

XEvent: An Event Notification System over Distributed Hash Table (DHT) Networks

Ruili Wang, Weixiong Rao and Chengqi Zhang

Abstract—In this paper, we propose a novel event notification system, named as XEvent. The system is built over a Distributed Hash Table (DHT)-based Peer-to-Peer (P2P) system by combining a content-based system with a specific event topic. XEvent can support basic topic-based message subscription by use of XML schema as event topic, and further filter the whole message contents using XPath as the filtering language. The unique features of XEvent include: (1) XEvent inherits the excellent features provided by DHT, including scalability, adaptation and self-maintenance; (2) XEvent provides the expressive topic-based and content-based event filtering; (3) XEvent can filter and deliver event messages to subscribers with high efficiency by building a delivery tree based on the subscriber's XPath filters, and (4) XEvent can provide fault-tolerance for node failure.

According to the results of simulation and tests on our XEvent prototype built over a new Bamboo-DHT system, XEvent can achieve the scalability, expressiveness, high efficiency and reliability for the event notification service.

Index Terms— Distributed algorithms, Computer networks, Query languages, Message systems

I. INTRODUCTION

AN event notification service can be implemented by a centralized event broker or by a network of distributed event brokers. The first case, every client application can act as a publisher, a subscriber, or both, and connect to a single central broker server, which acts simply as a message filtering and routing engine. Obviously, this centralized manner, lack of scalability, results in the problem of single point of failure. The second case, multiple servers can act as event brokers and are cooperatively to form a distributed, coherent message routing, matching and filtering engine for clients. This distributed manner can be implemented on a wide-area network (WAN) and be scalable to a large number of clients.

Recently, a new event filtering model, called a content-based event filtering model [1, 2], allows subscriptions to evaluate the entire content of an event message, not just a traditional event topic. The content-based event system provides a more powerful

and flexible event filtering than traditional topic-based event systems [1, 2]. However, what makes content-based event system challenging is that the system must be scalable to a large number of subscribers, publishers and event messages. Also, at the same time, the system must support a flexible and expressive event filtering mechanism, i.e. scalability and flexibility [1].

Distributed Hash Table (DHT) [3, 4, 5] provides unique features for a large scalable distributed system: scalability, adaptation and self-maintenance. Thus, the distributed applications built over DHT can inherit such features provided by DHT including completely decentralized, scalable, and self-organizing; and maintain guaranteed routing efficiency under the joining, departure and failure of nodes in a network. Although some event notification systems [9, 10] built over DHT systems [4, 5] have already been implemented, such systems are topic-based and cannot provide flexibility and expressiveness as they lack a content-based filtering mechanism.

The emergence of XML (eXtensible Markup Language) as a standard for information exchange on the Internet has led to an increased interest in content-based publish/subscribe system [6, 7, 8]. Using XML as the message format can allow structural information within message documents. Using XPath as an XML query language can provide the filtering on both the structure and the contents of published XML information. Although some XML-based filtering systems like [6, 7, 8] can provide a flexible message filtering and matching mechanism, such systems are all a centralized solution. They lack scalability and cannot be extended to a distributed environment.

In this paper, we propose a novel event notification system, named as XEvent. The system is built over a Peer-to-Peer (P2P) DHT networking by combining a content-based system with a specific event topic. XEvent can support basic topic-based message subscription by use of XML schema for event topic, and further filter the entire message contents using XPath as the filtering language. The unique features of XEvent include: (1) XEvent inherits the excellent features provided by DHT, including scalability, adaptation and self-maintenance; (2) XEvent provides the expressive topic-based and content-based event filtering; (3) XEvent can filter and deliver event messages to subscribers with high efficiency by building a delivery tree based on the subscriber's XPath filters, and (4) XEvent can provide fault-tolerance for node failure. According to the results of simulation and tests on our XEvent prototype built over a new Bamboo-DHT system [14] (developed by Berkeley), XEvent can achieve the scalability, expressiveness and high efficiency and reliability for the type-based event notification service.

R. Wang is with the Institute of Information Sciences and Technology, Massey University, New Zealand, Private Bag 11222, Palmerston North, New Zealand. (corresponding author's telephone: 0064-6-3569099 ext 2548; fax: 0064-6-3502259 e-mail: r.wang@massey.ac.nz).

W. Rao is with the Information System and Service Department, Shanghai General Motor, China 201102. (e-mail: rweixiong@yahoo.com.cn).

C. Zhang is with the Faculty of Information Technology, University of Technology, Sydney, Broadway NSW 2007, Australia. (e-mail: chengqi@it.uts.edu.au).

The remainder of the paper is organized as follows. Section II introduces the basic model of XEvent. Section III gives an overview of XEvent. In section IV we describe the basic operations in XEvent. Our XEvent prototype built on Bamboo[14] and experimental results are shown in section V. Related work is introduced in section VI, and section VII gives the conclusion.

