

SCRIBE: A large-scale and decentralized publish-subscribe infrastructure

Miguel Castro¹, Peter Druschel², Anne-Marie Kermarrec¹,
and Antony Rowstron¹

¹Microsoft Research Ltd., 7 J J Thomson Avenue,
Cambridge, CB3 0FB, UK.

²Rice University MS-132, 6100 Main Street,
Houston, TX 77005, USA.

PRELIMINARY DRAFT SUBMITTED FOR PUBLICATION

Abstract

This paper presents Scribe, a large-scale event notification infrastructure for topic-based publish-subscribe applications. Scribe supports large numbers of topics, with a potentially large number of subscribers per topic. Scribe is built on top of Pastry, a generic peer-to-peer object location and routing substrate overlaid on the Internet, and leverages Pastry's reliability, self-organization, and locality properties. Pastry is used to create a topic (group) and to build an efficient multicast tree for the dissemination of events within a topic with a potentially large set of subscribers (members). Scribe provides weak reliability guarantees, but we outline how an application can extend Scribe to provide stronger ones. Simulation results, based on a realistic network topology model, show that Scribe scales across a wide range of topics and subscribers per topic. Also, it balances the load on the nodes while achieving acceptable delay and link stress when compared to IP multicast.

Index terms: publish-subscribe, group communication, application-level multicast, peer-to-peer.

1 Introduction

Publish-subscribe has emerged as a promising paradigm for large-scale, Internet based distributed systems. In general, subscribers register their interest in a topic or a pattern of events and then asynchronously receive events matching their interest, regardless of the events' publisher. Topic-based publish-subscribe [1, 2, 3] is strongly similar to group-based communication; subscribing is equivalent

to becoming a member of a group. For such systems the challenge remains to build an infrastructure that can scale to, and tolerate the failure modes of the general Internet.

Techniques such as SRM (Scalable Reliable Multicast Protocol) [4] or RMTP (Reliable Message Transport Protocol) [5] have added reliability to network-level IP multicast [6, 7] solutions. However, tracking membership remains an issue in router-based multicast approaches and the lack of wide deployment of IP multicast limits their applicability. As a result, application-level multicast is gaining popularity. Appropriate algorithms and systems for scalable subscription management and scalable, reliable propagation of events are still an active research area [8, 9, 10, 11].

Recent work on peer-to-peer overlay networks offers a scalable, self-organizing, fault-tolerant substrate for decentralized distributed applications [12, 13, 14, 15]. Such systems offer an attractive platform for publish-subscribe systems that can leverage these properties. In this paper we present Scribe, a large-scale, decentralized event notification infrastructure built upon Pastry, a scalable, self-organizing peer-to-peer location and routing substrate with good locality properties [12]. Scribe provides efficient application-level multicast and is capable of scaling to a large number of subscribers, publishers and topics.

Scribe and Pastry adopt a fully decentralized peer-to-peer model, where each participating node has equal responsibilities. Scribe builds a multicast tree, formed by joining the Pastry routes from each subscriber to a rendez-vous point associated with a topic. Subscription maintenance and publishing in Scribe leverages the robustness, self-organization, locality and reliability properties of Pastry. Section 2 gives an overview of the Pastry routing and object location infrastructure. Section 3 describes the basic design of Scribe. We present performance results in Section 4 and discuss related work in Section 5.

2 Pastry

In this section we briefly sketch Pastry [12], a peer-to-peer location and routing substrate upon which Scribe was built. Pastry forms a secure, robust, self-organizing overlay network in the Internet. Any Internet-connected host that runs the Pastry software and has proper credentials can participate in the overlay network.

Each Pastry node has a unique, 128-bit nodeId. The set of existing nodeIds is uniformly dis-

tributed; this can be achieved, for instance, by basing the `nodeId` on a secure hash of the node’s public key or IP address. Given a message and a key, Pastry reliably routes the message to the Pastry node with a `nodeId` that is numerically closest to the key, among all live Pastry nodes. Assuming a Pastry network consisting of N nodes, Pastry can route to any node in less than $\lceil \log_{2^b} N \rceil$ steps on average (b is a configuration parameter with typical value 4). With concurrent node failures, eventual delivery is guaranteed unless $\lfloor l/2 \rfloor$ nodes with *adjacent* `nodeIds` fail simultaneously (l is a configuration parameter with typical value 16).

The tables required in each Pastry node have only $(2^b - 1) * \lceil \log_{2^b} N \rceil + 2l$ entries, where each entry maps a `nodeId` to the associated node’s IP address. Moreover, after a node failure or the arrival of a new node, the invariants in all affected routing tables can be restored by exchanging $O(\log_{2^b} N)$ messages. In the following, we briefly sketch the Pastry routing scheme. A full description and evaluation of Pastry can be found in [12].

For the purposes of routing, `nodeIds` and keys are thought of as a sequence of digits with base 2^b . A node’s routing table is organized into $\lceil \log_{2^b} N \rceil$ rows with $2^b - 1$ entries each. The $2^b - 1$ entries in row n of the routing table each refer to a node whose `nodeId` matches the present node’s `nodeId` in the first n digits, but whose $n + 1$ th digit has one of the $2^b - 1$ possible values other than the $n + 1$ th digit in the present node’s id. The uniform distribution of `nodeIds` ensures an even population of the `nodeId` space; thus, only $\lceil \log_{2^b} N \rceil$ levels are populated in the routing table. Each entry in the routing table refers to one of potentially many nodes whose `nodeId` have the appropriate prefix. Among such nodes, the one closest to the present node (according to a scalar proximity metric, such as the delay or the number of IP routing hops) is chosen in practice.

In addition to the routing table, each node maintains IP addresses for the nodes in its *leaf set*, i.e., the set of nodes with the $l/2$ numerically closest larger `nodeIds`, and the $l/2$ nodes with numerically closest smaller `nodeIds`, relative to the present node’s `nodeId`. Figure 1 depicts the state of a hypothetical Pastry node with the `nodeId` 10233102 (base 4), in a system that uses 16 bit `nodeIds` and a value of $b = 2$.

In each routing step, a node normally forwards the message to a node whose `nodeId` shares with the key a prefix that is at least one digit (or b bits) longer than the prefix that the key shares with the present node’s id. If no such node is found in the routing table, the message is forwarded to a node whose `nodeId` shares a prefix with the key as long as the current node, but is numerically closer to the key than the present node’s id. Such a node must be in the leaf set unless the message has already

NodeId 10233102			
Leaf set		SMALLER	LARGER
10233033	10233021	10233120	10233122
10233001	10233000	10233230	10233232
Routing table			
-0-2212102	1	-2-2301203	-3-1203203
0	1-1-301233	1-2-230203	1-3-021022
10-0-31203	10-1-32102	2	10-3-23302
102-0-0230	102-1-1302	102-2-2302	3
1023-0-322	1023-1-000	1023-2-121	3
10233-0-01	1	10233-2-32	
0		102331-2-0	
		2	
Neighborhood set			
13021022	10200230	11301233	31301233
02212102	22301203	31203203	33213321

Figure 1: State of a hypothetical Pastry node with nodeId 10233102, $b = 2$. All numbers are in base 4. The top row of the routing table represents level zero. The neighborhood set is not used in routing, but is needed during node addition/recovery.

arrived at the node with numerically closest nodeId. And, unless $\lfloor |L|/2 \rfloor$ adjacent nodes in the leaf set have failed simultaneously, at least one of those nodes must be live.

