are connected directly via individual overlay networks termed *überlays*, so that they can communicate with low latency. The überlay may only involve publishers and subscribers or may involve one or many brokers depending on entourage characteristics and resource availabilities of involved publishers, subscribers, brokers, and network links. In any case, these *direct connections* which are gradually established based on resource availabilities, will effectively reduce the latency of message transfers from publishers to subscribers.

1.5 Contributions

Several CPS systems have been proposed in the literature, and Atmosphere adopts several concepts presented therein. Thus, in the present paper, we focus on the following novel contributions of Atmosphere:

- 1) a protocol to dynamically identify topic-like entourages of publishers in a CPS system. Our technique hinges on precise advertisements. To not hamper flexibility, advertisements can be updated at runtime;
- 2) a protocol to efficiently and dynamically construct überlays interconnecting members of entourages with low latency based on resource availabilities;
- 3) the implementation of a scalable fault-tolerant CPS system for geo-distributed deployments named Atmosphere that utilizes our entourage identification and überlay construction techniques;
- 4) an evaluation of Atmosphere using real-life applications, including social networking, news feeds, the HDFS [3] distributed file-system, and the ZooKeeper [4] distributed lock service demonstrating the efficiency and versatility of Atmosphere through performance improvements over more straightforward approaches. In particular, we describe how Atmosphere was instrumental in building, without much effort, versions of HDFS and ZooKeeper operating efficiently across datacenters.

1.6 Roadmap

Section 2 provides background information. Section 3 presents our protocols. Section 4 introduces Atmosphere. Section 5 evaluates our solution. Section 6 presents related work. Section 7 draws conclusions.

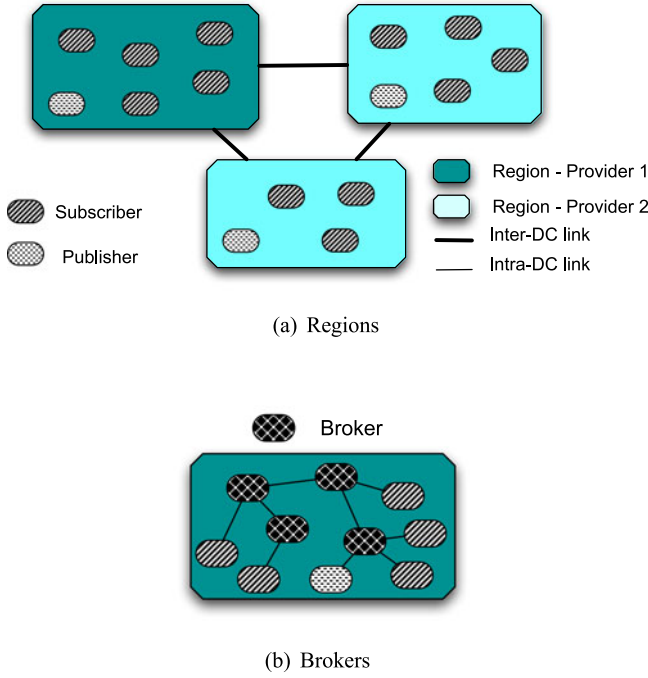


Fig. 1. Bird's-eye view.

2 BACKGROUND

This section presents background information and the model of system considered.

2.1 System Model

We assume a system G of processes communicating via unicast channels spanning g cloud datacenters or more generally *regions*. Regions may be operated by different cloud providers. Each region contains a number of components that produce messages (publishers) and/or that are interested in consuming messages (subscribers) produced. Fig. 1a shows an example system with three regions from two different providers where each region hosts a single producing and multiple consuming components.

2.2 CPS Communication

With *content-based publish/subscribe* (CPS), a message produced by a publisher contains a set of *property-value* pairs; inversely, components engage in consumption of messages by issuing subscriptions which consist in ranges of values—typically defined indirectly through operators such as \leq or \geq and corresponding threshold values.

A *broker overlay network* typically mediates the message distribution between publishers and subscribers. A broker, when receiving a message, analyzes the set of property-value pairs, and forwards the message to its neighbors accordingly. (For alignment with the terminology used in clouds we may refer to properties henceforth as *keys*.) Siena [10] is a seminal CPS framework for distributed wide-area networks that spearheaded the above-mentioned CPS overlay model. Siena's routing layer consists of broker nodes that maintain the interests of sub-brokers and end hosts connected to them in a *partially ordered set* (poset) structure. The root of the poset is sent to the *parent broker* to which a given broker is subscribed to. CPS systems like Siena

employ *subscription summarization* [11], [29] for brokers to construct a summary of the interests of the subscribers and brokers connected to it. This summary is sent to neighboring brokers. A broker that receives a published message determines the set of neighbors to which the message has to be forwarded by analyzing the corresponding subscription summaries. Summaries are continuously updated to reflect the changes to the routing network, occurring for instance through joins, leaves, and failures of subscribers or brokers.

2.3 Existing CPS System Limitations

When deployed naïvely, i.e., without considering topology, in the considered multi-region model (see Fig. 1a) CPS overlays will perform poorly especially if following a DAG as is commonly the case, due to the differences in latencies between intra- and inter-region links. To cater for such differences, a broker network deployed *across regions* could be set up such that (1) brokers in *individual regions* are *hierarchically* arranged and each subscriber/publisher is connected to exactly one broker (see Fig. 1b), and (2) root brokers of individual regions are connected (no DAG). The techniques proposed shortly are tailored to this setup.

However, the problem with such a deployment is still that—no matter how well the broker graph matches the network topology—routing will occur in most cases over multiple hops which is ineffective for multi-send scenarios where few subscribers only are interested in messages of some publisher. In the extreme case where messages produced by a publisher are consumed by a single subscriber there will be a huge overhead from application-level routing and filtering over multiple hops compared to a direct use of UDP or TCP. The same holds with multiple subscribers as long as the publisher has ample local resources to serve all subscribers over respective direct channels.

While several authors have proposed ways to identify and more effectively interconnect matching subscribers and publishers, these approaches are deployment-agnostic in that they do not consider network topology (or resource availabilities). Thus they trade *logical proximity* (in the message space) for *topological proximity*.

Majumder et al. [28] for instance show that using a single minimum spanning or a Steiner tree will not be optimal for subscriptions with differing interests. They propose a multiple tree-based approach and introduce an approximation algorithm for finding the optimum tree for a given type of publications. But unlike in our approach these trees are location-agnostic hence when applied to our model a given tree may contain brokers/subscribers from multiple regions and a given message may get transmitted across region boundaries multiple times, unnecessarily increasing the transmission latency.

Sub-2-Sub [37] uses gossip-based protocols to identify subscribers with similar subscriptions and interconnect them in an effective manner along with their publishers. In this process, network topology is not taken into account, which is paramount in a multi-region setup with varying latencies. Similarly, Tariq et al. [35] employ spectral graph theory to efficiently regroup and connect components with matching interests, but do not take network topology or latencies into account. Thus these systems cannot be readily

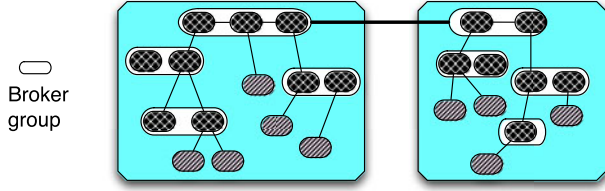


Fig. 2. Broker hierarchies.

deployed across regions. Publi+ [34] introduces a publish/subscribe framework optimized for bulk data dissemination. Similar to our approach, brokers of Publi+ identify publishers and their interested subscribers and instruct them to communicate directly for disseminating large bulk data. Publi+ uses a secondary content-based publish/subscribe network only to connect publishers and interested subscribers in different regions. Publi+ is not designed for dissemination of large amounts of small messages since the data dissemination between publishers and subscribers is always direct and the publish/subscribe network is only used to form these direct connections.

Further related work is discussed in Section 6.

3 ENTOURAGE COMMUNICATION

In this section, we introduce our solution for efficient communication between publishers and “small” sets of subscribers on a two-level geo-distributed CPS network of brokers with hierarchical deployments within individual regions as outlined in Fig. 2 for two regions. This solution can be adapted to existing overlay-based CPS systems characterized in Section 2.2.