II. BASIC MODEL IN XEVENT

View: In XEvent, we decompose an XML message of tree structure into multiple sequential elements from the root to all leaves. We call the sequential elements with the element content as Data View V_d , which includes an element path V_{dp} and element content V_{dc} . V_{dp} is the sequential elements from a root to a leaf, and all elements is separated by the parent-child connector ‘/’ to

<pre><SigmodRecord> <issue> <volume>11</volume> <articles> <article> <title>Process synchronization in Database Systems </title> <initPage>9</initPage> <endPage>29</endPage> <authors> <author position="00"> Philip A. Bernstein</author> <author position="01"> Marco A. Casanova</author> </authors> </article> </articles> </issue> </SigmodRecord></pre>	Element Path
	/SigmodRecord/issue/volume
	/SigmodRecord/issue/articles/article
	/SigmodRecord/issue/articles/article/initPage
	/SigmodRecord/issue/articles/article/endPage
	/SigmodRecord/issue/articles/article/authors/author
	Element Contents
	volume@text=11
	title@text=Process Synchronization in Database Systems
	initPage @text=9 endPage @text=11
author@position="00"@text=Philip A. Bernstein author@position="01"@text=Marco A.Casanova	
(a) an XML message	(b) XML view

Fig. 1. An XML message and its XML views

show the parent-child relationship between two elements. V_{dc} represents the element contents including the element texts and attributes. Fig. 1 shows an XML message and its V_d .

Similarly with V_d , we can define the schema view (V_s) based on an XML schema, which can be also modeled as a tree-like structure. An XML schema only specifies the structure information without content, while V_s only defines the sequential elements from the root to the leaves as the element paths (V_{sp}). Fig. 2 shows an XML schema (DTD) and view V_s .

XPath Model: we use XPath as the filtering language providing an expressive filtering against XML messages in publication. XPath can support both path structure and element content filtering. The path structure filtering in XPath supports a parent-child operator ‘/’ or an ancestor-descendant operator ‘//’. Also XPath allows the use of a wildcard operator ‘*’. XPath also allows element content filtering against both element attributes and element texts. We call the sequential elements in an XPath filter F from the head element of F to the first element followed by a relative path separator in F as a key K of F . For example,

- If $F_1 = /SigmodRecord/issue/articles/article/title$ then $K_1 = F_1 = /SigmodRecord/issue/articles/article/title$
- If $F_2 = /SigmodRecord/issue/articles//title$, then $K_2 = /SigmodRecord/issue/articles$

We define two operations for the XPath filter F : “+” and “-”. If $F_1 = E_1 \wedge \dots \wedge E_{k-1} \wedge E_k \wedge \dots \wedge E_n$ and $F_2 = E_k \wedge \dots \wedge E_n$, then $F_3 = F_1 - F_2 = E_1 \wedge \dots \wedge E_{k-1}$ and $F_3 + F_2 = F_1$. Also, if $F_1 = E_1 \wedge \dots \wedge E_{k-1} \wedge E_k \wedge \dots \wedge E_n$ and $F_3 = E_1 \wedge \dots \wedge E_{k-1}$, then $F_4 = F_1 - F_3 = E_k \wedge \dots \wedge E_n$ and $F_4 + F_3 = F_1$.

A covering relation between two XPath filters is defined that if F_1 is more specific than F_2 , then F_2 covers F_1 . The relation is also written as $F_2 \supseteq F_1$ or $F_1 \subseteq F_2$, if and only if both the path structure filtering P_2 and the element content filtering C_2 of F_2 covers the path structure filtering P_1 and the content filtering C_1 of F_1 . Based on the XPath covering relationship, we can construct an XPath covering graph. In the graph, a node represents an XPath filter and an edge represents the covering relationship between XPath filters. If $F_1 \subseteq F_2$, the node representing F_1 is the parent of the node representing F_2 . As a result, the most specific filter is the root of the graph, and the least specific filter is the leaf node of the graph. Given the filters shown in Fig. 3(b), we know the following XPath covering relations, and the XPath covering relation graph is constructed as Fig. 3(a):

- $F_1 \subseteq F_3 \subseteq F_5 \subseteq F_2 \subseteq F_6 \subseteq F_7 \subseteq F_8$;
- $F_4 \subseteq F_6 \subseteq F_7 \subseteq F_8$.

III. OVERVIEW OF THE XEVENT SYSTEM

In XEvent, we define the XML topic (schema) as the event topic. The event messages in XML format are instances of the

<!ELEMENT SigmodRecord (issue)* >	/SigmodRecord/issue/volume
<!ELEMENT issue (volume,number,articles) >	/SigmodRecord/issue/number
<!ELEMENT volume (#PCDATA)>	/SigmodRecord/issue/articles/article/ title
<!ELEMENT number (#PCDATA)>	/SigmodRecord/issue/articles/article/ initPage
<!ELEMENT articles (article)* >	/SigmodRecord/issue/articles/article/ endPage
<!ELEMENT article(title,initPage,endPage,authors)>	/SigmodRecord/issue/articles/article/ authors/author
<!ELEMENT title (#PCDATA)>	
<!ELEMENT initPage (#PCDATA)>	
<!ELEMENT endPage (#PCDATA)>	
<!ELEMENT authors (author)* >	
<!ELEMENT author (#PCDATA)>	
<!ATTLIST author position CDATA #IMPLIED>	
(a) an XML DTD	(b) view

Fig. 2. Event schema (DTD) and element path of schema views

event schema. XEvent supports subscriptions depending on both the event topic and event content filtering. The introduction of event schema into the event system can be useful for both publishers and subscribers. The schema can validate an XML message and enforce the validity of XML messages; Event subscribers can directly subscribe a certain topic of event messages based on an XML schema, and further subscribe the message using XPath as the filtering language.