2.1 Locality

Next, we discuss Pastry’s locality properties, i.e., the properties of Pastry’s routes with respect to the proximity metric. The proximity metric is a scalar value that reflects the “distance” between any pair of nodes, such as the number of IP routing hops, geographic distance, delay, or a combination thereof. It is assumed that a function exists that allows each Pastry node to determine the “distance” between itself and a node with a given IP address.

We limit our discussion to two of Pastry’s locality properties that are relevant to Scribe. The first property is the total distance, in terms of the proximity metric, that messages are traveling along Pastry routes. Recall that each entry in the node routing tables is chosen to refer to the nearest node, according to the proximity metric, with the appropriate nodeId prefix. As a result, in each step a message is routed to the nearest node with a longer prefix match. Simulations show that, given our network topology model, the average distance traveled by a message is less than 66% higher than the

distance between the source and destination in the underlying Internet [12].

Let us assume that two nodes within distance d from each other route messages with the same key, such that the distance from each node to the node with `nodeId` closest to the key is much larger than d . The second locality property is concerned with the “distance” the messages travel until they reach a node where their routes merge. Simulations show that, given our network topology model, the average distance traveled by each of the two messages before their routes merge is approximately equal to the distance between their respective source nodes. These properties have a strong impact on the locality properties of the Scribe multicast trees, as explained in Section 3.

2.2 Node addition and failure

A key design issue in Pastry is how to efficiently and dynamically maintain the node state, i.e., the routing table, leaf set and neighborhood sets, in the presence of node failures, node recoveries, and new node arrivals. The protocol is described and evaluated in [12].

Briefly, an arriving node with the newly chosen `nodeId` X can initialize its state by contacting a nearby node A (according to the proximity metric) and asking A to route a special message using X as the key. This message is routed to the existing node Z with `nodeId` numerically closest to X ¹. X then obtains the leaf set from Z , the neighborhood set from A , and the i th row of the routing table from the i th node encountered along the route from A to Z . One can show that using this information, X can correctly initialize its state and notify nodes that need to know of its arrival, thereby restoring all of Pastry’s invariants.

To handle node failures, neighboring nodes in the `nodeId` space (which are aware of each other by virtue of being in each other’s leaf set) periodically exchange keep-alive messages. If a node is unresponsive for a period T , it is presumed failed. All members of the failed node’s leaf set are then notified and they update their leaf sets to restore the invariant. Since the leaf sets of nodes with adjacent `nodeIds` overlap, this update is trivial. A recovering node contacts the nodes in its last known leaf set, obtains their current leaf sets, updates its own leaf set and then notifies the members of its new leaf set of its presence. Routing table entries that refer to failed nodes are repaired lazily; the details are described in [12].

¹In the exceedingly unlikely event that X and Z are equal, the new node must obtain a new `nodeId`.

2.3 Pastry API

In this section, we briefly describe the application programming interface (API) exported by Pastry to applications such as Scribe. The presented API is slightly simplified for clarity. Pastry exports the following operations:

nodeId = pastryInit(Credentials) causes the local node to join an existing Pastry network (or start a new one) and initialize all relevant state; returns the local node's nodeId. The credentials are provided by the application and contain information needed to authenticate the local node and to securely join the Pastry network. A full discussion of Pastry's security model is beyond the scope of this paper.

route(msg,key) causes Pastry to route the given message to the node with nodeId numerically closest to key, among all live Pastry nodes.

send(msg,IP-addr) causes Pastry to send the given message to the node with the specified IP address, if that node is live. The message is received by that node through the deliver method.

Applications layered on top of Pastry must export the following operations:

deliver(msg,key) called by Pastry when a message is received and the local node's nodeId is numerically closest to key among all live nodes, or when a message is received that was transmitted via *send*, using the IP address of the local node.

forward(msg,key,nextId) called by Pastry just before a message is forwarded to the node with nodeId = nextId. The application may change the contents of the message or the value of nextId. Setting the nextId to NULL will terminate the message at the local node.

newLeafs(leafSet) called by Pastry whenever there is a change in the leaf set. This provides the application with an opportunity to adjust application-specific invariants based on the leaf set.

In the following section, we will describe how Scribe is layered on top of the Pastry API. Other applications built on top of Pastry include PAST, a persistent, global storage utility [16, 17].

3 Scribe

Scribe is a scalable event notification infrastructure built on top of Pastry. Any Scribe node may create a *topic*; other nodes can then register their interest in the topic and become a *subscriber* to the topic. Any Scribe node with the appropriate credentials for the topic can then publish events, and Scribe disseminates these events to all the topic's subscribers. Scribe provides a best-effort dissemination of events, and specifies no particular event delivery order. However, stronger reliability guarantees and ordered delivery for a topic can be built on top of Scribe, as outlined in Section 3.2. Nodes can create, publish events and subscribe to many topics, and topics can have many publishers and subscribers. Scribe can support large numbers of topics with a wide range of subscribers per topic, and a high rate of subscriber turnover.

Scribe offers a simple API to its applications:

create(credentials, topicId) creates a topic with topicId. Throughout, the credentials are used for access control.

subscribe(credentials, topicId, eventHandler) causes the local node to subscribe to the topic with topicId. All subsequently received events for that topic are passed to the specified event handler.

unsubscribe(credentials, topicId) causes the local node to unsubscribe from the topic with topicId.

publish(credentials, topicId, event) causes the event to be published in the topic with topicId.

Scribe uses Pastry to manage topic creation, subscription, and to build a per-topic multicast tree used to disseminate the events published in the topic. Pastry and Scribe are fully decentralized, all decisions are based on local information, and each node has identical capabilities. Each node can act as a publisher, a root of a multicast tree, a subscriber to a topic, a node within a multicast tree, and any sensible combination of the above. Much of the scalability and reliability of Scribe and Pastry derives from this peer-to-peer model.

3.1 Scribe Implementation

A Scribe system consists of a network of Pastry nodes, where each node runs the Scribe application software. The Scribe software on each node provides the *forward* and *deliver* methods, which are

```

(1) forward(msg, key, nextId)
(2)   switch msg.type is
(3)     SUBSCRIBE : if !(msg.topic ∈ topics)
(4)                 topics = topics ∪ msg.topic
(5)                 route(msg,msg.topic)
(6)                 topics[msg.topic].children ∪ msg.source
(7)                 nextId = null // Stop routing the original message

```

Figure 2: Scribe implementation of forward.

```

(1) deliver(msg,key)
(2)   switch msg.type is
(3)     CREATE :      topics = topics ∪ msg.topic
(4)     SUBSCRIBE :  topics[msg.topic].children ∪ msg.source
(5)     PUBLISH :    ∇ node in topics[msg.topic].children
(6)                 send(msg,node)
(7)                 if subscribedTo(msg.topic)
(8)                 invokeEventHandler(msg.topic, msg)
(9)     UNSUBSCRIBE : topics[msg.topic].children =
(10)                  topics[msg.topic].children - msg.topic
(11)                 if (|topics[msg.topic].children| = 0)
(11)                  send(msg,topics[msg.topic].parent)

```

Figure 3: Scribe implementation of deliver.

invoked by Pastry whenever a Scribe message arrives. The pseudo-code for these Scribe methods, simplified for clarity, is shown in Figure 2 and Figure 3, respectively.