3.1 Broker Hierarchies

Consider the number of publishers in region r to be P_r and there to be m_r^p subscribers in region r interested in some messages from a given publisher p . Additionally, assume p to be publishing messages to G at rate f_p and the average size of these messages to be s_p .

Also take W_p to be the maximum bandwidth that a publisher p has at disposition when transferring messages to a remote region, and the percentage utilization of the processor (CPU) of p due to message transmission to be U_p . It is to be noted that while inter-region links usually have ample amounts of bandwidth, most of the time the user will be setting the value of W_p based on the economical factors (e.g., for the inter-region data transfer costs) instead of actual bandwidth limitations. Table 1 lists definitions of some of the notations used in this paper for easy reference.

3.2 Entourages

The range of messages published by a publisher p are identified by advertisements τ_p , which, as is customary in CPS, include the keys and the value ranges for each key. Analogously, the interest range of each subscriber or broker n is denoted by τ_n . $\tau_p \cap \tau_n$ denotes the common interest between a publisher p and a subscriber or broker n .

We define the interest match between a publisher p and a subscriber/broker n as a numerical value that represents the fraction of the publisher’s messages that the subscriber/

TABLE 1
Notations and Definitions

Symbol	Meaning
G	A system deployed across multiple regions
P_r	Number of publishers in region r
m_r^p	Number of subscribers in region r interested in publisher p
f_p	Publication rate of publisher p
s_p	Average size of messages from publisher p
W_p	Bandwidth quota of publisher p
U_p	Processor utilization of publisher p
Φ_p	A set of subscribers/brokers interested in publisher p
τ_p	Interest ranges of messages published by p
τ_n	Interest ranges of broker/subscriber n
ψ	Clustering factor
ch^r	Average churn of region r
ad_l	Average distance to the publisher from level l
cv_l	Coverage of level l
K_p^r	Max # of direct connections from publisher p to region r

broker is interested in assuming the publisher to have an equal probability of publishing a message with any given value within its range.

If the interest range of n is denoted by $\langle key_1, range_1 \rangle, \langle key_2, range_2 \rangle, \dots, \langle key_x, range_x \rangle$ and for each value range $range_i$ of n , the possible value range for the corresponding key key_i of messages published by p is denoted by $range'_i$, then the interest match is given by

$$\prod_{i=1}^x \frac{|range_i \cap range'_i|}{|range_i|}.$$

So, the interest match is defined to be the product of the intersection of the value ranges that correspond to the same key. Note that the published messages may have keys that the subscriber does not specify in which case we assume the subscriber to be interested in all possible values for those keys. We assume that every subscriber has at least one key specified (i.e., $x \geq 1$). The key without a value range has a wildcard for the corresponding value range.

A publisher p and a set Φ_p of subscribers/brokers form a ψ -close entourage if each member of Φ_p has at least a ψ interest match with p where $0 \leq \psi \leq 1$. ψ is a parameter that defines how close the cluster is to a topic. If $\psi = 1$, each member of the cluster is interested in every message published by p , hence the cluster can be viewed as a topic.

3.3 Solution Overview

Next we describe our solution to efficient cross-cloud communication in entourages. The solution consists of three main parts which we describe in turn.

- 1) A decentralized protocol that can be used to identify entourages in a CPS system.
- 2) A mechanism to determine the maximum number K_p of direct connections a given publisher p can maintain without adversely affecting message transmission.
- 3) A protocol to efficiently establish auxiliary networks termed *überlays* between publishers and their respective subscribers.

3.4 Entourage Identification

We describe the *DCI* (dynamic entourage identification) protocol that can be used to identify entourages in a CPS-based

Executed by every broker n

```

1: super                                     {ID of the parent broker}
2: subbrokers                               {Sub-brokers}
3: subscribers                             {Subscribers directly connected to the broker}
4: wait  $\leftarrow 0$                          {# of records to be received by sub-brokers}
5: results  $\leftarrow \emptyset$                  {Results to be sent to the parent broker}

6: when RECEIVE(COUNT,  $p, \tau_p$ ) from  $n'$ 
7:   if rejectCounts( $n'$ ) = true then       {If Count should be rejected}
8:     SEND(COUNTREJECT,  $p$ ) to super       {Send rejection}
9:   end  $\leftarrow$  false                     {Whether will be forwarding COUNT}
10:  for all  $node \in subbrokers \cup subscribers$  do
11:    if interestMatch( $\tau_{node}, \tau_p$ )  $\geq \psi$  then {Sufficient interest}
12:      if  $node \in subbrokers$  then
13:        SEND(COUNT,  $p, \tau_p$ ) to subbroker {Forward COUNT}
14:        wait  $\leftarrow$  wait + 1           {# of results to wait for}
15:      else
16:        results  $\leftarrow$  results  $\cup \{(node, 1)\}$  {Add node}
17:      end
18:    end  $\leftarrow$  true
19:  if |wait| + |results| = 0 then           {No matching nodes found}
20:    end  $\leftarrow$  true
21:  if end = true then
22:    results  $\leftarrow \emptyset$                {Resetting records; any responses discarded}
23:  reply  $\leftarrow$  false
24:  if end = true or (|results| > 1 and wait = 0) then
25:    reply  $\leftarrow$  true                     {Send the COUNTREPLY to parent broker}
26:  if reply = true or |results| + wait > 1 then
27:    results  $\leftarrow$  results  $\cup \{(n, 0)\}$  {Adding current broker}
28:  if reply = true then
29:    SEND(COUNTREPLY,  $p, results$ ) to super {Send reply}

30: when RECEIVE(COUNTREPLY,  $p, results'$ ) from  $n'$ 
31:  if  $\exists \langle n', * \rangle \in results'$  then
32:    for all  $\langle n'', depth \rangle \in results'$  do
33:      results  $\leftarrow$  results  $\cup \{(n'', depth + 1)\}$  {Depth + 1}
34:      wait  $\leftarrow$  wait - 1                 {Have to wait for 1 less record}
35:    if wait = 0 then
36:      SEND(COUNTREPLY,  $p, results$ ) to super {Got all replies}

37: when RECEIVE(COUNTREJECT,  $p$ ) from  $n'$ 
38:  results  $\leftarrow \emptyset$                  {Resetting records; any responses discarded}
39:  results  $\leftarrow$  results  $\cup \{(node, 0)\}$  {Add current node}
40:  SEND(COUNTREPLY,  $p, results$ ) to super {Send reply to parent}

```

Fig. 3. DCI protocol.

application. The protocol assumes the brokers in region i to form a hierarchy, starting from one or more root brokers. An abstract version of the protocol is given in the Fig. 3.

The protocol works by disseminating a message named COUNT along the message dissemination path of publishers. A message initiated by a publisher p contains τ_p and ψ values. Once the message reaches a root node of the publisher's region, it is forwarded to all remote regions.

The brokers implement two main event handlers, (1) to handle COUNT messages (line 6) and (2) to handle replies to COUNT messages – COUNTREPLY messages (line 30).

COUNT messages are embedded into advertisements and carry the keys and value ranges of the publisher. When a broker receives a COUNT message via event handler (1), it first determines the subscribers/brokers directly attached to it that have an interest match of at least ψ with the publisher p (line 11). If there is at least one subscriber/broker with a non-zero interest match that is smaller than ψ then the COUNT message is not forwarded to any child. Otherwise the COUNT (line 13) message is forwarded to all interested children. This is because children with less than ψ interest match are not

considered to be direct members of the p 's entourage and yet messages published by p have to be transmitted to all interested subscribers including ones that have less than ψ interest match. In such a situation, instead of creating direct connections with an ancestor node and some of the descendants, we choose to only establish direct connections with the ancestor node since establishing direct connections with both an ancestor and its descendent will result in duplicate message delivery and unfair latency advantages to a portion of the subscribers. A subscriber or a broker that does not forward a COUNT message immediately creates a COUNTREPLY (line 29) message with its own information and sends it back to the parent.