We build XEvent (an event notification system) over DHT by arranging a set of cooperating peers in a distributed topology. Each peer in XEvent can act as three roles: publisher, subscriber or

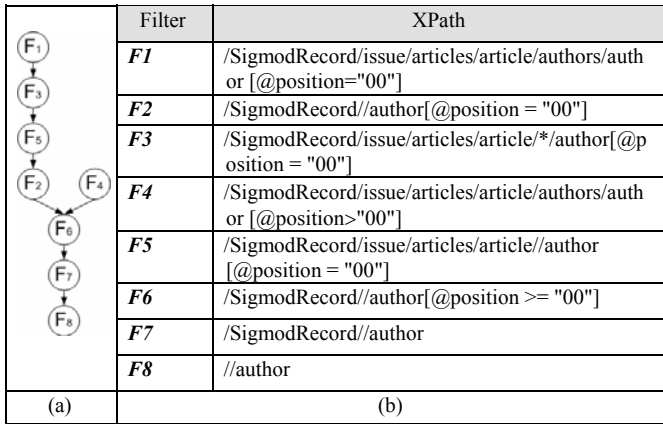


Fig. 3. XPath Relation Graph

broker. The basic operations in XEvent include: schema registration, event subscription and message publication. These operations are the three steps for an event message to be delivered from a publisher to subscribers. Firstly, the registration of an XML schema as event topic into XEvent is mandatory. Secondly, the subscribers can subscribe an interested event by use of the event schema as event topic and an XPath filter as content filter. Thirdly, when a message producer publishes an XML message after the validation of the XML schema, the event message is filtered over XEvent and forwarded by intermediate peers until the message is delivered to its interested subscribers.

In XEvent, we build two layers of overlay network over the IP network (Fig. 4), i.e. event sketch tree and delivery tree. The event sketch tree (described in section 4.1) is built based on the event schema (topic) view. The event delivery tree (described in section 4.2) is organized logically based on the XPath relation graph while physically within the event sketch tree. Using the delivery tree, the event is directly forwarded to its most specific and matched XPath filter node in the event delivery tree (described in section 4.3). If the message matches the most specific XPath filter, all descendant filters of the most specific filter are also matched with the message. Such a delivery tree can greatly reduce the overhead of message routing and matching. Concerning the reliability, an event sketch tree can self-heal its node failure and provide the reliability for the upper event delivery tree (described in section 4.4), and the event delivery tree can devote itself to message filtering and delivery without care of node maintenance.

IV. BASIC OPERATIONS IN XEVENT

In this section, we give the description of XEvent operations: schema registration, event subscription, message publication and maintenance of event delivery tree.

A. Schema Registration

During schema registration, an event skeleton tree T_k is built. The event topic (schema) S , which is a tree structure, can generate views V_s as defined in section 2. The destination host of Hash(S) in DHT namespace is termed as H_t . As the definition of XPath key in section 2, we term F_{root} as the element paths V_{sp} , and K_{root} as

the key of V_{sp} . As a result, $K_{root} = F_{root} = E_1/.../E_n$ since all path separators in $F_{root} = E_1/.../E_n$ are '/'. For an event schema S in XEvent and the element paths F_{root} of V_s , we introduce the destination of Hash(K_{root}) in DHT namespace as the root node of T_k , termed as H_{root} .

When building T_k , F_{root} can be further decomposed into multiple sub-filters F_s and each of F_s is the sequential elements from the head element of F_{root} to other elements of F_{root} . For

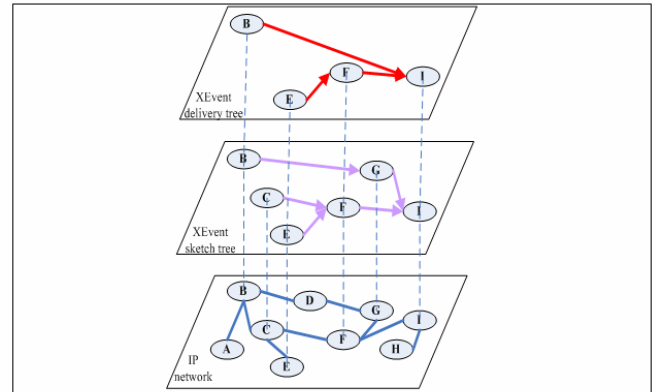


Fig. 4. XEvent Overlay Network

$F_{root} = E_1/.../E_n$, the sub-filters F_s can be $E_1/.../E_{n-1}/E_n$, $E_1/.../E_{n-1}$, ..., E_1/E_2 , E_1 . For each sub-filter F_s , the key K of $F_s = F_s$ because all elements in sub-filter F_s are connected by '/'. Then the skeleton tree T_k , rooted at F_{root} , is built based on the filter covering the relation graph of all sub-filters F_s . Each sub-filter F_s is physically located the destination host of Hash(K) in the DHT namespace. The connection of two sub-filters represents the covering relationship of two sub-filters, the parent sub-filter is more specific than its children, and the root sub-filter is the most specific. Fig. 5 shows the event sketch tree and its