Recall that the forward method is called whenever a Scribe message is routed through a node. The deliver method is called when a Scribe message arrives at the node with nodeId numerically closest to the message’s key, or when a message was addressed to the local node using the Pastry *send* operation. The possible message types in Scribe are SUBSCRIBE, CREATE, UNSUBSCRIBE and PUBLISH; the roles of these messages are described in the next sections.

The following variables are used in the pseudocode: *topics* is the set of topics that the local node is aware of, *msg.source* is the nodeId of the message’s source node, *msg.event* is the published event (if present), *msg.topic* is the topicId of the topic and *msg.type* is the message type.

3.1.1 Topic Management

Each topic has a unique *topicId*. The Scribe node with a *nodeId* numerically closest to the *topicId* acts as the *rendez-vous point* for the associated topic. The rendez-vous point forms the root of a multicast tree created for the topic.

To create a topic, a Scribe node asks Pastry to route a CREATE message using the *topicId* as the key (e.g. `route(CREATE,topicId)`). Pastry delivers this message to the node with the *nodeId* numerically closest to *topicId*. The Scribe deliver method adds the topic to the list of topics it already knows about (line 3 of Figure 3). It also checks the credentials to ensure that the topic can be created, and stores the credentials in the topics set. This Scribe node becomes the rendez-vous point for the topic.

The *topicId* is the hash of the topic's textual name concatenated with its creator's name. The hash is computed using a collision resistant hash function (e.g. SHA-1 [18]), which ensures a uniform distribution of *topicIds*. Since Pastry *nodeIds* are also uniformly distributed, this ensures an even distribution of topics across Pastry nodes. A *topicId* can be generated by any Scribe node using only the textual name of the topic and its creator, without the need for an additional naming service. Of course, proper credentials are necessary to subscribe or publish in the associated topic.

3.1.2 Membership management

Scribe creates a multicast tree, rooted at the rendez-vous point, to disseminate the events published in the topic. The multicast tree is created using a scheme similar to reverse path forwarding [19]. The tree is formed by joining the Pastry routes from each subscriber to the rendez-vous point. Subscriptions to a topic are managed in a decentralized manner to support large and dynamic sets of subscribers.

Scribe nodes that are part of a topic's multicast tree are called *forwarders* with respect to the topic; they may or may not be subscribers to the topic. Each forwarder maintains a *children table* for the topic containing an entry (IP address and *NodeId*) for each of its children in the multicast tree.

When a Scribe node wishes to subscribe to a topic, it asks Pastry to route a SUBSCRIBE message with the topic's *topicId* as the key (e.g. `route(SUBSCRIBE,topicId)`). This message is routed by Pastry towards the topic's rendez-vous point. At each node along the route, Pastry invokes Scribe's forward method. Forward (lines 3 to 7 in Figure 2) checks its list of topics to see if it is currently a forwarder; if so, it accepts the node as a child, adding it to the children table. If the node is not already a forwarder, it creates an entry for the topic, and adds the source node as a child in the associated children table. It

then becomes a forwarder for the topic by sending a SUBSCRIBE message to the next node along the route from the original subscriber to the rendez-vous point. The original message from the source is terminated; this is achieved by setting `nextId = null`, in line 7 of Figure 2.

Figure 4 illustrates the subscription mechanism. The circles represent nodes, and some of the nodes have their *nodeId* shown. For simplicity $b = 1$, so the prefix is matched one bit at a time. We assume that there is a topic with *topicId* 1100 whose rendez-vous point is the node with the same identifier. The node with *nodeId* 0111 is subscribing to this topic. In this example, Pastry routes the SUBSCRIBE message to node 1001; then the message from 1001 is routed to 1101; finally, the message from 1101 arrives at 1100. This route is indicated by the solid arrows in Figure 4.

Let us assume that nodes 1001 and 1101 are not already forwarders for topic 1100. The subscription of node 0111 causes the other two nodes along the route to become forwarders for the topic, and causes them to add the preceding node in the route to their children tables. Now let us assume that node 0100 decides to subscribe to the same topic. The route that its SUBSCRIBE message would take is shown using dot-dash arrows. However, since node 1001 is already a forwarder, it adds node 0100 to its children table for the topic, and the SUBSCRIBE message is terminated.

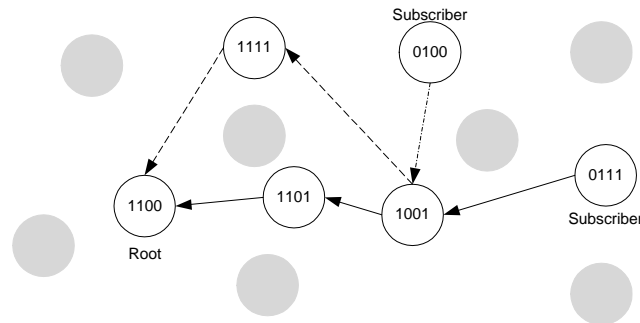


Figure 4: Base Mechanism for Subscription and Multicast Tree Creation.

When a Scribe node wishes to unsubscribe from a topic, a node locally marks the topic as no longer required. If there are no other entries in the children table, it sends a UNSUBSCRIPTION message to its parent in the multicast tree, as shown in lines 9 to 11 in Figure 3. The message proceeds recursively up the multicast tree, until a node is reached that still has entries in the children table after removing the departing child. It should be noted that nodes in the multicast tree are aware of their parent's *nodeId* only after they have received an event from their parent. Should a node wish to unsubscribe before receiving an event, the implementation transparently delays the unsubscribe until the first event is

received.

The subscriber management mechanism is efficient for topics with different numbers of subscribers, varying from one to all Scribe nodes. The list of subscribers to a topic is distributed across the nodes in the multicast tree. Pastry's randomization properties ensure that the tree is well balanced and that the forwarding load is evenly balanced across the nodes. This balance enables Scribe to support large numbers of topics and subscribers per topics. Subscription requests are handled locally in a decentralized fashion. In particular, the rendez-vous point does not handle all subscription requests.

The locality properties of Pastry ensure that for most nodes subscribing to a topic, (1) the node is further from the rendez-vous point, according to Pastry's proximity metric, than its parent in the multicast tree; and (2), the node's parent in the multicast tree is close to the node, according to the proximity metric. These properties follow from the locality properties of Pastry, discussed in Section 2.1. Simulation results quantifying the locality properties of the Scribe multicast tree will be presented in Section 4.