In the latter event handler (2), a broker aggregates COUNTREPLY messages from its children that have at least a ψ interest match with p (line 33), and sends this information to its respective parent broker through a new COUNTREPLY message (line 36). When performing this aggregation, a broker does not add its own information to a COUNTREPLY message if it previously forwarded the COUNT message to exactly one child. This is because a broker that is only used to transfer traffic between two other brokers or a broker and a subscriber has a child that has the same interest match with p but is hop-wise closer to the subscribers. This child is a better match when establishing an entourage. Aggregated COUNTREPLY messages are ultimately sent to p .

To stop the COUNTREPLY messages from growing indefinitely, a broker may truncate COUNTREPLY messages that are larger than a predefined size M . When truncating, entries from the lowest levels of the hierarchy are removed first. When removing an entry, entries of all its siblings (i.e., entries that have the same parent) are also removed. This is because as mentioned before, our entourage establishment protocol does not create direct connections with both an ancestor node and one of its descendants.

A subscriber or a broker may decide to respond to its parent with a COUNTREJECT (line 8) message instead of a COUNTREPLY either due to policy decisions or local resource limitations. A broker that receives a COUNTREJECT from at least one of its children will discard COUNTREPLY messages for the same publisher from the rest of its children (line 38) and will send COUNTREPLY to the parent that only mentions the current node (line 40).

As a publisher's range of values in published messages evolves, it will have to send new advertisements with COUNT messages to keep its entourage up to date. This is supported in our system Atmosphere presented in the next section by exposing an advertisement update feature in the client API.

3.5 Entourage Size

We devise a heuristic to determine the maximum number of direct connections a given publisher can maintain to its entourage without adversely affecting the performance of transmission of messages.

3.5.1 Factors and Challenges

Capabilities of any node connected to a broker network are limited by a number of factors. A node obviously has to spend processor and memory resources to process and publish a stream of messages. The bandwidth between the node

and the rest of the network could also become a bottleneck if messages are significantly large, or transmitted at a significantly high rate. This is particularly valid in a multi-tenant cloud environment. The transport protocols used by the publisher and latencies to the receivers could limit the rate at which the messages are transmitted.

If the implementation is done in a smart enough way, the increase in memory footprint and the increase in latency due to transport deficiencies can be minimized. The additional memory required for creating data-structures for new connections is much smaller compared to the memory available in today's computers (note that we do not consider embedded devices with significantly low memory capacities). The latencies could become a significant factor if the transport protocol is implemented in a naïve manner, e.g., with a single thread that sends messages via TCP directly to many nodes, one by one. The effect could be minimized by using smarter implementation techniques, e.g., by using features such as multi-threaded transport layers, custom built asynchronous transport protocols, and message aggregation.

Conversely, the processor and bandwidth consumption could significantly increase with the number of unicast channels maintained by a publisher as every message has to be repeatedly transmitted over all connections and every transmission requires CPU cycles and network bandwidth.

3.5.2 Number of Connections

First we determine the increase in processor usage of a given publisher due to establishing direct connections with subscribers or brokers. With each new direct connection, a publisher has to repeatedly send its messages along a new transport channel. So a safe worst case assumption is to suppose that the amount of processing power needs to be proportional to the number of connections over which messages are transmitted.

Additionally, as mentioned previously, a given publisher p will have a bandwidth quota of W_p when communicating with remote regions. Considering both these factors, the number of direct connections K_p which publisher p can establish can be approximated by the expression $\min(\frac{1}{U_p}, \frac{W_p}{f_p \times s_p})$. Please refer to Table 1 for the definitions of the notations used.

This requires the publishers to keep track of their processor utilization; in most of the operating systems, processor utilization can be determined by using system services (e.g., the `top` command in Unix). The above bound on the number of directly connected nodes is not an absolute bound, but rather a initial measure used by any publisher to prevent itself from creating an unbounded number of connections. A publisher that establishes K_p connections and needs more connections will reevaluate its processor and bandwidth usage and will create further direct connections using the same heuristic, i.e., assuming the required processor and bandwidth usage to be proportional to the number of connections established.

3.6 Überlay Establishment

We use information obtained through the techniques described above to dynamically form "over-overlays"

termed *überlays* between members of identified entourage so that they can communicate efficiently and with low latency.

3.6.1 Graph Construction

A publisher first constructs a graph data structure with the information received from the DCI protocol. This graph will give the publisher an abstract view of the way its subscribers are connected to the brokers. There are three important differences between the graph constructed by the publisher ($G1$) and a graph constructed by globally observing the way subscribers are actually networked with the brokers ($G2$):

1. $G1$ only shows brokers that distribute the publisher's traffic to two or more child brokers in the broker hierarchy while $G2$ will also show any broker that simply forwards traffic between two other brokers or a broker and a subscriber.
2. $G1$ may have been truncated to show only a number of levels starting from the first broker that distributes the publishers traffic into two children while $G2$ will show all the brokers and subscribers that receive the publishers traffic.
3. $G1$ will only show brokers/subscribers that have at least a ψ interest match with the publisher while $G2$ will show all brokers/subscribers that show interest in some of the publishers messages.

Figs. 4a and 4b show an example graph constructed by a publisher and an actual network of brokers and subscribers that will result in the graph, respectively. The broker $B5$ was not included in the former due to 1. above, and subscribers $S4$ and $S5$ may not have been included either due to 2. or 3. (i.e., because the graph was truncated after three levels, or because $S4$ or $S5$ did not have at least ψ interest match with the publisher p) or simply because $S4$ or $S5$ decided to reject the `COUNT` message from their parent due to one of the reasons given before.

3.6.2 Connection Establishment

Once graphs are established for each remote region a publisher can go ahead and establish überlays. The publisher determines, the number of direct connections it can establish with each remote region r (K_p^r) by dividing K_p among regions proportionally to the sizes (i.e., number of nodes) of respective $G1$ graphs.

For each region r the publisher tries to decide if it should create direct connections with brokers/subscribers in one of the levels of the graph, and if so with which level. The former question is answered based on the existence of a non-empty graph. If the graph is empty, this means that none of the brokers/subscribers had at least ψ interest match with the publisher and hence forming an entourage for distributing messages of p is not viable. To answer the latter question, i.e., the level of a graph with which direct connections should be created, we compare two properties of a graph.

- ad_l – the average distance from a level l , to the interested subscribers. While executing the DCI protocol, we keep track of the average distance from each broker to the subscribers directly below it and the number of subscribers directly below each broker.

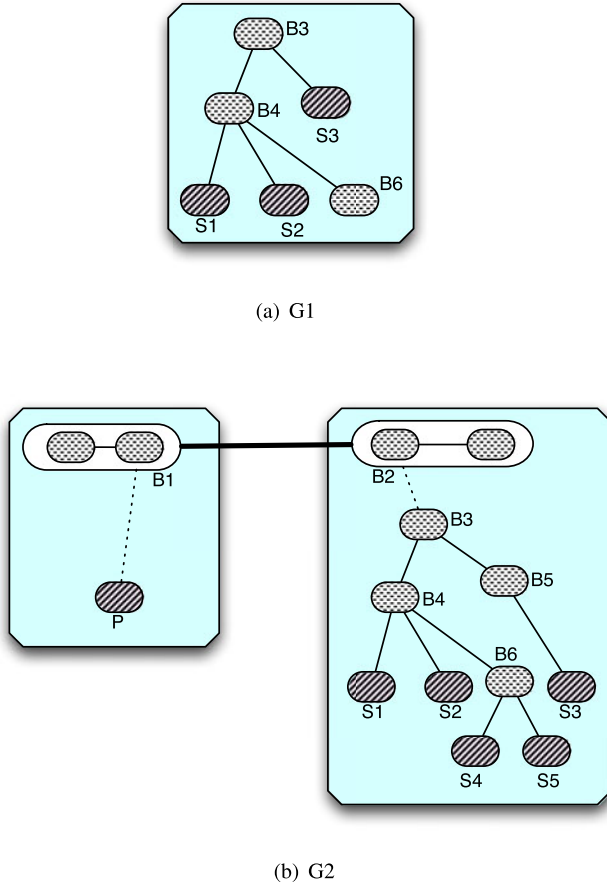


Fig. 4. Graph versus overlay.

To determine the average distance for a level of a graph, we take the full average of the distances from each subscriber to the broker closest to it at the considered level.

cv_l – the portion of the total overlay of the region that will be covered by the selected level. This is the ratio between the number of brokers covered by a given level (brokers in and below the level) and the total number of brokers in the considered region.

