Artificial Intelligence formulated this projection for compatibility purposes from the original article published at Global Journals. However, this technology is currently in beta. *Therefore, kindly ignore odd layouts, missed formulae, text, tables, or figures.*

1 2	A Multicast Protocol For Content-Based Publish-Subscribe Systems
3	Prof. Vishal Nagar ¹ , Yashpal Singh ² and D.C. Dhubkarya ³
4	¹ CSe Jhnai-India,UPTU Lucknow-India
5	Received: 10 February 2010 Accepted: 3 March 2010 Published: 15 March 2010

7 Abstract

6

The publish/subscribe (or pub/sub) paradigm is a simple and easy to use model for
interconnecting applications in a distributed environment. Many existing pub/sub systems are

- ¹⁰ based on pre-defined subjects, and hence are able to exploit multicast technologies to provide
- scalability and availability. An emerging alternative to subject-based systems, known as
- ¹² content-based systems, allow information consumers to request events based on the content of
- ¹³ published messages. This model is considerably more flexible than subject-based pub/sub,
- ¹⁴ however it was previously not known how to efficiently multicast published messages to
- ¹⁵ interested content-based subscribers within a network of broker (or router) machines. In this
- ¹⁶ paper, we develop and evaluate a novel and efficient technique for multicasting within a
- ¹⁷ network of brokers in a content-based subscription system, thereby showing that
- ¹⁸ content-based pub/sub can be deployed in large or geographically distributed settings.
- 19

20 Index terms— Multicast, Publishers, Subscribers, information spaces, OMG, content-based routing

²¹ 1 INTRODUCTION

22 he publish/subscribe paradigm is a simple, easy to use and efficient to implement paradigm for interconnecting applications in a distributed environment. Pub/sub based middleware is currently being applied for application in-23 24 tegration in many domains including financial, process automation, transportation, and mergers and acquisitions. 25 Pub/sub systems contain information providers, who publish events to the system, and information consumers, who subscribe to particular categories of events within the system. The system ensures the timely delivery of 26 27 published events to all interested subscribers. A pub/sub system also typically contains message brokers that are responsible for routing messages between publishers and subscribers. The earliest pub/sub systems were subject-28 based. In these systems, each unit of information (which we will call an event) is classified as belonging to one 29 of a fixed set of subjects (also known as groups, channels, or topics). Publishers are required to label each event 30 with a subject; consumers subscribe to all the events within a particular subject. For example a subject-based 31 pub/sub system for stock trading may define a group for each stock issue; publishers may post information to the 32 appropriate group, and subscribers may subscribe to information regarding any issue. In the past decade, systems 33 supporting this paradigm have matured significantly resulting in several academic and Yashpal Singh 1, Vishal 34 35 Nagar 2, D.C.Dhubkarya 1 yash_biet@yahoo.com, vishal_biet@yahoo.com, dcd3580@yahoo.com 1.BIET Jhansi, 36 2. CSE Jhansi , India industrial strength solutions [4][10] [12][13] [15]. A similar approach has been adopted 37 by the OMG for CORBA event channels [11]. An emerging alternative to subject-based systems is content-based subscription systems [6] [14]. These systems support a number of information spaces, each associated with an 38 event schema defining the type of information contained in each event. Our stock trade example (shown in Figure 39 ??) may be defined as a single information space with an event schema defined as the tuple [issue: string, price: 40 dollar, volume: integer]. A contentbased subscription is a predicate against the event schema of an information 41 space, such as (issue=?IBM? & price < 120 & volume > 1000) in our example. With content-based pub/sub, 42 subscribers have the added flexibility of choosing filtering criteria along as many dimensions as event attributes, 43

$\mathbf{2}$ THE MATCHING ALGORITHM

without requiring pre-definition of subjects. In our stock trading example, the subject-based subscriber is forced 44 to select trades by issue name. In contrast, the content-based subscriber is free to use an orthogonal criterion, 45 such as volume, or indeed a collection of criteria, such as issue, price and volume. Furthermore, content-based 46 pub/sub removes the administrative overhead of defining and maintaining a large number of groups, thereby 47 making the system easier to manage. Finally, content-based pub/sub is more general in that it can be used to 48 implement subjectbased pub/sub, while the reverse is not true. While contentbased pub/sub is the more powerful 49 paradigm, efficient and scalable implementations of such systems have previously not been developed. In order 50 to efficiently implement a content-based pub/sub system, two key problems must be solved: 51

i. The problem of efficiently matching an event against a large number of subscribers on a single message 52 broker. ii. 53