3.1.3 Event dissemination

Publishers use Pastry to locate the rendez-vous point of a topic. If the publisher is aware of the rendez-vous point's IP address then the PUBLISH message can be sent straight to the node. If the publisher does not know the IP address of the rendez-vous point, then it uses Pastry to route to that node (e.g. `route(PUBLISH, topicId)`), and asks the rendez-vous point to return its IP address to the publisher. Events are disseminated from the rendez-vous point along the multicast tree in the obvious way (lines 5 and 6 of Figure 3).

The caching of the rendez-vous point's IP address is an optimization, to avoid repeated routing through Pastry. If the rendez-vous point fails then the publisher can route the event through Pastry and discover the new rendez-vous point. If the rendez-vous point has changed because a new node has arrived, then the old rendez-vous point can forward the publish message to the new rendez-vous point and ask the new rendez-vous point to forward its IP address to the publisher.

There is a single multicast tree for each topic and all publishers use the above procedure to publish events. This allows the rendez-vous node to perform access control.

3.2 Reliability

Publish/subscribe applications may have diverse reliability requirements. Some topics may require reliable and ordered delivery of events, whilst others require only best-effort delivery. Therefore, Scribe provides only best-effort delivery of events but it offers a framework for applications to implement stronger reliability guarantees.

Scribe uses TCP to disseminate events reliably from parents to their children in the multicast tree, and it uses Pastry to repair the multicast tree when a forwarder fails.

Repairing the multicast tree. Periodically, each non-leaf node in the tree sends a heartbeat message to its children. When events are frequently published on a topic, most of these messages can be avoided since events serve as an implicit heartbeat signal. A child suspects that its parent is faulty when it fails to receive heartbeat messages. Upon detection of the failure of its parent, a node calls Pastry to route a SUBSCRIBE message to the topic's identifier. Pastry will route the message to a new parent, thus repairing the multicast tree.

For example, in Figure 4, consider the failure of node 1101. Node 1001 detects the failure of 1101 and uses Pastry to route a SUBSCRIBE message towards the root through an alternative route. The message reaches node 1111 who adds 1001 to its children table and, since it is not a forwarder, sends a SUBSCRIBE message towards the root. This causes node 1100 to add 1111 to its children table.

Scribe can also tolerate the failure of multicast tree roots (rendez-vous points). The state associated with the rendez-vous point, which identifies the topic creator and has an access control list, is replicated across the k closest nodes to the root node in the nodeId space (where a typical value of k is 5). It should be noted that these nodes are in the leaf set of the root node. If the root fails, its immediate children detect the failure and subscribe again through Pastry. Pastry routes the subscriptions to a new root (the live node with the numerically closest nodeId to the topicId), which takes over the role of the rendez-vous point. Publishers likewise discover the new rendez-vous point by routing via Pastry.

Children table entries are discarded unless they are periodically refreshed by an explicit message from the child, stating its continued interest in the topic.

This tree repair mechanism scales well: fault detection is done by sending messages to a small number of nodes, and recovery from faults is local; only a small number of nodes ($O(\log_{2^b} N)$) is involved.

Providing additional guarantees. By default, Scribe provides reliable, ordered delivery of events only if the TCP connections between the nodes in the multicast tree do not break. For example, if some nodes in the multicast tree fail, Scribe may fail to deliver events or may deliver them out of order.

Scribe provides a simple mechanism to allow applications to implement stronger reliability guarantees. Applications can define the following upcall method, which are invoked by Scribe.

forwardHandler(msg) is invoked by Scribe before the node forwards an event, `msg`, to its children in the multicast tree. The method can modify `msg` before it is forwarded.

subscribeHandler(msg) is invoked by Scribe after a new child is added to one of the node's children tables. The argument is the SUBSCRIBE message.

faultHandler(msg) is invoked by Scribe when a node suspects that its parent is faulty. The argument is the SUBSCRIBE message that is sent to repair the tree. The method can modify `msg` to add additional information before it is sent.

For example, an application can implement ordered, reliable delivery of events by defining the upcalls as follows. The `forwardHandler` is defined such that the root assigns a sequence number to each event and such that recently published events are buffered by the root and by each node in the multicast tree. Events are retransmitted after the multicast tree is repaired. The `faultHandler` adds the last sequence number, n , delivered by the node to the SUBSCRIBE message and the `subscribeHandler` retransmits buffered events with sequence numbers above n to the new child. To ensure reliable delivery, the events must be buffered for an amount of time that exceeds the maximal time to repair the multicast tree after a TCP connection breaks.

To tolerate root failures, the root needs to be replicated. For example, one could choose a set of replicas in the leaf set of the root and use an algorithm like Paxos [20] to ensure strong consistency.

4 Experimental evaluation

This section presents results of simulation experiments to evaluate the performance of a prototype Scribe implementation. These experiments compare the performance of Scribe and IP multicast along three metrics: the delay to deliver events to subscribers, the stress on each node, and the stress on

each physical network link. We also ran our implementation in a real distributed system with a small number of nodes.

4.1 Experimental Setup

We developed a simple packet-level, discrete event simulator to evaluate Scribe. The simulator models the propagation delay on the physical links but it does not model either queueing delay or packet losses because modeling these would prevent simulation of large networks. We did not model any cross traffic during the experiments.

The simulations ran on network topologies generated using the Georgia Tech [21] random graph generator. We used a transit-stub model to generate the *core* of the network and we attached a LAN with a star topology to each node in the core. There were an average of 100 nodes in each LAN and each node in a LAN was directly attached by a link to the core node (as was done in [22]). Scribe nodes were assigned randomly to LAN nodes with uniform probability and Pastry was configured with a leaf set size of 16 ($|L| = 16$), a neighborhood set size of 32 ($|M| = 32$), and $b = 4$.

It is difficult to decide on a realistic set of parameters for these topologies and on a realistic assignment of Scribe nodes to topology nodes. Therefore, we studied two very different topologies: a topology with 600 nodes in the core (60,000 LAN nodes), and a topology with 5050 nodes in the core (505,000 LAN nodes). In both cases, we populated the network with 60,000 Scribe nodes. There is a Scribe node assigned to every LAN node in the small topology but only an average of 12% of the nodes in each LAN are Scribe nodes in the large topology. To be conservative, we present results obtained with the large topology because they are less favorable to Scribe.

We used the routing policy weights generated by the Georgia Tech random graph generator [21] to perform IP unicast routing. IP multicast used reverse path forwarding from the source like DVMRP [6]. The delay of each LAN link was set to 1ms and the average delay of core links (computed by the graph generator) was 40.5ms.

Scribe was designed to be a generic infrastructure capable of supporting multiple concurrent applications with varying requirements. Therefore, we ran experiments with a large number of topics and with a wide range of subscriptions per topic. Since there are no obvious sources of real-world trace data to drive these experiments, we adopted a Zipf-like distribution for the number of subscribers to each topic.