By creating direct connections closer to the subscribers, the entourage will be able to deliver messages with low latency. By creating direct connections at higher levels, the direct connections will cover a larger portion of the region's broker network, hence reducing the likelihood of having to recreate the direct connections due to new subscriber joins. This is especially important in the presence of high levels of churn (ch^r for region r). Additionally the publisher can create direct connections which are also bounded by the value of K_p^r for the considered region. The publisher proceeds by selecting the level to which it will establish direct connections (L_p) based on the following heuristic:

$$\frac{cv_{L_p} \times ch^r + 1}{ad_{L_p} + 1} \geq \frac{cv_l \times ch^r + 1}{ad_l + 1} \quad \forall l \in \{1 \dots \lfloor \log K_p^r \rfloor\}.$$

Essentially the heuristic determines the level which gives the best balance between the coverage and the average distance to subscribers. The importance of coverage depends

on the churn of the system. Each factor of the heuristic is incremented by one so that the heuristic gives a non-zero and deterministic value when either churn or distance is zero. To measure the churn, each broker keeps track of the rate at which subscribers join/leave it. This information is aggregated and sent upwards towards the roots where the total churn of the region is determined.

If there are more than K_p^r nodes at the selected level then the publisher will first establish connections with K_p^r randomly selected nodes there. The publisher will keep sending messages through its parent so that the rest of the nodes receive the published messages. Any node that already establishes direct connections with the publisher will discard any message from the publisher received through the node's parent. Once these connections are established the publisher, as mentioned, re-evaluates its resource usage and creates further direct connections as necessary.

If a new subscriber that is interested in messages from the publisher joins the system, initially it will get messages routed via the CPS overlay. The new subscriber will be identified, and a direct connection may be established in the next execution of the DCI protocol. If a node that is directly connected to the publisher needs to discard the connection, it can do so by sending a COUNTRJECT message directly to the publisher. A publisher upon seeing such a message will tear down the direct connection established with the corresponding node.

4 ATMOSPHERE

In this section, we describe Atmosphere, our CPS framework for multi-region deployments which employs the DCI protocol and overlay establishment introduced previously. The core implementation of Atmosphere in Java has approximately 3,200 lines of code.

4.1 Overlay Structure

Atmosphere uses a two-level overlay structure based on *broker* nodes. Every application node that wishes to communicate with other nodes has to initially connect to one of the brokers which will be identified as the node's *parent*. A set of peer brokers form a *broker group*. Each broker in a group is aware of other brokers in that group. Broker-groups are arranged to form *broker-hierarchies*. Broker-hierarchies are illustrated in Fig. 2. As the figure depicts, a broker-hierarchy is established in each considered region. A region can typically represent a LAN, a datacenter, or a zone within a datacenter. At the top (root) level broker-groups of hierarchies are connected to each other. The administrator has to decide on the number of broker-groups to be formed in each region and the placement of broker-groups.

Atmosphere employs subscription summarization to route messages. Each broker summarizes the interests of its subordinates and sends the summaries to its parent broker. Root-level brokers of a broker-hierarchy share their subscription summaries with each other. At initiation, the administrator has to provide each root-level group the identifier of at least one root-level broker from each of the remote regions.

4.2 Fault Tolerance and Scalability

Atmosphere employs common mechanisms for fault tolerance and scalability. Each broker group maintains a strongly consistent membership, so that each broker is aware of the live brokers within its group. A node that needs to connect to a broker-group has to be initially aware of at least one live broker (which will become the node's parent). Once connected, the parent broker provides the node with a list of live brokers within its broker-group and keeps the node updated about membership changes. Each broker, from time to time, sends heartbeat messages to its *children*.

If a node does not receive a heartbeat from its parent for a predefined amount of time, the parent is presumed to have failed, and the node connects to a different broker of the same group according to the last membership update from the failed parent. A node that wishes to leave, sends an *unsubscription* message to its parent broker. The parent removes the node from its records and updates the peer brokers as necessary.

Atmosphere can be scaled both horizontally and vertically. Horizontal scaling is achieved by adding more brokers to groups. Additionally, Atmosphere can be vertically scaled by increasing the number of levels of the broker-hierarchy. Nodes may subscribe to a broker at any level.

4.3 Flexible Communication

Atmosphere implements the DCI protocol of Section 3. To this end, each publisher sends *COUNT* messages to its broker. These messages are propagated up the hierarchy and once the root brokers are reached, are distributed to the remote regions to identify entourages. Once suitable entourages are identified, overlays are established which are used to disseminate messages to interested subscribers with low latency.

When changes in subscriptions (e.g., joining/leaving of subscribers) arrive at brokers these may propagate corresponding notifications upstream even if subscriptions are covered by existing summaries; when arriving at brokers involved in direct connections these can notify publishers directly of changes, prompting them to re-trigger counts.

Fig. 5 shows the major entities of a Atmosphere deployment and corresponding interactions that result in messages being exchanged through a direct connection.

4.4 Advertisements

By wrapping it with the client library of Atmosphere the DCI protocol for publishers/subscribers is transparent to application components, at the exception of *advertisements* which publishers can optionally issue to make effective use of direct connections.

Advertisements are supported in many overlay-based CPS systems, albeit not strictly required. Similarly, publishers in Atmosphere are not forced to issue such advertisements in Atmosphere, although effective direct connection establishment hinges on accurate knowledge of publication ranges. Atmosphere can employ runtime monitoring of published messages if necessary. For such inference, the client library of Atmosphere compares messages published by a given publisher against the currently stored advertisement

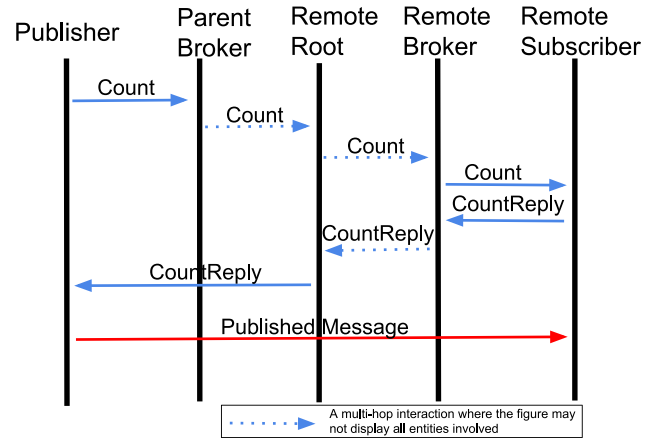


Fig. 5. Entities and interactions.

and adapts the advertisement if required. When witnessing significant changes, the new advertisement is stored and the DCI protocol is re-triggered.

Note that messages beyond the scope of a current advertisement are nonetheless propagated over the direct connections in addition to the overlay. The latter is necessary to deal with joining subscribers in general as mentioned, while the former is done for performance reasons—the directly connected nodes might be interested in the message since the publisher's range of publications announced earlier can be a subset of the ranges covered by any subscriptions.

The obvious downside of obtaining advertisements only by inference is that overlay creation is delayed and thus latency is increased until the ideal connections are established. To avoid constraining publishers indefinitely to previously issued advertisements, the Atmosphere client library offers API calls to issue explicit advertisement updates. Such updates can be viewed as the publisher-side counterpart of *parametric subscriptions* [24] whose native support in a CPS overlay network have been shown to not only have benefits in the presence of changing subscriptions, but also to improve upstream propagation of changes in subscription summaries engendered via unsubscriptions and new subscriptions.

5 EVALUATION

We demonstrate the efficiency and versatility of Atmosphere via both microbenchmarks and real-life applications.

5.1 Setup

We use two datacenters for our experiments, both from Amazon EC2. The datacenters are located in US east coast and US west coast. From each of these datacenters we lease 10 *small* EC2 instances with 1.7 GB of memory and one virtual core and 10 *medium* EC2 instances with 3.7 GB of memory and two virtual cores each.