The problem of efficiently multicasting events within a network of message brokers. This problem becomes 54 crucial in two settings: 1) when the pub/sub system is geographically distributed and message brokers are 55 connected via a relatively low speed WAN (compared to high-speed LANs), and 2) when the pub/sub system has 56 to scale to support a large number of publishers, subscribers and events. In both cases it becomes crucial to limit 57 the distribution of a published event to T only those brokers that have subscribers interested in that event. One of 58 59 the strengths of subject-based pub/sub systems is that both problems are trivial to solve: the matching problem 60 is solved using a mere table lookup; the multicast problem is solved by defining a multicast group per subject, 61 and multicasting each event to the appropriate multicast group. For content-based pub/sub systems, however, 62 previous literature does not contain solutions to either problem, matching or multicasting. In this paper we present the first efficient solution to the multicast problem for content-based pub/sub. In a companion paper 63 [2] we present an efficient soution to the matching problem for these systems. There are two straightforward 64 approaches to solving the multicasting problem for content-based systems: (1) in the match-first approach, the 65 event is first matched against all subscriptions, thus generating a destination list and the event is then routed to 66 all entries on this list; and (2) in the flooding approach, the message is broadcast or flooded to all destinations 67 using standard multicast technology and unwanted messages are filtered out at these destinations. The match-68 first approach works well in small systems, but in a large system with thousands of potential destinations, the 69 increase in message size makes the approach impractical. Further, with this approach we may have multiple 70 copies of the same message going over the same network link on its way to multiple remote subscribers. The 71 flooding approach suffers when, in a large system, only a small percentage of clients want any single message. 72 73 Furthermore, the flooding technique cannot exploit locality of information requests, i.e., when clients in a single 74 geographic area are, for many applications, likely to have similar requests for data. The central contribution of this paper is a new protocol for content-based routing, an efficient solution to 75 the multicast problem for content-based pub/sub systems. With this protocol, called link matching, each broker 76

partially matches events against subscribers at each hop in the network of brokers to determine which brokers to 77 send the message. Further, each broker forwards messages to its subscribers based on their subscriptions. The 78 disadvantages of the match-first approach are avoided since no additional information is appended to the message 79 headers. Further, at most one copy of a message is sent on each link. The disadvantages of the flooding approach 80 are avoided as the message is only sent to brokers and clients needing the message, thus exploiting locality. We 81 illustrate, using a network simulator, that flooding overloads the network at significantly lower publish rates 82 than link matching. We also describe our implementation of a distributed Java based prototype of content-based 83 pub/sub brokers. The remainder of this paper is organized as follows. In section 2, we present a solution to the 84 matching problem (i.e., the case when the network consists of a single broker). 85

In section 3, we discuss how to extend the solution to the matching problem into a solution to the content-86 based routing problem in a multi-broker network. In section 4, we evaluate the performance of this approach and 87 compare it to the flooding approach. 88 II.

89

$\mathbf{2}$ THE MATCHING ALGORITHM 90

This section summarizes a non-distributed algorithm for matching an event against a set of subscriptions, and 91 returning the subset of subscriptions that are satisfied by the event. (A more detailed presentation of matching 92 along with experimental and analytic measures of performance are the subject of our companion paper [2].) This 93 matching algorithm is the basis of our distributed multicast protocol, presented in the following section. 94

Our approach to matching is based on sorting and organizing the subscriptions into a parallel search tree 95 (or PST) data structure, in which each subscription corresponds to a path from the root to a leaf. The 96 97 matching operation is performed by following all those paths from the root to the leaves that are satisfied by the 98 event. Intuitively, this data structure yields a scaleable algorithm because it exploits the commonality between 99 subscriptions as shared prefixes of paths from root to leaf. Figure ?? shows an example of a matching tree for 100 an event schema consisting of five attributes a1 through a5. These attributes could represent, for example, the stock issue, price, or volume attributes mentioned above. The root of the tree corresponds to a test of the value 101 of attribute a1, the nodes at the next level correspond to a test of attribute a2, etc. The branches are labeled 102 with the values of the attributes being tested. In the example, we only show equality tests (although range tests 103 are also possible), so the right branch of the root represents the test $a_1 = 1$. The left branch of the root, with 104 label *, means that the subscriptions along that branch do not care about the value of the attribute. Each leaf 105

is labeled with the identifiers of all the subscribers wishing to receive events matching the predicate, i.e., all the tests from the root to the leaf. For example, in Figure ??, the rightmost leaf corresponds to a subscription whose predicate is (a1 = 1 & a2 = 2 & a3 = 3 & a5 = 3)