The distribution is defined by $Sub(r) = \lfloor Nr^{-1.25} + 0.5 \rfloor$, where r is the rank of the topic, N is the number of Scribe nodes, and $Sub(r)$ is the number of subscribers for the topic with rank r . The number of topics was fixed at 1,500 and the number of Scribe nodes (N) was fixed at 60,000, which were the maximum numbers that we were able to simulate. The exponent 1.25 was chosen to ensure a minimum number of subscribers of six; this number appears to be typical of Instant Messaging applications, which is one of the targeted publish/subscribe applications. Figure 5 shows the number of subscribers versus topic rank. The maximum number of subscribers per topic is 60,000 (topic rank 1) whilst the minimum is 6 (topic rank 1500), and the total number of subscribers to all the topics is 237,154.

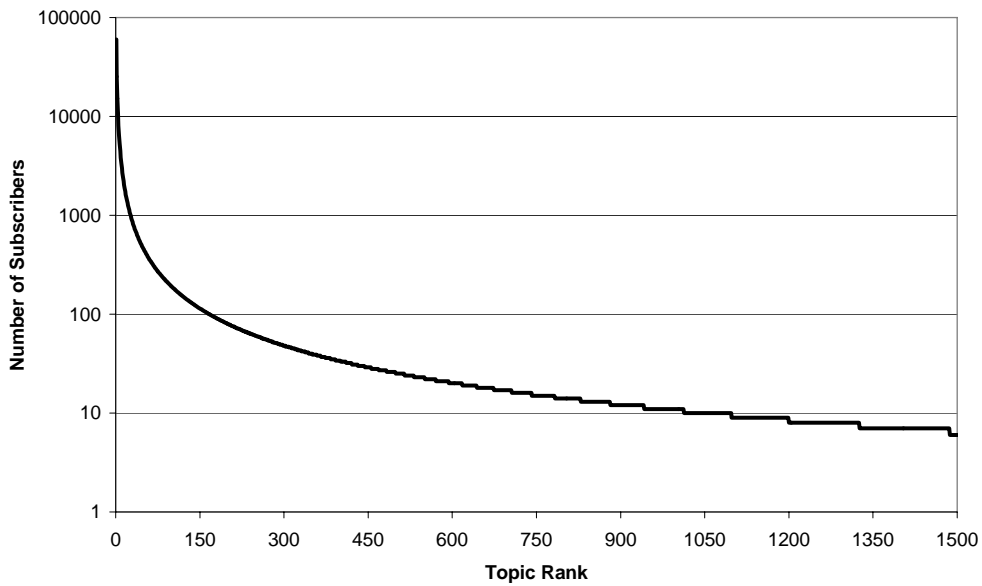


Figure 5: Distribution of the number of subscribers to a topic versus topic rank.

The subscribers to each topic were selected randomly with uniform probability from the set of Scribe nodes. Distributions with better network locality would improve the performance of Scribe.

4.2 Delay penalty

The first set of experiments compares the delay to deliver events using Scribe and IP multicast. For each topic, we measured the delay to deliver an event published at the rendez-vous point to each of the subscribers to the topic. Then, we computed the maximum and mean of this distribution for each topic. We divided the mean and the maximum for each topic in Scribe by the corresponding IP

multicast values to compute the delay penalty.

Figure 6 shows the cumulative distribution of delay penalty per topic for both the mean and the maximum. The y-value of a point represents the number of topics with delay penalty less than or equal to the point's x-value. Scribe's performance is good because it leverages Pastry locality properties. For most topics, it increases the mean delay relative to IP multicast by at most 60% and the maximum delay by at most 50%. In the worst case, it increases the mean delay by 3.7 and the maximum delay by 4.2.

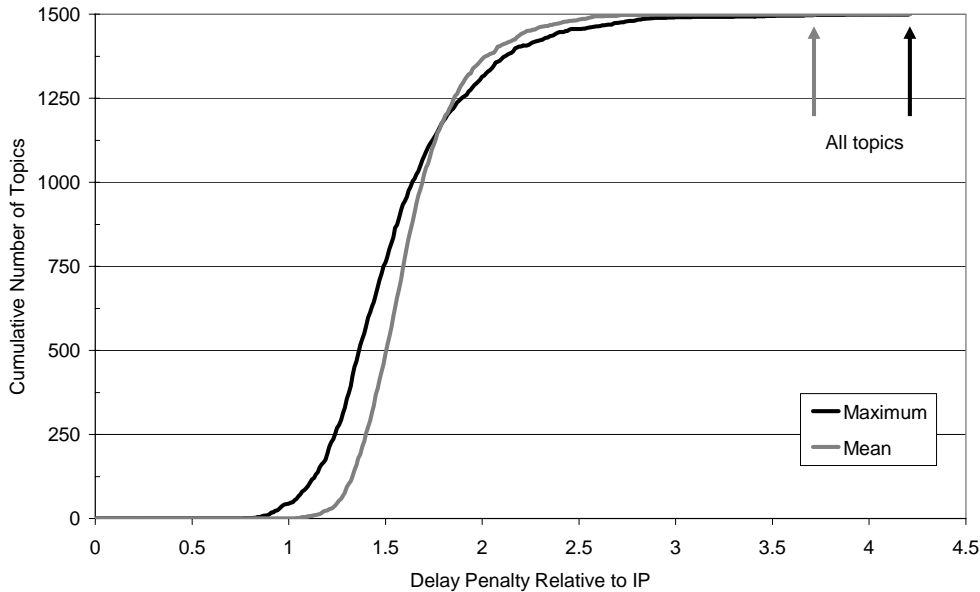


Figure 6: Cumulative distribution of delay penalty relative to IP multicast per topic.

Figure 7 presents a detailed view of the distribution of delays for the 60,000 subscribers of the topic with rank 1. We computed the *relative delay penalty* (RDP) [22] for each subscriber, which is equal to the ratio of the delay to deliver the event to the subscriber using Scribe and IP multicast. The figure plots the cumulative distribution of RDP per subscriber. The y-value of a point is the fraction of subscribers with delay penalty less than or equal to the point's x-value.

The results show that Scribe performs well: the mean RDP is 1.66, the median is 1.56, more than 80% of the subscribers have RDP less than 2, and more than 98% have RDP less than 3. It is important to note that we obtained even better results with the smaller topology where all LAN nodes are Scribe nodes. For this topology, the mean RDP was 1.35, the median was 1.4, and over 97% of the subscribers had an RDP less than 2.

IP routing does not always produce minimum delay routes because it is performed using the routing policy weights from the Georgia Tech model [21]. As a result, Scribe can sometimes find routes with lower delay than IP multicast and 3% of the subscribers in Figure 7 have an RDP less than 1.