Our experiments are conducted using three publish/subscribe systems: (1) Atmosphere with DCI protocol disabled, representing a pure CPS system (referred to as *CPS* in the following); (2) Atmosphere with DCI protocol enabled (*Atmosphere*); (3) Apache ActiveMQ topic-based messaging system [1] (*TPS*). ActiveMQ is configured for

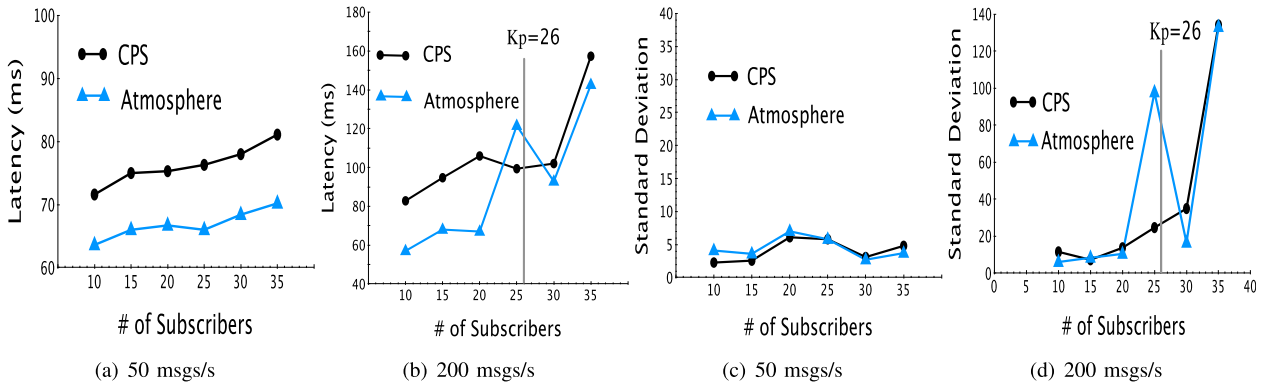


Fig. 6. Latency and standard deviation of latency.

fair comparison to use TCP just like Atmosphere and to not persist messages. All code is implemented in Java.

5.2 Microbenchmarks

We first assess the performance benefits of Atmosphere via micro-benchmarks.

5.2.1 Latency

We conduct experiments to observe the message transmission latency of Atmosphere with and without DCI protocol enabled. The experiment is conducted across two datacenters and uses *small* EC2 instances. A single publisher is deployed in the first datacenter, while between 10 and 35 subscribers are deployed in the second datacenter. Each datacenter maintains three root brokers.

Figs. 6a and 6b show the latency for message rates 50 and 200 msg/s while Figs. 6c and 6d show the standard latency deviations for the same rates. We separate latency from its standard deviation for clarity.

As the graphs clearly show, when the number of interested subscribers is small, maintaining unicast channels between the publisher and the subscribers pays off, even considering that the relatively slow connection to the remote datacenter is always involved, and only local hops are avoided. This helps to dramatically reduce both the average message transmission latency and the variance of latency across subscribers. For message rates 50 and 200, when the number of subscribers is 10, maintaining direct connections decreases the latency by 11 and 31 percent respectively.

For message rates 50 and 200, the value of K_p is determined to be 50 and 26 respectively. The Figs. 6b and 6d show that both the message transmission latency and its variation considerably increases when the publisher reaches this limit. Also the figures show the benefit of not using an overlay after the number of subscribers exceed K_p . For example, as shown in Fig. 6b when publishers move from maintaining an overlay with its entourage to communicating using CPS (25 to 30 subscribers) the average message transmission latency decrease by 24 percent. The increase in latency at 35 subscribers is due to brokers being overloaded, which can be avoided in practical systems by adding more brokers to the overlay and distributing the subscribers among them. Also note that the broker overlay used for this experiment consists of only two levels which is the case where entourage overlays exhibit least benefits.

5.2.2 Number of Subscribers in an Entourage

We conduct experiments using three publisher setups: (1) a publisher uses a *small* EC2 instance (1 core) and sends messages of size 4 KB (p_1); (2) a publisher uses a *medium* EC2 instance (2 cores) and sends messages of size 4 KB (p_2); (3) a publisher uses a *small* EC2 instance and sends messages of size 8 KB (p_3). Subscribers and publishers are placed in two different datacenters as previously. Publishers produce at the *highest possible* rates here. Figs. 7a, 7b, and 7c show how message latency, throughput, and standard latency deviation, respectively, vary for these setups as the number of subscribers changes.

The throughput of p_2 is significantly higher than that of p_1 . This is expected since the rate at which messages can be transmitted increases with the processing power within the relevant confines. Interestingly though, the average message transmission latency for p_2 is higher than the average transmission latency of messages published by p_1 . One hypothesis to explain this result is that the latency depends on the throughput and not directly on the processing power. It may be argued that the variation in latencies could be due to the load characteristics and the placement of the EC2 nodes used, but the possibility for this is low since we used EC2 instances from the same EC2 availability zone as publishers and both experiments used the same set of subscribers. We plan to further investigate this as part of future work.

The size of the transmitted messages has a substantial effect on both throughput and latency. The latter effect becomes significant as the number of subscribers increases. Additionally, Fig. 7c shows that the variation in transmission latency can be significantly decreased by increasing the processing power of the publisher or by decreasing the size of the transmitted messages (e.g., by using techniques such as compression).

5.2.3 Effect of ψ

To study the effects of the clustering factor (ψ) on latency, we deploy a system of one publisher and multiple subscribers. We generated subscribers with interest ranges (of size 20) starting randomly from a fixed set of 200 interests. The publisher publishes a message to one random interest at specific intervals. Brokers are organized into a fully complete binary tree with three levels and 40 subscribers are connected to leaf level brokers. On this setup, latency

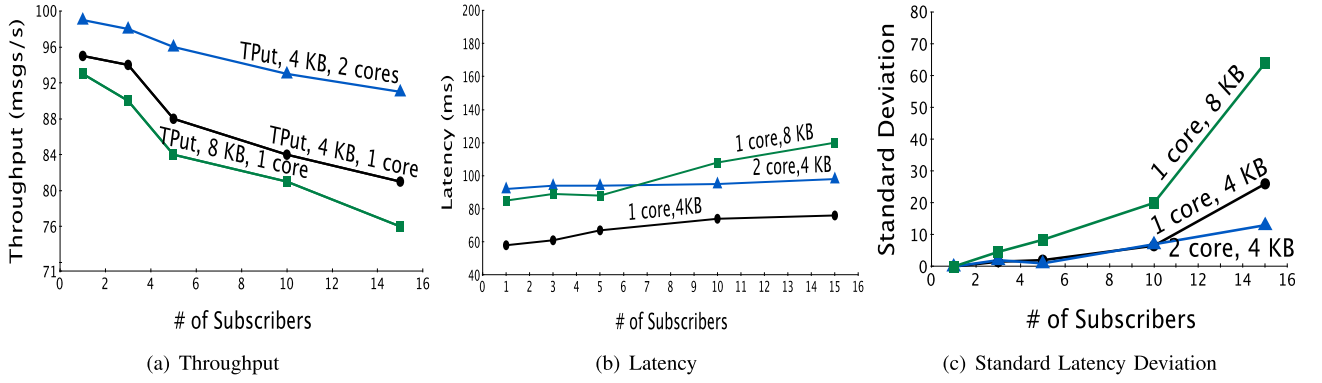


Fig. 7. Effects of resource usage.

measurements are taken with different ψ values. Fig. 8 shows the results. When ψ is high, entourage is not created because no broker has an interest match as high as ψ . This means messages get delivered to root-level brokers which causes higher delays as the messages need to travel through all the levels in the broker network. For a lower value of ψ , an entourage is established, reducing latency.

For communications that require low response times such as chat messaging and real time streaming, ψ should be set to lower values. On the other hand for communication that do not demand low response times such as file sharing or news message exchanging ψ may be set to higher values.

5.3 Case Studies

We developed four applications to show how Atmosphere can be used to make real-world applications efficient.