). Since a4 does not appear in this subscription, it is represented by a label * in the PST. Given this tree 109 representation of subscriptions, the matching algorithm proceeds as follows. We begin at the root, with current 110 attribute a 1. At any non-leaf node in the tree, we find the value v j of the current attribute a j. We then traverse 111 any of the following edges that apply: (1) the edge labeled v j if there is one, and (2) the edge labeled * if there 112 is one. This may lead to either 0, 1, or 2 successor nodes (or more in the general case where the tests are not all 113 strict equalities). We initiate parallel subsearches at each successor node. When any of the parallel subsearches 114 reaches a leaf, all subscriptions at that leaf are added to the list of matched subscriptions. For example, running 115 the matching algorithm with the matching tree of Figure ?? and the event $a = \langle 1, 2, 3, 1, 2 \rangle$ will visit all the 116 nodes marked with dark circles and will match four subscription predicates, corresponding to the dark circles at 117 leaf nodes. 118

The way in which attributes are ordered from root to leaf in the PST can be arbitrary. In our experience, however, performance seems to be better if the attributes near the root are chosen to have the fewest number of subscriptions labeled with a *.

In the companion paper [2], we have analytically shown that the cost of matching using the above algorithm increases less than linearly as the number of subscriptions increase.

3 A. Optimizations

A number of optimizations may be applied to the parallel search tree to decrease matching time –these optimizations are explained fully in [2].

Factoring-Some search steps can be avoided, at the cost of increased space, by factoring out certain attributes. That is, certain attributes —preferably those for which the subscriptions rarely contain -don't care? tests —are selected as indices. A separate subtree is built for each possible value (or for ranges, each distinguished value range) of the index attributes.

Trivial Test Elimination-Nodes with a single child which is reached by a *-branch may be eliminated. Delayed Branching-Following *-branches may be delayed until after a set of tests have been applied. This optimization prunes paths from that *-branch which are inconsistent with the tests.

134 It is worth noting that, under certain circumstances, after applying optimizations, the parallel search tree will 135 no longer be a tree but instead a directed acyclic graph.

136 III.

137 4 THE LINK MATCHING ALGORITHM

The previous section described a non-distributed algorithm for matching events to subscriptions. This section presents the central contribution of this paper –an extended matching algorithm for a network of brokers, publishers, and subscribers (as shown in Figure ??). The problem, in this case, is to efficiently deliver an event from a publisher to all distributed subscribers interested in the event.

One straightforward solution to this problem is to perform the matching algorithm of the previous section 142 at the broker nearest to the publisher, producing a destination list consisting of the matched subscribers. This 143 destination list may be undesirably long in a large network with thousands of subscribers, and it may be infeasible 144 to transmit and process large messages containing long destination lists throughout the network. Link matching 145 is our strategy for multicasting events without using destination lists. After receiving an event, each broker 146 147 receiving an event performs just enough matching steps to determine which of its neighbors should receive it. As shown in Figure ??, neighbors may be brokers or clients (this figure shows a spanning tree derived from the actual 148 non-tree broker network). That is, each broker, rather than determining which subset of all subscribers is to 149 receive the event, instead computes which subset of its neighbors is to receive the event, i.e., it determines those 150 links along which it should transmit the message. Intuitively, this approach should be more efficient because the 151 number of links out of a broker is typically much less than the total number of subscribers in the system. 152

To perform link matching, we use the parallel search tree (PST) structure of the previous section, where each path from root to leaf represents a subscription. We augment the PST with vectors of trits, where the value of each trit is either -Yes,? (Y) -N o,? (N) or -M aybe? (M). We begin by annotating leaf nodes in the PST with a trit vector of size equal to the number of links out of that broker.

For each link out of a broker, a position in a trit vector determines whether to send matched events down that link, based on whether there exists a subscription reachable via that link. Leaf annotations are then propagated to non-leaf nodes in a bottom-up manner. A -Yes? in a trit annotation means that (based on the tests performed so far) the event will be matched by some subscriber that is best reached by sending the message along the given link; -No? means that the event will definitely not be matched by any subscriber along that link; and -M aybe? means that further searching must take place to determine whether or not there is such a subscriber. Annotations are described in more detail below.