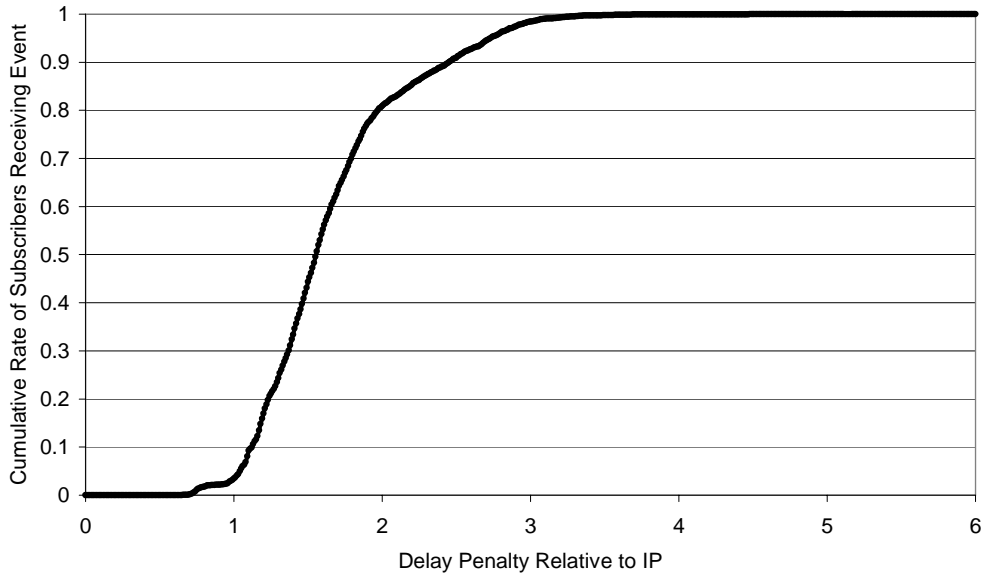


Figure 7: Cumulative distribution of delay penalty relative to IP for the subscribers of the topic with rank 1.

The maximum RDP is 6 whereas the maximum penalty in Figure 6 was less than 4.5. This is because the RDP is computed for each individual subscriber and the delay to the subscriber with RDP 6 is small compared to the delays to other subscribers. Recall that the results in Figure 6 were computed using the mean and maximum of the distribution of delays to all the subscribers of a topic.

4.3 Node stress

In Scribe, end nodes are responsible for maintaining subscription information and for forwarding and duplicating packets whereas routers perform these tasks in IP multicast. To evaluate the stress imposed by Scribe on each node, we measured the number of topics with non-empty children tables and the number of entries in children tables in each Scribe node. Figure 8 shows the number of populated children tables per Scribe node, whilst Figure 9 shows the total number of children table entries per Scribe node.

Even though there are 1,500 topics, the mean number of topics with non-empty children tables per node is only 2.5 and the median number is only 1. Figure 8 shows that the distribution has a long tail

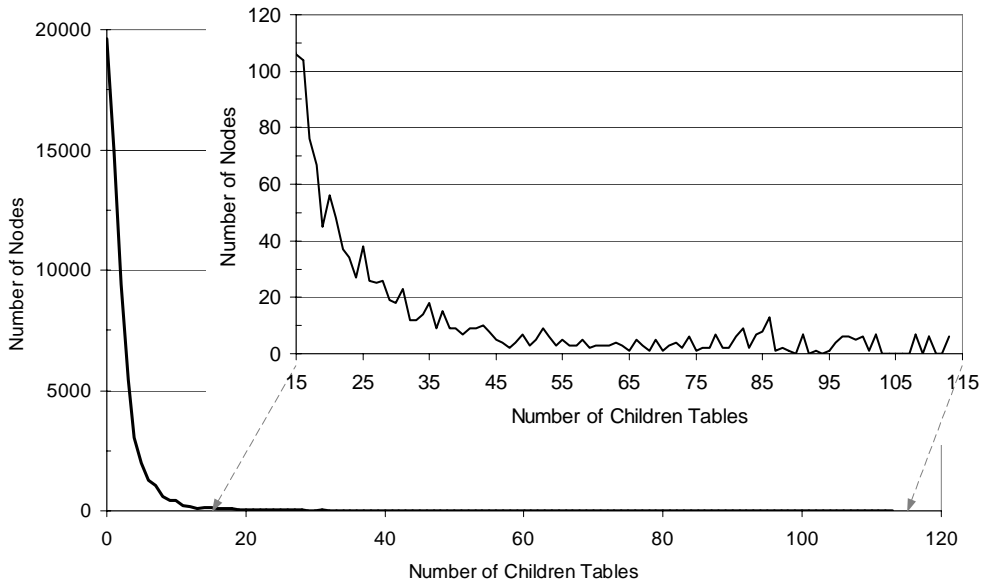


Figure 8: Number of children tables per Scribe node.

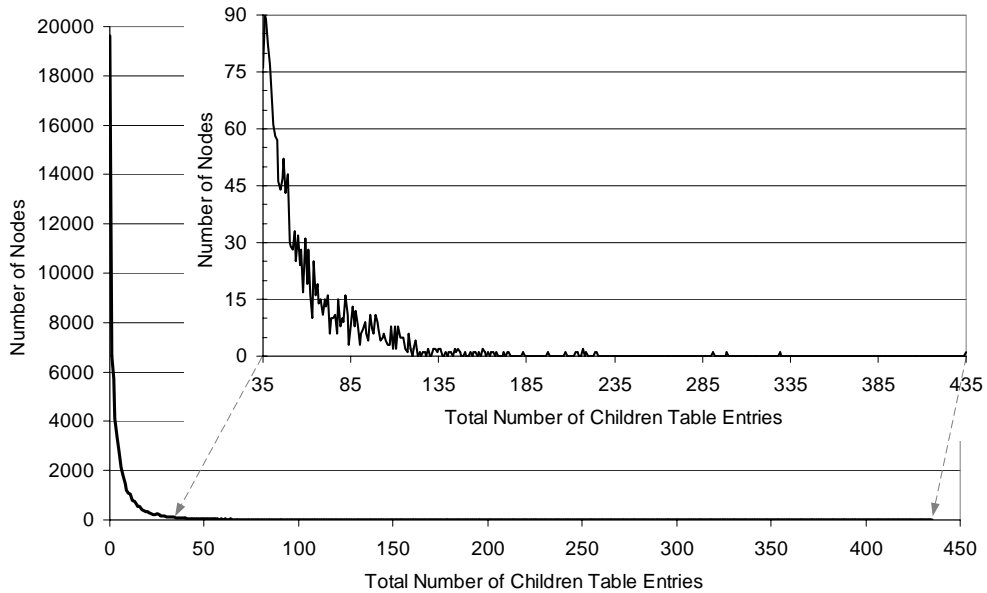


Figure 9: Number of children table entries per Scribe node.

with a maximum number of children tables per node of 113. But this tail is extremely thin as shown in the enlarged section of the graph. Similarly, the mean number of entries in all the children tables of any single Scribe node is only 6.2 and the median is only 2. This distribution also has a long thin tail with a maximum of 435.

These results indicate that Scribe does a good job of partitioning and distributing forwarding load

over all nodes: each node is responsible for forwarding events for a small number of topics, and it forwards events only to a small number of nodes. This is important to achieve scalability with the number of topics and subscribers.

4.4 Link stress

The final set of experiments compares the stress imposed by Scribe and IP multicast on each directed link in the network topology. We computed the stress by counting the number of packets that are sent over each link when an event is published to each of the 1,500 topics. Figure 10 shows the distribution of link stress for both Scribe and IP multicast with the results for zero link stress omitted.