5.3.1 Social Network

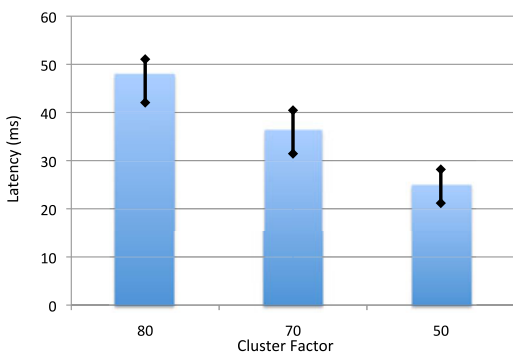
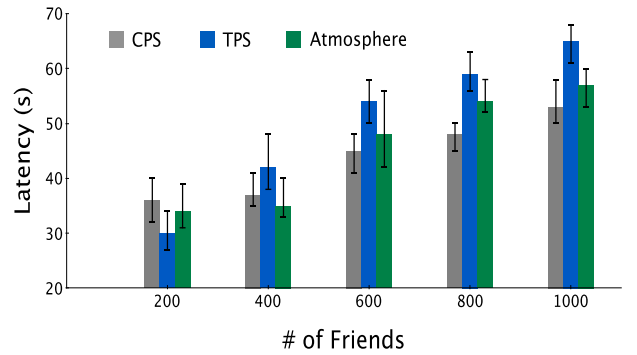
Typical IM clients attached to social networking sites support the following two operations: (1) *status* updates, in which the current status (Busy/Active/Idle or a custom message) of a user is propagated to all users in his/her friend list; (2) the ability to start a conversation with another user in the friend list. Even when explicit status updates are infrequent, IM clients automatically update user status to Idle/Active generating a high number of status updates. We developed an instant messaging service that implements this functionality either on top of Atmosphere or ActiveMQ.

Fig. 9a shows latency measurements for status updates. Figs. 9b shows latency measurements for a randomly

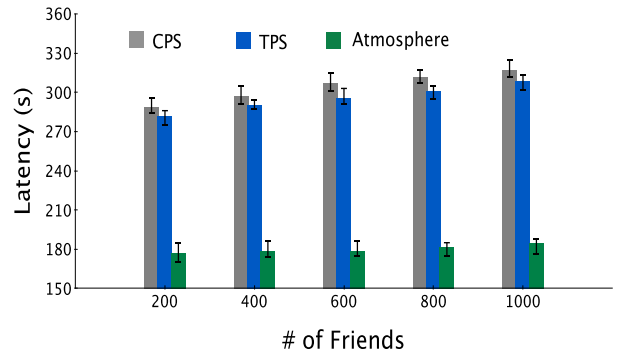
selected friend. For conversations, we use actual conversation logs posted by users of Cleverbot [5]. We evaluate this type of communication on Atmosphere, pure CPS with Atmosphere, and ActiveMQ. The results show that our system is 40 percent faster than pure CPS and 39 percent faster than ActiveMQ in delivering instant messages. For delivering status messages, in the worst case, Atmosphere is on par with both systems because our system distinguishes between the communication types required for status updates and instant message exchange and dynamically forms entourage overlays for delivering instant messages only.

5.3.2 News Service

We developed an Atmosphere-based news feed application that delivers news to subscribed clients. Our news service

Fig. 8. Effect of ψ .

(a) Status



(b) Friend

Fig. 9. Evaluation of our social network App.

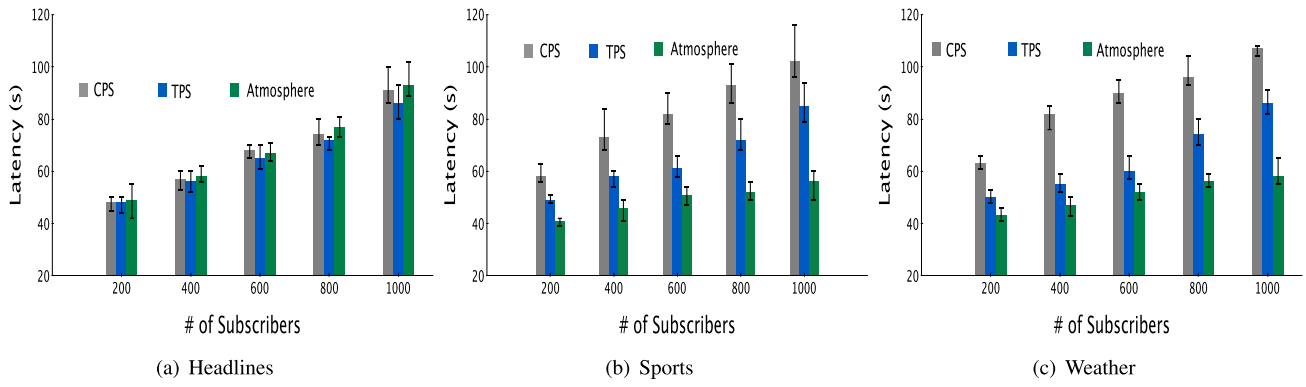


Fig. 10. News App evaluation.

generates two types of messages: (1) messages containing news headlines categorized according to the type of news (e.g., sports, politics, weather); (2) messages containing detailed news items of a given category. This service can also operate on top of either Atmosphere or ActiveMQ.

In Figs. 10a, 10b, and 10c we explore latency of the news application for three different communication patterns. The total number of subscribers varies from 200 to 1,000 with a subset of 30 subscribers interested in sports-based news and a subset of 20 subscribers interested in weather reports. We measure the average latency for delivering sports news and weather reports to these 30 and 20 subscribers. Other subscribers receive all news. Here again our system delivers sports and weather reports 35 percent faster than a pure CPS system and around 25 percent faster than ActiveMQ. This is because Atmosphere automatically creates an entourage for delivering these posts.

5.3.3 Geo-Distributed Storage System

We developed a storage service based on the Apache Hadoop Distributed File System (HDFS) [3]. HDFS is a system for storing large files across multiple hosts in a consistent and highly available manner. HDFS consists of a single master node called the *name node* and multiple slave nodes called the *data nodes*. The name node is only used to determine the data nodes that hold blocks of a file so it incurs no data transfer overhead. The geo-distributed storage system presented here and the geo-distributed lock service presented next were each developed by a graduate student

(with no prior experience in using Atmosphere) in approximately two weeks.

We used Atmosphere to connect HDFS deployments in multiple datacenters, and to make them work as a single geo-distributed storage service (1,200 Java lines of code). Each datacenter maintain its own HDFS instance (with its own name node and data nodes) while still allowing access to files in all the datacenters. The solution is simple because we delegate decisions related to fault-tolerance, concurrent writes and consistency to the underlying HDFS instances.

Figs. 11a and 11b show the latency of the geo-distributed storage system (*Atmosphere*) compared to two naïve HDFS deployments, namely deployed (a) in a single datacenter (*Centralized*), resulting in many remote accesses, and (b) without Atmosphere in multiple datacenters (*Distributed*), resulting in random distribution. Atmosphere-based deployment performs both file create+write and read operations faster than the other two deployments. This is because it always stores data in the HDFS instances closest to the client creating the file, which may not be the case in the other two.

5.3.4 Geo-Distributed Lock Service

We implemented a geo-distributed lock service that can be used to store system configuration information in a consistently replicated manner for fault-tolerance. The service is based on Apache ZooKeeper [4] a system for maintaining distributed configuration and lock services. ZooKeeper guarantees scalability and strong consistency by replicating

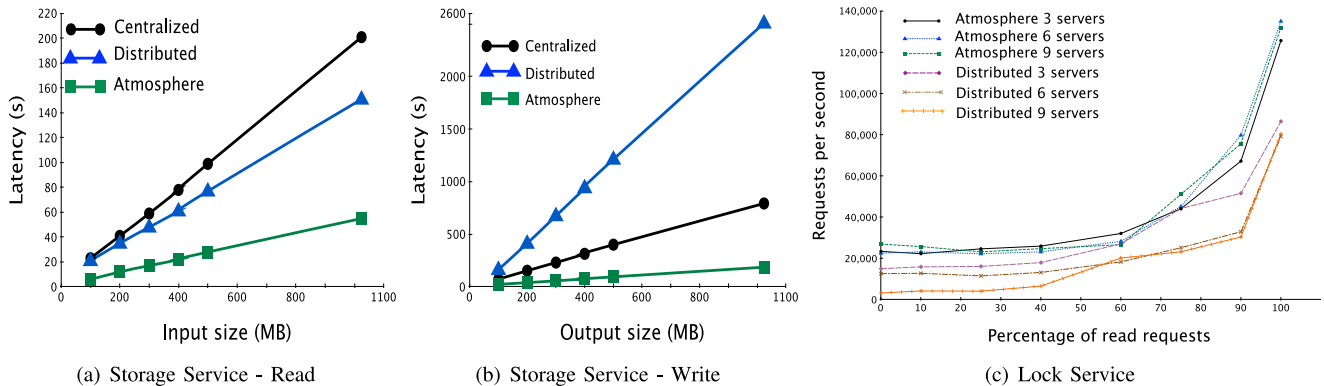


Fig. 11. File system and lock service evaluation.

data across a set of nodes called an *ensemble* and by executing consensus protocols among these nodes.