TABLE I
XPath INSTANTIATION

XPath Filter	XPath instantiation Result
F1 //articles//title	/SigmodRecord/issue/articles// title
F2 /*//issue//article/*//author	/SigmodRecord/issue//article/*// author

sub-filters of the event schema (topic) in Fig. 2.

B. Event Subscription

In XEvent, event subscribers can subscribe their interested message by an event schema (topic) and by an XPath filter. As the topic subscription is simple, we only discuss the XPath subscription in this paper. The XPath subscription can be implemented in two phases: XPath instantiation and event delivery tree building.

XPath Instantiation: we have defined the absolute path separator ('/'), and the relative path separator ('/*/' and '/*//'). Depending on whether an absolute path separator '/' at the head of an XPath filter or not, we category an XPath filter F into F_a or F_r , where F_a has an absolute path separator '/' at the head of the XPath filter, and F_r has a relative path separator '/*/' or '/*//' at the head of the XPath filter. For example, //author or /*//issue/articles/

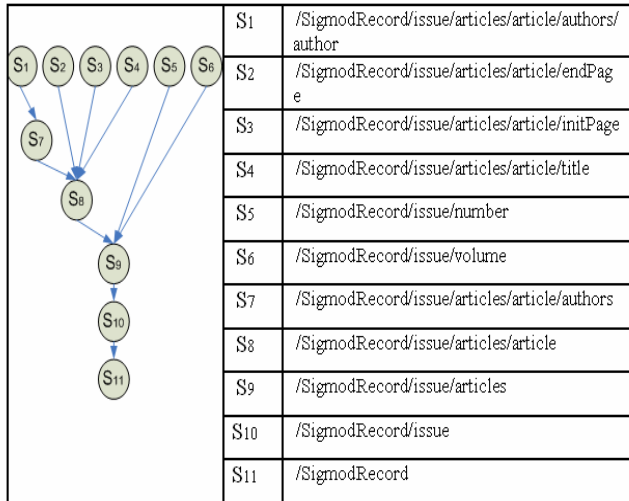


Fig. 5. Event sketch tree and its sub-Filters

are F_r . F_r can decompose into F_r except the relative path separator at the head, i.e. $F_r = \{\wedge Fa\}$ where $\wedge = '/*/'$ or $'//'$. The XPath instantiation actually replaces a relative path separator at the head of F_r into an absolute path separators, i.e. replace the $'/*/'$ or $'//'$ with $'/'$. For these relative path separators, which are located at the intermediate or the end part of an XPath filter, like $/SigmodRecord//author$ or $/SigmodRecord/issue/articles/article/*//author$, no instantiation is needed.

XPath instantiation can be done as follows. Given an even schema S and XPath filter $F_r = \{\wedge Fa\}$ where $\wedge = '/*/'$ or $'//'$, the instantiation result of F_r against S , termed as F_i , is to replaced F_r with $F_i = \{Fs \wedge Fa\}$ where $\wedge = '/'$, if there exists one element path V_{ep} is satisfied with F_r and a sub-filter F_s of V_{ep} where $(V_{ep} - S)^*$ is satisfied with F_a . XPath instantiation can be implemented in the similar way as [6] by use of XPath as Finite State Machines (FSM) and V_{ep} as the message. Table 1 shows two XPath filter, and the result of XPath instantiation against the event schema in Fig. 2.

Building event delivery tree: When an XPath filter F is submitted to XEvent by a subscriber, the XPath filter F is first checked that F is F_a or F_r . If F is F_a , F is directly used to build an event delivery tree; otherwise F needs to instantiate locally. After instantiation, the XPath instantiated result F_i is used to build an event delivery tree.

The key of building an event delivery tree is to make use of the fact: for event schema S , F_{root} of S , and an XPath filter F (F_i or F_a), there always exists a sub-filter F_s of F_{root} , i.e. a sub-filter F_s is XPath key K of F ($K = F_s$) if F is satisfied with the XML event message that is instance of S . For such an XPath filter F , the destination of $Hash(K)$ in DHT namespace, as a result, is the same as the destination of $Hash(F_s)$ where is located at the event sketch tree. The event delivery tree T_d for an event topic (schema) S consists up the XPath filters for the event topic S , which are organized based on the XPath filter relation graph. As a result, T_d , logically organized by an XPath filter covering relation graph, is

physically located within T_k . For a new XPath filter F , the subscription into XEvent can be done in both folders: physically inserted into T_k and logically inserted into T_d . For the former, it can be easily implemented by locating the destination host of $Hash(K)$ since $K = F_s$; for the latter, it can be implemented by finding F 's parent (child) filter at T_d , termed as F_p (F_c), and inserting F between F_p and F_c .