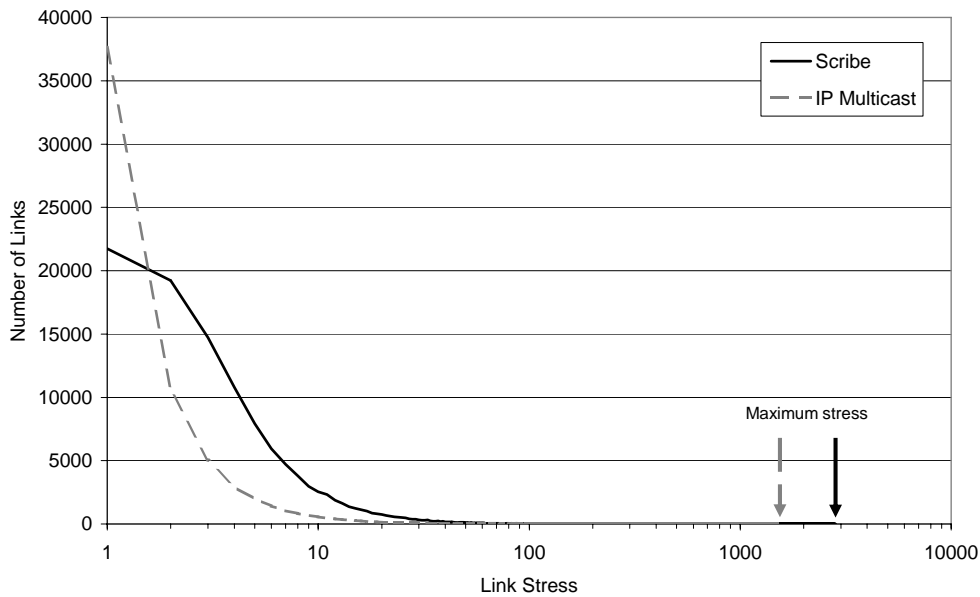


Figure 10: Link stress for publishing an event to each of 1,500 topics.

The total number of links is 1,035,168 and the total number of messages over these links is 1,467,178 for Scribe and 488,547 for IP multicast. The mean number of messages per link in Scribe is 1.4 whilst for IP multicast it is 0.5. The median is 0 for both. The maximum link stress in Scribe is 2798, whilst for IP multicast the maximum link stress is 1500. This means that the maximum link stress induced by Scribe is less than twice that for a IP multicast on this experiment. The results are good because Scribe distributes load across nodes (as shown before) and because it leverages Pastry’s locality properties. Subscribers that are close in the network tend to be children of the same parent in the multicast tree that is also close to them. This reduces link stress because the parent receives a

single copy of the event message and forwards copies to its children along short routes.

It is interesting to compare Scribe with a naïve multicast that is implemented by performing a unicast transmission from the source to each subscriber. This naïve implementation would have a maximum link stress greater than or equal to 60,000 (which is the maximum number of subscribers per topic).

Figure 10 shows the link stress for publishing an event to each topic. The link stress for subscribing is identical because the process we use to create the multicast tree for each topic is the inverse of the process used to disseminate events.

4.5 Tree Management

Results presented in the previous section showed that the load, both in terms of subscribers and topics, is well distributed across the nodes. However, the actual node stress does not only depend on the number of topics and the number of entries in its children tables, but also on the rate of events being published. Thus, it is possible that a node becomes overloaded. To cope with such situations, we have designed and implemented an algorithm that enables a node to off-load children to other nodes.

The algorithm works as follows. When a node is overloaded and wishes to reduce the number of entries in the children table of a topic, it selects the child that is farthest away, according to the proximity metric. The parent drops the chosen child by sending it a message containing a list of the other members of the children table for the topic along with their distance from the parent. When the child receives the message, it performs the following operations: *(i)* It measured the distance between itself and the members of the list it received from the parent; *(ii)* then for each node, it computes the total distance between itself and the old parent via each of the nodes; *(iii)* finally, it sends a subscribe message to the node that provides the smallest combined distance. That way, it minimizes the distance to reach its parent through one of its previous siblings.

We also considered a number of algorithms to reduce the height of the multicast tree. For topics with few subscribers, the basic subscription mechanism produces deep trees that have long paths with no branching. This may unnecessarily increase the number of forwarders, and the latency for event dissemination. We evaluated an algorithm that collapses the Scribe multicast tree by removing nodes which have only one entry in a topic's children table. With the Zipf-like distribution, it did not alter the maximum link stress during publishing, and the global link stress was reduced by approximately two

for every node collapsed. Therefore, we didn't adopt this optimization in the Scribe implementation. However, if Scribe was being used predominantly for topics with small groups of subscribers, there might be benefits in pursuing a tree height reduction algorithm further.

5 Related work

Like Scribe, Overcast [23] and Narada [22] implement multicast using a self-organizing overlay network, and they assume only unicast support from the underlying network layer. Overcast builds a source-rooted multicast tree using end-to-end bandwidth measurements to optimize bandwidth between the source and the various group members. Narada uses a two step process to build the multicast tree. First, it builds a mesh per group containing all the group members. Then, it constructs a spanning tree of the mesh for each source to multicast data. The mesh is dynamically optimized by performing end-to-end latency measurements and adding and removing links to reduce multicast latency. The mesh creation and maintenance algorithms assume that all group members know about each other and, therefore, do not scale to large groups.

Scribe builds a multicast tree on top of a Pastry network, and relies on Pastry to optimize the locality of nodes routed through based on some metric (e.g. IP hops or latency). The main difference is that the Pastry network can scale to an extremely large number of nodes because the algorithms to build and maintain the network have space and time costs of $O(\log_2^b N)$. This enables support for extremely large groups and sharing of the Pastry network by a large number of groups.

The recent work on Bayeux [11] is the most similar to Scribe. Bayeux is built on top of a scalable peer-to-peer object location system called Tapestry [13] (which is similar to Pastry). Like Scribe, it supports multiple groups, and it builds a multicast tree per group on top of Tapestry but this tree is built quite differently. Each request to join a group is routed by Tapestry all the way to the node acting as the root. Then, the root records the identity of the new member and uses Tapestry to route another message back to the new member. Every Tapestry node (or router) along this route records the identity of the new member. Requests to leave the group from the topic are handled in a similar way.

Bayeux has two scalability problems when compared to Scribe. Firstly, it requires nodes to maintain more group membership information. The root keeps a list of all group members, the routers one hop away from the route keep a list containing on average $\frac{S}{b}$ members (where b is the base used in

Tapestry routing), and so on. Secondly, Bayeux generates more traffic when handling group membership changes. In particular, all group management traffic must go through the root. Bayeux proposes a multicast tree partitioning mechanism to ameliorate these problems by splitting the root into several replicas and partitioning members across them. But this only improves scalability by a small constant factor.