Hence, similar to the geo-distributed storage service, we maintain a ZooKeeper ensemble per datacenter and interconnect the ensembles (i.e., handle the application requests over the ensembles) using Atmosphere.

We compare the Atmosphere-based lock service with a naïve distributed deployment of ZooKeeper (*Distributed*) where all ZooKeeper nodes participated in a single geo-distributed ensemble. This experiment uses three datacenters. For each run, a constant number of ZooKeeper servers are started at each datacenter. Our system provides the same guarantees as naïve ZooKeeper except the unlikely scenario of complete datacenter failures (in this case *Atmosphere* deployment may lose a part of the stored data).

We vary the percentage of read requests and observe the loads that the systems can handle with 3, 6, and 9 total nodes forming ensembles. Fig. 11c shows the results of the experiment for Atmosphere-based (*Atmosphere*) and a distributed deployment of ZooKeeper where all ZooKeeper nodes participated in a single geo-distributed ensemble (*Distributed*); it shows that by establishing overlays, *Atmosphere* deployment can handle a higher load.

These case studies demonstrate the general applicability of Atmosphere and its underlying protocols.

6 RELATED WORK

This section elaborates on further related work and differences to Atmosphere.

6.1 Other CPS Approaches

Chand and Felber [13] introduce another CPS architecture for distributed systems. The approach is similar to Siena but has support for dynamic subscription and cancellation of interests. This is done by keeping record of the sub-interface through which a filter was received and not just the filter itself. The employed protocol ensures perfect routing which means that a message will be received by all the interested parties and no others. HyperCBR [12] introduces a partition-based approach for content-based routing. Subscriptions and events are propagated in a multidimensional space. The system is partitioned such that every subscription partition interacts with at least one node of every message partition and vice versa. This model can be extended by adding more subscription dimensions. The authors use simulations and analytical studies to show that the model can scale *up* to a large number of nodes.

Many adaptive communication systems have been proposed in the literature. Yoneki and Bacon [38] propose an adaptive publish/subscribe system suitable for mobile ad-hoc networks that represents summarized subscriptions in bloom filters. Rodrigues et al. [32] introduce a mechanism to adapt gossip-based broadcast algorithms according to the amount of resources available to nodes and the global level of congestion. These adaptations focus on scaling *up*.

Among the first efforts to offer content-based addressing as core network service is *data-oriented network architecture* (DONA) [25]. The proposal focuses on encoding of content and content-addressing into low-level network packets. DONA has a broader scope than the present

work, but does not address scaling up and down. *Content-centric networking* (CCN) [30] is a related thrust. Its scope is yet more broadly defined.

6.2 Other Solutions for Cloud Communication

Cloud service providers such as Microsoft and Amazon have introduced *content delivery networks* (CDNs) for communication between their datacenters. Microsoft Azure CDN caches Azure blob content at strategic locations to make them available around the globe. A special variable, TTL, decides the time after which the cached blob content should expire. Therefore the content requests that come only once during a TTL amount of time will not have any performance and cost advantages over a request that is always retrieved remotely. Amazon's CloudFront is a CDN service that can be used to transfer data across Amazon's datacenters. CloudFront can be used to transfer both static and streamed content using a global network of *edge locations*. Objects are organized according to unique domain names and when an end-user requests an object, it is automatically transferred to the nearest edge location. CDNs focus on stored large multimedia data rather than on live communication. In addition, both above-mentioned CDN networks can be used only within their respective service provider boundaries and individual regions.

Volley [7] strategically partitions geo-distributed *stored* data so that the individual data items are placed close to the global "centroid" of the past accesses. For this Volley extracts information about users and data item accesses from log files of the application and needs to be configured with other information such as the number of datacenters, their geographic location and a policy to determine cost of communication between two given datacenters. Volley periodically runs an iterative algorithm to determine the optimum location for data items and places them in a datacenter close to this location depending on space availabilities. The global resource allocation problem for geo-distributed cloud-based applications addressed by solutions such as Volley is orthogonal to the problem of efficient communication within geo-distributed cloud applications which we focus on. The solutions are complementary.

Use of IP Multicast has been restricted in some regions and across the Internet due to difficulties arising with multicast storms or multicast DOS attacks. Dr. Multicast [36] is a protocol that can be used to mitigate these issues. The idea is to introduce a new logical group addressing layer on top of IP Multicast so that access to physical multicast groups and data rates can be controlled with an acceptable user policy. This way system administrators can place caps on the amount of data exchanged in groups and the members that can participate on a group. Dr. Multicast specializes on *intra*-datacenter communication and does not consider *inter*-datacenter communication.

7 CONCLUSIONS

Developing and composing applications executing in the cloud-of-clouds requires generic communication mechanisms. Existing CPS frameworks—though providing generic communication abstractions—do not operate efficiently across communication patterns and regions,

exhibiting large performance gaps to more specific solutions. In contrast, existing simpler TPS solutions cover fewer communication patterns but more effectively—in particular scenarios with few publishers and many subscribers which are wide-spread in cloud-based computing.

We introduced the DCI protocol, a mechanism that can be used to adapt existing solutions to efficiently support different communication patterns; presented its implementation in Atmosphere, a scalable and fault-tolerant CPS framework suitable for multi-region-based deployments such as cross-cloud scenarios. We illustrated the benefits of our approach through different experiments evaluating multi-region deployments of Atmosphere.

We are currently working on complementary techniques that will further broaden the range of efficiently supported communication patterns, for example the migration of subscribers between brokers guided by resource usage on these brokers. Additionally we are exploring the use of Atmosphere as the communication backbone for other systems including our [21] framework for efficiently executing Pig/PigLatin workflows in geo-distributed cloud setups and our G-MR [23] system for efficiently executing sequences of MapReduce jobs on geo-distributed data sets.

ACKNOWLEDGMENTS

The authors are very grateful to Amazon and to Larry Peterson for making it possible to evaluate our research in EC2 and VICCI respectively. This work is financially supported by DARPA Grant N11AP20014 “Large-Scale Cloud-based Data Analysis”, Purdue Research Foundation Grant 204533 “Seamless Cloud Computing”, Google Research award “Geo-Distributed Big Data Processing”, and Cisco Research award “A Fog Architecture”. This paper extends our prior work published at Middleware 2013 [22].