The process of finding F_p is described as follows. Given an XPath filter F , the destination H_k of $Hash(K)$ in DHT is located at T_k . Based on the definition of XPath covering relation and the building of T_k , F_p is physically located at H_k , or at parent node of H_k in T_k , or at ancestor node of H_k in T_k , and impossible at children nodes of H_k in T_k . Then F_p can be found by first traversing all existing XPath filters at H_k to determine the least specific XPath filter as F_p . If no XPath filters exist at H_k , then the same procedure is done to find F_p at parent node of H_k , ..., ancestor node of H_k until the root of T_k . If no XPath filters exist at the parent node of H_k , ancestor node of H_k and the root of the event skeleton tree T_k , F_p returns null, which means F being the most specific filter and becoming the new root of T_k . Finding F_c can be done in the similar way as finding F_c . After insert, the link between F and F_p is established. If both F and F_p are located at the same host, the link within the local host is actually a pointer from F_p to F ; otherwise, a physical link between different hosts is needed. Since multiple XPath filters share the same XPath filter key, these XPath filters with the same key are located at one node in T_k . A physical link can contain multiple links between XPath filters. As a result a message (for example a Heartbeat message) from a parent node to a child node in T_k can aggregate the information of multiple XPath filter links to decrease the communication traffic.

In order to improve the performance of message delivery, we introduce a repair mechanism for the event delivery tree T_d (see section 4.4). Since there exists the case that the root XPath filter of T_d is not located at the root node of T_k , in the repair mechanism a virtual XPath filter located at the root node of T_k is used to link the current root XPath filter of T_d ; Otherwise, the repair mechanism is not needed when the current root XPath filter of T_d is not located at the root node of T_k .

Concerning the efficiency of finding F_p , the IP hops of finding F_p (F_c) is limited within the depth of T_d not the number of XPath filters, which can greatly decrease the communication cost. Thus, the cost of finding the least specific filter in a local host is negligible compared with the communication cost in DHT-based P2P system. Fig. 6(a) shows the XPath filter covering relation graph after instantiation of the XPath filters in Fig. 3(b), and Fig. 6(b) shows the event delivery tree T_d and its physical location in T_k . Here note: (1) F_8 ($//author$) after XPath instantiation becomes $/SigmodRecord/issue/articles/article/authors/author$. Then the new XPath covering relation graph in Fig. 6(a) is twisted from the one shown in Fig. 3(a). (2) The filters in T_d are located within T_k , but the physical link between filters in T_d does not follow the physical link in T_k . S_1 and S_8 in T_d is directly linked not via S_7 .

C. Message Publication

In XEvent, only the valid XML message against the event topic (schema) can be published into XEvent. A valid XML message

* About definition of -, see the XPath model in section 2

can be decomposed into XML views V_d with element contents V_{dc} and element path V_{dp} . For V_d , the destination of $\text{Hash}(V_{dp})$ in the DHT is a root node H_{root} of the event sketch tree T_k for the event schema as described in section 4.1. Based on the repair mechanism of the event delivery tree T_d , there is always a root XPath filter F_r at H_{root} . If F_r is virtual, the incoming V_d is delivered to its child XPath filter for further message filtering. Otherwise, V_d is matched with F_r . If V_d is satisfied with the F_r , V_d is also satisfied with all of its children. As a result, the XML message can be directly routed to the subscribers who subscribe F_r and all children of F_r . If the XML view V_d is unsatisfied with F_r , the traversing the event delivery tree from F_r to its children is needed until the XML view V_d is satisfied with an XPath filter F in T_d . Then, we can determine that the XPath filter F and its all children filters in T_d are interested in the XML message. Also, the XML message is delivered to the subscribers who subscribe F and its children filters. If an XPath filter F receives a previous duplicated V_d , the XML view V_d will not be forwarded and it will be discarded. In section 4.2, we introduce the repair mechanism of T_d . We argue that if a virtual XPath filter exists at the root of T_k , the matching can directly be forwarded to its child filter without traversing the child node to find the root filter of T_d .

About efficiency of message publication, since XPath filters with the same key are located at the same node of T_k , the forwarding cost within the local node is negligible compared with the forwarding cost between physical nodes. Also, the forwarding hops are limited within the depth of T_k not of T_d , which can greatly decrease the communication cost.

D. Maintenance of an Event Delivery Tree

The requirement of reliability needs XEvent to be robust against node failure. To repair the node failure in T_k , a node in T_k is required to periodically exchange Heartbeat messages with its parent/child nodes to detect failure. Each XPath filter in the node has a soft state with leases for its parent/child XPath filter.