In Scribe, the expected amount of group membership information kept by each node is small, as the subscribers are distributed over the nodes, and furthermore, can be bounded by a constant independent of the number of group members, using the hot spot algorithm. Additionally, group join and leave requests are handled locally. This allows Scribe to scale to extremely large groups and to deal with rapid changes in group membership efficiently.

The mechanisms for fault resilience in Bayeux and Scribe are also very different. All the mechanisms for fault resilience proposed in Bayeux are sender-based whereas Scribe uses a receiver-based mechanism. In Bayeux, routers proactively duplicate outgoing packets across several paths or perform active probes to select alternative paths. Both these schemes have some disadvantages. The mechanisms that perform packet duplication consume additional bandwidth, and the mechanisms that select alternative paths require replication and transfer of group membership information across different paths. Scribe relies on heartbeats sent by parents to their children in the multicast tree to detect faults, and children use Pastry to reroute to a different parent when a fault is detected. Additionally, Bayeux does not provide a mechanism to handle root failures whereas Scribe does.

6 Conclusions

We have presented Scribe, a large-scale and fully decentralized event notification system built on top of Pastry, a peer-to-peer object location and routing substrate overlayed on the Internet. Scribe is designed to scale to large numbers of subscribers and topics, and supports multiple publishers per topic.

Scribe leverages the scalability, locality, fault-resilience and self-organization properties of Pastry. The Pastry routing substrate is used to maintain topics and subscriptions, and to build an efficient multicast tree associated with each topic. Scribe's randomized placement of topics and multicast roots balances the load among participating nodes. Furthermore, Pastry's properties enable Scribe to exploit

locality to build an efficient multicast tree and to handle subscriptions in a decentralized manner.

Fault-tolerance in Scribe is based on Pastry's self-organizing properties. The default reliability scheme in Scribe ensures automatic adaptation of the multicast tree to node and network failures. Event dissemination is performed on a best-effort basis; consistent ordering of delivered events is not guaranteed. However, stronger reliability models can be easily layered on top of Scribe.

Our simulation results, based on a realistic network topology model, indicate that Scribe scales well. Scribe is able to efficiently support a large number of nodes, topics, and a wide range of subscribers per topic. Hence Scribe can concurrently support applications with widely different characteristics. Results also show that it balances the load among participating nodes, while achieving acceptable delay and link stress, when compared to network-level (IP) multicast.

References

- [1] Talarian Corporation. *Everything You need to know about Middleware: Mission-Critical Interprocess Communication (White Paper)*. <http://www.talarian.com/>, 1999.
- [2] TIBCO. *TIB/Rendezvous White Paper*. <http://www.rv.tibco.com/whitepaper.html>, 1999.
- [3] P.T. Eugster, P. Felber, R. Guerraoui, and A.-M. Kermarrec. The many faces of publish/subscribe. Technical Report DSC ID:2000104, EPFL, January 2001.
- [4] S. Floyd, V. Jacobson, C.G. Iu, S. McCanne, and L. Zhang. A reliable multicast framework for lightweight sessions and application level framing. *IEEE/ACM Transaction on networking*, pages 784–803, December 1997.
- [5] J.C. Lin and S. Paul. A reliable multicast transport protocol. In *Proc. of IEEE INFOCOM'96*, pages 1414–1424, 1996.
- [6] S. Deering and D. Cheriton. Multicast Routing in Datagram Internetworks and Extended LANs. *ACM Transactions on Computer Systems*, 8(2), May 1990.
- [7] S. Deering, D. Estrin, D. Farinacci, V. Jacobson, C. Liu, and L. Wei. The PIM Architecture for Wide-Area Multicast Routing. *IEEE/ACM Transactions on Networking*, 4(2), April 1996.
- [8] K.P. Birman, M. Hayden, O. Ozkasap, Z. Xiao, M. Budiu, and Y. Minsky. Bimodal multicast. *ACM Transactions on Computer Systems*, 17(2):41–88, May 1999.
- [9] Patrick Eugster, Sidath Handurukande, Rachid Guerraoui, Anne-Marie Kermarrec, and Petr Kouznetsov. Lightweight probabilistic broadcast. In *Proceedings of The International Conference on Dependable Systems and Networks (DSN 2001)*, July 2001.
- [10] Luis F. Cabrera, Michael B. Jones, and Marvin Theimer. Herald: Achieving a global event notification service. In *HotOS VIII*, May 2001.

- [11] Shelly Q. Zhuang, Ben Y. Zhao, Anthony D. Joseph, Randy H. Katz, and John Kubiatawicz. Bayeux: An Architecture for Scalable and Fault-tolerant Wide-Area Data Dissemination. In *Proc. of the Eleventh International Workshop on Network and Operating System Support for Digital Audio and Video (NOSSDAV 2001)*, June 2001.
- [12] Antony Rowstron and Peter Druschel. Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems. In *Proc. IFIP/ACM Middleware 2001*, Heidelberg, Germany, November 2001.
- [13] Ben Y. Zhao, John D. Kubiatawicz, and Anthony D. Joseph. Tapestry: An infrastructure for fault-resilient wide-area location and routing. Technical Report UCB//CSD-01-1141, U. C. Berkeley, April 2001.
- [14] I. Stoica, R. Morris, D. Karger, M. F. Kaashoek, and H. Balakrishnan. Chord: A scalable peer-to-peer lookup service for Internet applications. In *Proc. ACM SIGCOMM'01*, San Diego, CA, August 2001.
- [15] S. Ratnasamy, P. Francis, M. Handley, R. Karp, and S. Shenker. A Scalable Content-Addressable Network. In *Proc. of ACM SIGCOMM*, August 2001.
- [16] Peter Druschel and Antony Rowstron. PAST: A persistent and anonymous store. In *HotOS VIII*, May 2001.
- [17] Antony Rowstron and Peter Druschel. Storage management and caching in PAST, a large-scale, persistent peer-to-peer storage utility. In *Proc. ACM SOSR'01*, Banff, Canada, October 2001.
- [18] FIPS 180-1. Secure hash standard. Technical Report Publication 180-1, Federal Information Processing Standard (FIPS), National Institute of Standards and Technology, US Department of Commerce, Washington D.C., April 1995.
- [19] Yogen K. Dalal and Robert Metcalfe. Reverse path forwarding of broadcast packets. *Communications of the ACM*, 21(12):1040–1048, 1978.
- [20] L. Lamport. The Part-Time Parliament. Report Research Report 49, Digital Equipment Corporation Systems Research Center, Palo Alto, CA, September 1989.
- [21] E. Zegura, K. Calvert, and S. Bhattacharjee. How to model an internetwork. In *INFOCOM96*, 1996.
- [22] Yang hua Chu, Sanjay G. Rao, and Hui Zhang. A case for end system multicast. In *Proc. of ACM Sigmetrics*, pages 1–12, June 2000.
- [23] John Jannotti, David K. Gifford, Kirk L. Johnson, M. Frans Kaashoek, and James W. O'Toole. Overcast: Reliable Multicasting with an Overlay Network. In *Proc. of the Fourth Symposium on Operating System Design and Implementation (OSDI)*, pages 197–212, October 2000.