The link matching algorithm consists of the following three steps. First, at each broker, the parallel search tree is annotated with a trit vector encoding link routing information for the subscriptions in the broker network. 166 Second, an initialization mask of trits must be computed at each broker for each spanning tree used for message 167 routing.

(Collectively, the masks for a single spanning tree across all the brokers encode the spanning tree in the network.) Third, at match time the initialization mask for a given spanning tree (based on the publisher) is refined until the broker can determine whether or not to send a message on each link, that is, until all values in the mask are either Yes or No. These three steps are described in detail in the following three subsections respectively.

¹⁷³ 5 A. Annotating the PST

Each broker in the network has a copy of all the subscriptions, organized into a PST as discussed in the previous 174 section, and illustrated in Figure ??. Note that the approach we describe here for computing tree annotations 175 is limited to trees with only equality tests and don't care branches. A more general solution requires the use of 176 a parallel search graph and is not described here to conserve space. Each broker annotates each node of the PST 177 with a trit vector annotation. This annotation vector contains m trits, one per outgoing link from the broker. As 178 mentioned earlier, the trit is Yes when a search reaching that node is guaranteed to match a subscriber reachable 179 through that link, No when a search reaching that node will have no subsearch leading to a subscriber reachable 180 through that link, and Maybe otherwise. Annotation is a recursive process starting with the leaves of the parallel 181 search tree, which represent the subscriptions. We label each leaf node trit in link position l with Y if one of that 182 leaf node's subscribers is located at a destination reached through link l, and N otherwise. After all the leaves 183 have been annotated, we propagate the annotations back toward the root of the PST using two operators: An 184 example is shown in Figure ??. 185

¹⁸⁶ 6 B. Computing The Initialization Mask

We assume that each broker knows the topology of the broker network as well as the best paths between each broker and each destination. To simplify the discussion, we ignore alternative routes for load balancing or recovery from failure and congestion. Instead, we assume that events always follow the shortest path. From this topology information, each broker constructs a routing table mapping each possible destination to the link which is the next hop along the best path to the destination. We also assume that the broker knows the set of spanning trees, only one of which will ever be used by each publisher.

In the case where the broker network is acyclic (Figure ??), computation of the spanning tree is straightforward. 193 If the broker topology is not a tree, then computing the spanning tree is more complex. However, even in this 194 case, there will be a relatively small set of different spanning trees. At worst, there will be one spanning tree for 195 196 each broker that has publisher neighbors and in most practical cases, where the broker network is -t ree-like?, there will be significantly fewer spanning trees. Using these best paths and spanning trees, each broker computes 197 the downstream destinations for each spanning tree. A destination is downstream from a broker when it is a 198 descendant of the broker on the spanning tree. Based upon the above analysis, each broker associates each unique 199 spanning tree with an initialization mask, one trit per link. The trit at link l has the value Maybe if at least 200 one of the destinations routable via l is a descendant of the broker in the spanning tree; and No if none of the 201 destinations routable via l are descendants of the broker 1. The significance of the mask is that an event arriving 202 at a broker should only be propagated along those links leading to descendant destinations -that is, those links 203 whose mask bit is M and will eventually be refined to a Y via matching, described below. 204

²⁰⁵ 7 C. Matching Events

When an event originating at a publisher is received at a broker, the following steps are taken using the annotated search tree:

i. A mask is created and initialized to the initialization mask associated with the publisher's spanning tree. ii.
 Starting with the root node, the mask is refined using the trit vector annotation at the current node. During
 refinement, any M in the mask is replaced by the corresponding trit vector annotation. If the mask is now fully

refined —that is, it has no M trits —then the search terminates, returning the refined mask. Otherwise, step 3
is executed. iii.
The designated test is performed and, 0, 1, or 2 children are found for continuing the search as mentioned in

The designated test is performed and, 0, 1, or 2 children are found for continuing the search as mentioned in Section 2. A subsearch is executed at each such child using a copy of the current mask.

On the return of each subsearch, all Maybe trits in the current mask for which a Yes trit exists in the subsearch mask, are converted to Yes trits. After all the children have been searched, the remaining Maybe trits in the current mask are made No trits. The current mask is returned. iv.

The top-level search terminates and sends a copy of the event to all links corresponding to Yes trits in the returned mask.

220 This concludes the description of the link matching algorithm.

221 IV.

222 8 IMPLEMENTATION AND PERFORMANCE

We have implemented the matching algorithms described above and tested them on a simulated network topology as well as on a real LAN, as explained in the following two subsections respectively.