REFERENCES

- [1] Active MQ, <http://activemq.apache.org/>
- [2] Amazon Simple Notification Service, <http://aws.amazon.com/sns/>
- [3] Apache HDFS, <http://hadoop.apache.org>
- [4] Apache ZooKeeper, <http://hadoop.apache.org/zookeeper>
- [5] Cleverbot, <http://zookeeper.apache.org>
- [6] Websphere MQ, <http://www-01.ibm.com/software/integration/wmq/>
- [7] S. Agarwal, J. Dunagan, N. Jain, S. Saroiu, A. Wolman, and H. Bhogan, “Volley: Automated data placement for Geo-distributed cloud services,” in *Proc. 7th USENIX Conf. Netw. Syst. Des. Implementation*, 2010, pp. 17–32.
- [8] M. K. Aguilera, R. E. Strom, D. C. Sturman, M. Astley, and T. D. Chandra, “Matching events in a content-based subscription system,” in *Proc. 18th ACM Symp. Principles Distrib. Comput.*, 1999, pp. 53–62.
- [9] A. Bessani, M. Correia, B. Quaresma, F. André, and P. Sousa, “DepSky: Dependable and secure storage in a cloud-of-clouds,” in *Proc. 6th Eur. Conf. Comput. Syst.*, 2011, pp. 31–46.
- [10] A. Carzaniga, D. Rosenblum, and A. Wolf, “Design and evaluation of a wide-area event notification service,” *ACM Trans. Comput. Syst.*, vol. 19, no. 3, pp. 332–383, 2001.
- [11] A. Carzaniga, D. S. Rosenblum, and A. L. Wolf, “Achieving scalability and expressiveness in an internet-scale event notification service,” in *Proc. 19th ACM Symp. Principles Distrib. Comput.*, 2000, pp. 219–227.
- [12] S. Castelli, P. Costa, and G. P. Picco, “HyperCBR: Large-scale content-based routing in a multidimensional space,” in *Proc. 27th IEEE Int. Conf. Comput. Commun.*, 2008, pp. 1714–1722.
- [13] R. Chand and P. A. Felber, “A scalable protocol for content-based routing in overlay networks,” in *Proc. 2nd IEEE Int. Symp. Netw. Comput. Appl.*, 2003, pp. 123–130.
- [14] S. Das, D. Agrawal, and A. El Abbadi, “G-Store: A scalable data store for transactional multi key access in the cloud,” in *Proc. 1st ACM Symp. Cloud Comput.*, 2010, pp. 163–174.
- [15] J. Dean and S. Ghemawat, “MapReduce: Simplified data processing on large clusters,” *Commun. ACM*, vol. 51, no. 1, pp. 107–113, 2008.
- [16] G. DeCandia, D. Hastorun, M. Jampani, G. Kakulapati, A. Lakshman, A. Pilchin, S. Sivasubramanian, P. Voshall, and W. Vogels, “Dynamo: Amazon’s highly available key-value store,” in *Proc. 21st ACM Symp. Operating Syst. Principles*, 2007, pp. 205–220.
- [17] S. Deering and D. Cheriton, “Multicast routing in datagram internetworks and extended LANs,” *ACM Trans. Comput. Syst.*, vol. 8, no. 2, pp. 85–110, May 1990.
- [18] L. Fiege, F. C. Gärtner, O. Kasten, and A. Zeidler, “Supporting mobility in content-based publish/subscribe middleware,” in *Proc. 4th ACM/IFIP/USENIX Int. Middleware Conf.*, 2003, pp. 103–122.
- [19] S. G. Grivas, K. T. Uttam, and H. Wache, “Cloud broker: Bringing intelligence into the cloud,” in *Proc. 3rd IEEE Int. Conf. Cloud Comput.*, 2010, pp. 544–545.
- [20] (2012.) Apache Zookeeper: Hedwig. [Online]. Available: <https://cwiki.apache.org/ZOOKEEPER/hedwig.html>
- [21] C. Jayalath and P. Eugster, “Efficient Geo-distributed data processing with rout,” in *Proc. 33rd IEEE Int. Conf. Distrib. Comput. Syst.*, 2013, pp. 470–480.
- [22] C. Jayalath, J. Stephen, and P. Eugster, “Atmosphere: A universal cross-cloud communication infrastructure,” in *Proc. 14th ACM/IFIP/USENIX Int. Middleware Conf.*, 2013, pp. 163–182.
- [23] C. Jayalath, J. Stephen, and P. Eugster, “From the cloud to the atmosphere: Running mapreduce across datacenters,” *IEEE Trans. Comput.*, vol. 63, no. 1, pp. 74–87, Jan. 2014.
- [24] K. R. Jayaram, C. Jayalath, and P. Eugster, “Parametric subscriptions for content-based publish/subscribe networks,” in *Proc. 11th ACM/IFIP/USENIX Int. Middleware Conf.*, 2010, pp. 128–147.
- [25] T. Koponen, M. Chawla, B.-G. Chun, A. Ermolinskiy, K. H. Kim, S. Shenker, and I. Stoica, “A data-oriented (and beyond) network architecture,” in *Proc. ACM Conf. Appl., Technol. Arch., Protocols Comput. Commun.*, 2007, pp. 181–192.
- [26] J. Kreps, N. Narkhede, and J. Rao, “Kafka: A distributed messaging system for log processing,” in *Proc. 6th Int. Workshop Netw. Meets Databases*, <http://research.microsoft.com/en-us/um/people/srikanth/netdb11/netdb11papers/netdb11-final12.pdf>, 2011.
- [27] M. Li, F. Ye, M. Kim, H. Chen, and H. Lei, “A scalable and elastic publish/subscribe service,” in *Proc. IEEE Int. Parallel Distrib. Process. Symp.*, 2011, pp. 1254–1265.
- [28] A. Majumder, N. Shrivastava, R. Rastogi, and A. Srinivasan, “Scalable content-based routing in pub/sub systems,” in *Proc. 28th IEEE Int. Conf. Comput. Commun.*, 2009, pp. 567–575.
- [29] P. Triantafyllou and A. A. Economides, “Subscription summarization: A new paradigm for efficient publish/subscribe systems,” in *Proc. 24th Int. Conf. Distrib. Comput. Syst.*, 2004, pp. 562–571.
- [30] Palo Alto Research Center, Project CCNx. [Online]. Available: <http://www.ccnx.org>
- [31] P. Pietzuch and J. Bacon, “Hermes: A distributed event-based middleware architecture,” in *Proc. 22nd Int. Conf. Distrib. Comput. Syst., Workshops*, 2002, pp. 611–618.
- [32] L. Rodrigues, S. B. Handurukande, J. Pereria, R. Guerraoui, and A.-M. Kermmarrec, “Adaptive gossip-based broadcast,” in *Proc. IEEE/IFIP Int. Conf. Dependable Syst. Netw.*, 2003, pp. 47–56.
- [33] G. Li, S. Hou, and H.-A. Jacobsen, “A unified approach to routing, covering and merging in publish/subscribe systems based on modified binary decision diagrams,” in *Proc. 25th IEEE Int. Conf. Distrib. Comput. Syst.*, 2005, pp. 447–457.
- [34] R. Sherafat and H.-A. Jacobsen, “Publi+: A peer-assisted publish/subscribe service for timely dissemination of bulk content,” in *Proc. 32nd IEEE Int. Conf. Distrib. Comput. Syst.*, 2012, pp. 345–354.
- [35] M. A. Tariq, B. Koldehofe, G. G. Koch, and K. Rothermel, “Distributed spectral cluster management: A method for building dynamic publish/subscribe systems,” in *Proc. 6th ACM Int. Conf. Distrib. Event-Based Syst.*, 2012, pp. 213–224.

- [36] Y. Vigfusson, H. Abu-Libdeh, M. Balakrishnan, K. Birman, R. Burgess, G. Chockler, H. Li Haoyuan, and Y. Tock, "Dr. Multicast: Rx for data center communication scalability," in *Proc. 5th Eur. Conf. Comput. Syst.*, 2010, pp. 349–362.
- [37] S. Voulgaris, E. Riviere, A.-M. Kermarrec, and M. van Steen, "Sub-2-Sub: Self-organizing content-based publish subscribe for dynamic large scale collaborative networks," in *Proc. 5th Int. Workshop Peer-To-Peer Syst.*, <http://iptps06.cs.ucsb.edu/papers/Voulgaris-sub06.pdf>, 2006.
- [38] E. Yoneki and J. Bacon, "An adaptive approach to content-based subscription in mobile ad hoc networks," in *2nd IEEE Conf. Pervasive Comput. Commun. Workshops*, 2004, pp. 92–97.



Chamikara Jayalath received the BSc degree from the University of Morutawa, Sri Lanka. He is working toward the PhD degree in computer science at Purdue University. His research interests include distributed systems and especially cloud computing. He is a committer and Project Management Committee member for the Apache Web services project.



Julian Stephen received the BTech degree from Mahatma Gandhi University Kottayam, India, in 2004. He is working towards the PhD degree in computer science at Purdue University. His research interests include cloud computing, most notably security in cloud computing.



Patrick Eugster received the MS and PhD degrees from EPFL. He is an associate professor in computer science at Purdue University. His research interests include distributed systems and programming. He is a recipient of a US National Science Foundation CAREER award in 2007 and a member of US Defense Advanced Research Projects Agency 2011 computer science study group.

▷ For more information on this or any other computing topic, please visit our Digital Library at www.computer.org/publications/dlib.