We describe failures and their repair mechanisms as follows. Firstly, failure of H_t can be detected by the Heartbeat message from the root node H_{root} in the event sketch tree T_k , but the failure cannot cause any hamper of message publication or event subscription except for those XPath filters which need

while without re-building the entire event delivery tree T_d . Secondly, failure of H_{root} in T_k can be detected by the Heartbeat message from H_t , and can be healed by H_t to re-lookup the destination host of $\text{Hash}(F_{root})$ at DHT as the new H_{root} . Thirdly, failure of any other node in T_k except H_{root} can be detected by its parent node. The parent node, destination of $\text{Hash}(E_1/\dots/E_{k-1}/E_k)$ in DHT, can re-lookup the destination of Hashing $(E_1/\dots/E_{k-1})$ to heal the failed child.

Even before the node in T_k is recovered, XEvent provides a reliable message forwarding mechanism. The XML view V_d can be first routed to the children node of T_k , not the failure node where its children XPath filter is located. This begins the message matching against the XPath filters at the children nodes.

V. EXPERIMENTAL EVALUATION

In order to evaluate XEvent performance, we implement the prototypes of basic SCRIBE [9] and XEvent using the discrete event simulator over Bamboo [14], a recent DHT-based P2P system. The simulation environment is a network topology with 5050 routers, which were generated by the Georgia Tech random graph generator using the transit-stub model [15]. The prototype codes including SCRIBE and XEvent run on 100,000 end nodes that were randomly assigned to routers in the core with uniform probability. Each end system was directly attached by a LAN link to its assigned router. During the evaluation, we designed the following test scenario: (1) *SigmodRecord.dtd*. [16] as an event topic (schema) is registered into XEvent; (2) 4096 XPath filters are generated by the IBM XPath query generator[17] to subscribe their interested messages; (3) XML-based event messages, which are validate with the XML event topic (schema), are published into XEvent continuously as a message stream. Finally, XEvent routes the published messages to all satisfied subscribers. From the time that the messages are published to the time that all satisfied subscribers get the messages, we evaluate XEvent and SCRIBE. For prototype implementation of SCRIBE, we filter the coming message content at the end host of SCRIBE, which is the destination of Hash (XPath key).

To facilitate our comparison, we use two matrices: Relative

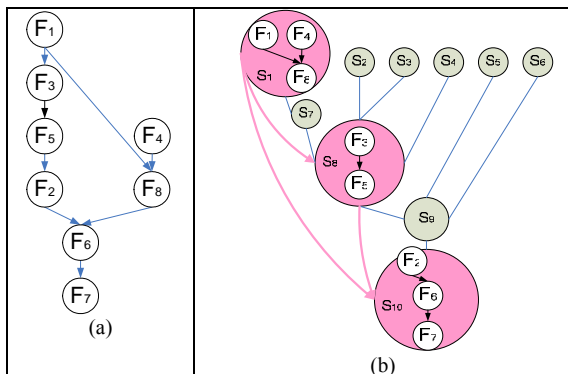


Fig. 6. Event Delivery Tree

instantiation at H_t . The failed H_t can be healed by any publisher, subscriber or third-party to re-register the event topic (schema)

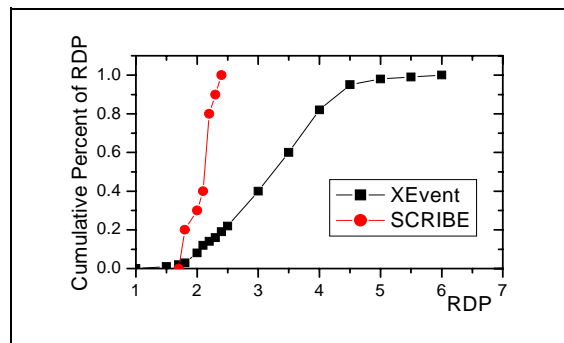


Fig. 7. The cumulative distribution of Relative Delay Penalty (RDP) of both XEvent and SCRIBE.

Delay Penalty (RDP) and physical link stress, which were originally proposed in [13] and were sequentially used in [10]. In addition to these two matrices, stability is also evaluated.

Relative Delay Penalty (RDP) is a measure of the increase in delay that applications incur while using overlay routing [13, 10]. Using RDP on both XEvent and SCRIBE [10], it is the ratio of DHT overlay unicast routing distances to IP unicast routing distances. Assuming symmetric routing, IP Multicast and naive unicast both have a RDP of 1. We use the ratio of the amount of total overlay DHT unicast routing hops to the amount of all naive unicast hops as the RDP. Fig. 7 shows the cumulative distribution of RDP of both XEvent and SCRIBE. The X-axis represents the RDP value and y-axis represents the cumulative percent of RDP. XEvent has an average RDP value of 1.9 and SCRIBE has the value of 4.2. It can be explained that all communication in XEvent is within the event sketch tree, while the communication in SCRIBE is scattered through out the whole network.