225 9 A. Simulation Results

226 The goals of our simulations were twofold i.

To measure the network loading characteristics of the link matching protocol and compare it to that of the flooding protocol. ii.

229 To measure the processing time taken by the link matching algorithm at individual broker nodes and compare

230 it to that of centralized matching (i.e., the non-trit matching algorithm described in Section 2).

²³¹ 10 i. Simulation Setup

The simulated broker network topology is shown in Figure ??. The topology has 39 brokers and 10 subscribing 232 clients per broker, each client with potentially multiple subscriptions. In addition, there is an unspecified number 233 of publishing clients –three of these publishers, shown as P1, P2, and P3 in the figure, publish events that are 234 tracked by the simulator and the rest simply load the brokers by publishing messages that take up CPU time 235 at the brokers. As shown in Figure ??, the 39 brokers form three trees of 13 brokers each. The root of each of 236 these three trees are connected to the roots of the other two. Also, as shown, there are a small number of lateral 237 links between non-root nodes in the trees to allow messages from some publishers to follow a different path than 238 other publishers. This topology is intended to model a real-world wide-area network with each of the three rooted 239 from each other (intercontinental), but the trees distributed far 240 brokers within a tree closer to each other (interstate). The top-level brokers are modeled to have a one-way hop 241 delay of about 65 ms, links from them to their next level neighbors is 25ms, the third level hop delay is about 242 10ms, and the hop delay to clients is 1ms. 243

The broker network simulates an information space with several control parameters, such as the number of attributes in the event schema, the number of values per attribute and the number of factoring levels (i.e., the preferred attributes of Section 2.1). Subscriptions are generated randomly, but one of the control parameters is the probability that each attribute is a * (i.e., don't care). For non-* attributes, the values are generated according to a zipf distribution. In addition, we simulate -l ocality of interest? in subscriptions by having subscribers within each subtree of the broker topology have similar distributions of interested values whereas subscriptions across from the other two subtrees have different distributions.

Events are also generated randomly, with attribute values in a zipf distribution. Events arrive at the publishing brokers according to a Poisson distribution. The mean arrival rate of published events, which is a key parameter, is controlled by a user specified parameter.

In the simulation, time is measured in -t icks? of a virtual clock, with each tick corresponding to about 12 microseconds. The virtual clock, used only for simulation purposes, is implemented as synchronized brokers' clocks.

Each event carries with it its -current? virtual time from the beginning of the simulation. An event spends time traversing a link (-hop delay?), waiting at an incoming broker queue, getting matched, and being sent (software latency of the communication stack).

²⁶¹ 11 Network Loading Results

As mentioned earlier, the purpose of this simulation run was to determine, for the link matching and the flooding 262 protocols, the event publish rate at which the broker network becomes -ov erloaded? (or congested), for a varying 263 number of subscriptions. A broker is overloaded when its input message queue is growing at a rate higher than 264 the broker processor can handle. This simulation run was performed with the following parameters. The event 265 schema has 10 attributes (with 2 attributes used for factoring), and each attribute has 5 values. The subscriptions 266 are generated randomly in such a way that the first attribute is non-* with probability 0.98, and this probability 267 decreases at the rate of 85% as we go from the first to the last attribute. This means that subscriptions are very 268 269 selective –on average, each event matches only about 0.1% of subscriptions. The number of events published is 270 500.

The results from the simulation run are shown in Chart 1.

The chart shows that a broker network running the flooding protocol saturates at significantly lower event publish rates than the link matching protocol for any number of subscriptions. In particular, when each event is destined to only a small percentage of all clients, link matching dramatically outperforms flooding. In the case where events are distributed quite widely, the difference is not as great, since most links are used to distribute events in the link matching protocol. This result illustrates that link matching is well-suited to the type of selective multicast that is typical of pub/sub systems deployed on a WAN.

²⁷⁸ 12 iii. Matching Time Results

As mentioned earlier, the purpose of this simulation run was to measure the cumulative processing time taken by 279 the link matching algorithm and the centralized (non-trit) matching algorithm. The processing time taken per 280 event in the link matching algorithm is the sum of the times for all the partial matches at intermediate brokers 281 along the way from publisher to subscriber. This simulation run was performed with the following parameters. 282 The event schema has 10 attributes (with 3 attributes used for factoring), and each attribute has 3 values. The 283 subscriptions are generated randomly in such a way that the first attribute is non-* with probability 0.98, and 284 this probability decreases at the rate of 82% we go from the first to the last attribute. Again, this means that 285 subscriptions are very selective -on average, each event matches only about 1.3% of subscriptions. The number 286 of events published is 1000. 287

The results from the simulation run are shown in Chart 2.