Link Stress is a measure of how effective SCRIBE and XEvent prototypes are in a distributed network load across different physical links. It refers to the number of identical copies of a packet carried by a physical link. IP multicast has a stress of 1, and naive unicast has a worst case stress equal to number of receivers. We calculate the ratio of the total number of messages that are sent over links to the total number of links as the average link stress. Fig. 8 shows the link stress comparison of XEvent and Scribe. The X-axis represents the Log scale of Link stress, i.e. the duplicate message copies across the physical link. The Y-axis represents the Log scale of link number. Obviously, XEvent has much less link stress than SCRIBE. It can be easily explained that the link between the nodes in XEvent event sketch tree can contain multiple logical links between XPath filters of XEvent event delivery tree. Also the logical link of XPath filters within the same local host does not require any physical communication.

Stability: To evaluate the stability of XEvent, we need to design an unstable scenario by re-designing step 3 in normal scenarios. We separately design 10%, 20% and 30% nodes crashed in the network model from normal message publication in the time interval of 1000 seconds, then we count the ratio of the nodes, which can receive the message, to all live nodes. In Fig. 9, X-axis represents time and Y-axis represents the ratio in percent. From the figure, the time interval for XEvent (with the self-healing mechanism) to recover 100% message receiving is shorter than for SCRIBE. Also, the failure makes the data lost of some nodes in SCRIBE since we simply implement the basic SCRIBE without the group maintenance mechanism.

VI. RELATED WORKS

Recent work on the scalable design of structured P2P overlay networks has introduced a new class of structured networks called Distributed Hash Tables (DHT). Well known representatives include [3, 4, 5]. All of these systems were built to allow efficient key lookups. Nodes in a P2P network send messages to each other based on a unique name, generated from a secure one-way hash of some unique string. Assigned a unique ID to each node by DHT, a message is delivered to the destination host in a fault-resilient fashion. The P2P routing and locating infrastructure can provide scalability, fault-tolerance, self-maintenance and adaptation for

upper applications.

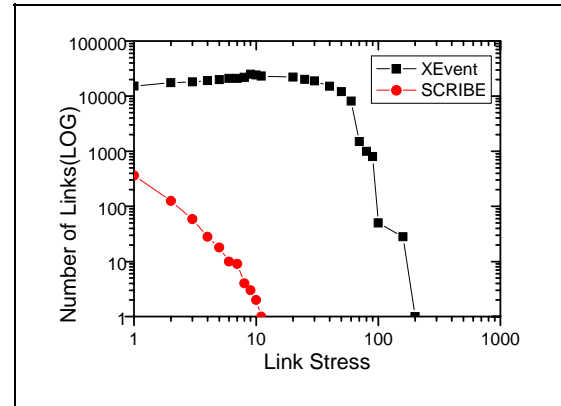


Fig. 8. The link stress comparison of XEvent and Scribe.

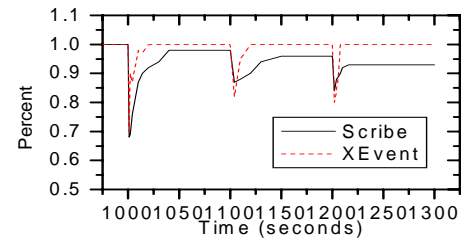


Fig. 9. XEvent stability

Several event notification systems have already developed over DHT. SCRIBE [9] is a topic-based event notification system that is built on top of “Pastry” [4]. “Pastry” is a DHT system developed at Rice University. Subscribers join topics of interest, where each topic is identified with a Pastry-level key. Bayeux [10] is built on top of Tapestry [5]. However, neither SCRIBE [9] nor Bayeux [10] supports content-based subscriptions and publications. Hermes [11] provides a type- and attribute-based routing schema over DHT, but the introduction the rendezvous node into Hermes can result in limited scalability and a single point of failure, also the content filtering is rather limited. An extension of traditional Bloom filters, called multi-level Bloom filters [12], can be used to route path queries in a P2P system, and build content-based overlay networks by linking together peers with similar content. However, the P2P system is totally different from the DHT-based structured P2P systems. Therefore, the system will not or cannot implement a flexible filtering mechanism like XEvent.

Some pioneer researchers in database field have incorporated XML into publish /subscribe system [6, 8]. An XML document filtering system, for Selective Dissemination of Information (SDI) has been proposed in [6], which allows users to define their interests using XPath query language. Using a modified Finite State Machine (FSM), a novel indexing mechanism and matching algorithms are developed. YFilter, combines all of the queries into a single Non-deterministic Finite Automaton (NFA) [8]. The approach exploits commonality among path expressions by merging common prefixes of the query paths so that they are

processed at once. All of these approaches are based a single centralized filtering engine and are hard to scale to a large distributed network.

TABLE II
ABBREVIATIONS LIST

V_d	XML Data View
V_{dp}	Element path in V_d
V_{dc}	Element content in V_d
V_s	Schema View
V_{sp}	Element path in V_s
F	XPath filter
S	Event schema (topic)
K	Key of F
F_{root}	Element paths V_{sp}
K_{root}	Key of F_{root}
H_{root}	Destination Host of F_{root} in DHT namespace
T_k	Event sketch tree
T_d	Event delivery tree
F_s	XPath Sub-Filter of V_{sp}
F_a	XPath Absolute Filter
F_r	XPath Relative Filter
F_i	XPath Instantiation result
F_p	XPath Parent filter in T_d
F_c	XPath child filter in T_d
H_t	Destination Host of Hash(S) in DHT namespace

VII. CONCLUSION

We have presented XEvent, a novel hybrid event notification system built over Distributed Hash Table (DHT)-based Peer-to-Peer (P2P) system. XEvent can support both event topic and event content subscription, which can provide expressiveness for subscribers by the use of XML as event message format and use of XPath as subscription filter.

XEvent make use of DHT to build two layers of an overlay network: event sketch tree and event delivery tree. An event sketch tree can provide the stability for the upper event delivery tree. The event delivery tree can devote itself to swiftly match and deliver messages to the interested subscribers, without considering the churn in DHT.

The initial simulation results indicate that XEvent can provide an efficient event notification system with quick recovery mechanisms and low Relative Delay Penalty and link stress compared with the SCRIBE [9] system.

REFERENCES

[1] A. Carzaniga, D. S. Rosenblum, and A. L. Wolf, "Content-based addressing and routing: A general model and its application," Department of Computer

- Science, University of Colorado, USA, Technical Report CU-CS-902-00, Jan. 2000.
- [2] R. Chand and P. Felber, "A scalable protocol for content-based routing in overlay networks," in *IEEE International Symposium on Network Computing and Applications (NCA'03)*, Cambridge, Massachusetts, Apr. 2003.
- [3] I. Stoica, R. Morris, D. Karger, F. Kaashoek, and H. Balakrishnan, "Chord: A scalable peer-to-peer lookup service for Internet applications," in *Proceedings of ACM SIGCOMM'01*, San Diego, CA, Aug. 2001.
- [4] A. Rowstron and P. Druschel, "Pastry: Scalable, distributed object location and routing for large-scale peer-to-peer systems," in *MiddleWare*, Germany, Nov. 2001.
- [5] B. Y. Zhao, J. D. Kubiatowicz, and A. D. Joseph, "Tapestry: An infrastructure for fault-tolerant wide-area location and routing," University of California Berkeley, USA, Technical Report CSD-01-1141, 2001.
- [6] M. Altinel and M. J. Franklin, "Efficient Filtering of XML Documents for Selective Dissemination of Information," in *Proceedings of the 26th International Conference on Very Large Data Bases*, Cairo, 2000, pp 53–64.
- [7] A. K. Gupta and D. Suci, "Stream processing of XPath queries with predicates," in *Proceedings of the ACM SIGMOD International Conference on Management of Data*, California, USA, 2003, pp. 419–430.
- [8] Y. Diao and M. Franklin, "Query Processing for High-Volume XML Message Brokering," in *Proceedings of the 29th International Conference on Very Large Data Bases*, Berlin, Germany, 2003, pp. 261–272.
- [9] M. Castro, P. Druschel, A. Kermarrec, and A. Rowstron, "SCRIBE: A large-scale and decentralized application-level multicast infrastructure," *IEEE Journal on Selected Areas in Communications*, vol. 20, no. 8, pp. 1489–1499, Oct. 2002.
- [10] S. Q. Zhuang, B. Y. Zhao, A. D. Joseph, R. H. Katz, and J. Kubiatowicz, "Bayeux: An Architecture for Scalable and Fault-tolerant WideArea Data Dissemination," in *Proceedings of the 11th International Workshop on Network and Operating System Support for Digital Audio and Video*, USA, June 2001.
- [11] P. R. Pietzuch and J. Bacon, "Peer-to-Peer Overlay Broker Networks in an Event-Based Middleware," in *Proceedings of the 2nd International Workshop on Distributed Event-based Systems*, USA, 2003.
- [12] G. Koloniari, Y. Petrakis and E. Pitoura, "Content-Based Overlay Networks of XML Peers Based on Multi-Level Bloom Filters," in *International Workshop on Databases, Information Systems and Peer-to-Peer Computing*, Germany, 2003.
- [13] Y. Chu, S. G. Rao, and H. Zhang, "A case for end system multicast (keynote address)," in *Proceedings of the 2000 ACM SIGMETRICS International Conference on Measurement and Modeling of Computer Systems*, USA, 2000.
- [14] S. Rhea, D. Geels, T. Roscoe, and J. Kubiatowicz, "Handling Churn in a DHT," in *Proceedings of the USENIX Annual Technical Conference (USENIX)*, Boston, Massachusetts, June 2004.
- [15] Georgia tech internet topology model.
<http://www.cc.gatech.edu/fac/Ellen.Zegura/graphs.html>
- [16] <http://www.cs.washington.edu/research/xmldatasets/>
- [17] <http://www.alphaworks.ibm.com/tech/xmlgenerator>