For the link matching algorithm, six lines, -LM 1 hop? through -LM 6 hops?, are shown -these correspond to the number of hops an event had to traverse on its way from a publishing broker to a subscriber. On the Y axis, the chart shows the number of -m atching steps? performed on average. A matching step is the visitation of a single node in the matching tree. Although our current implementation has traded off time efficiency in favor of space efficiency, we estimate that a time efficient implementation can execute a matching step in the order of a few microseconds.

The chart shows that the cumulative matching steps for up to four hops using the link matching algorithm is 295 not more than the number of matching steps taken by the centralized algorithm. For more than four hops the link 296 matching algorithm takes more matching steps, however the link matching protocol is still a better choice over 297 the centralized algorithm because (1) the extra processing time for link matching (of the order of much less than 298 1ms) is insignificant compared to network latency (of the order of tens of ms), (2) the gain in latency to regional 299 publishers and subscribers obtained by distributing brokers is significant, and (3) for really large numbers of 300 subscribers (i.e., much beyond 10000), the slopes of the lines in Chart 2 indicate that centralized matching may 301 take more steps than link matching. 302

³⁰³ 13 B. System Prototype

We have implemented the matching algorithms in a network of broker nodes where brokers are connected using a 304 specified topology. A broker network may implement multiple information spaces by specifying an event schema 305 (one per information space) defining the type of information contained in each event. Clients subscribe to an 306 information space by first connecting to a broker node, then providing subscription information which includes 307 a predicate expression of event attributes. This section describes the implementation of such a broker node. As 308 illustrated in Figure ??, each broker node consists of a matching engine, client and broker protocols, a connection 309 manager and a transport layer. The matching engine which implements one of the matching algorithms described 310 earlier, consists of a subscription manager, and an event parser. A subscription manager receives a subscription 311 from a client, parses the subscription expression, and adds the subscription to the matching tree. An event parser 312 first parses a received event, then un-marshals it according to the pre-defined event schema. The matchine engine 313 then uses the implemented matching algorithm to get a list of subscibers interested in the un-marshaled event. 314 The broker to client protocol is implemented by the client protocol object, whereas the broker to broker 315

protocol is implemented using the broker protocol object. These protocol objects are robust enough to handle transient failures of connections by maintaining an event log per client. Once a client re-connects after a failure, the client protocol object delivers the events received while the client was dis-connected. A garbage collector periodically cleans up the log. The connection manager object maintains the connections to clients and the other brokers in the network.

The transport layer sends and receives messages to and from clients and other brokers in the network. To improve scalability, it implements an asynchronous -sen d? operation by maintaining a set of outgoing queues, one per connection.

A broker thread sends a message by en-queueing it in the appropriate queue. A pool of sending threads is responsible for monitoring these queueues for outgoing messages, and sending them to destinations using the underlying network protocol.

Currently, broker nodes are implemented in Java using TCP/IP as the network protocol. In an experimental setup where a 200 MHz pentium pro PC is used as a broker node, and low end PCs (using 133 MHz pentium processors) are used as clients connected using a 16MB token ring network, the current implementation of the broker can deliver upto 14,000 events/sec. Also, as shown in Chart 3 for the pure matching algorithm, brokers can perform matching very quickly, at the rate of about 4ms for 25,000 subscribers. In fact, our matching algorithms are so efficient that the transport system and network costs of a broker outweigh the cost of matching at a broker. V.

334 14 RELATED WORK

As mentioned earlier, alternatives to the link matching approach were either to (1) first compute a destination list for events by matching at or near the publisher and then distributing the event using the distribution list, or (2) to multicast the event to all subscribers which would then filter the event themselves.

Computing a destination list is a good approach for small systems involving only a few subscribers. For these 338 cases, the matching algorithm presented in section 2 provides a good solution. However, scalability is essential if 339 contentbased systems are to fill the same infrastructure requirements as subject-based publish/subscribe systems. 340 341 In cases where destination lists may grow to include hundreds or thousands of destinations, the match-first approach becomes impractical. Multicasting an event and then filtering also has its disadvantages. Lack of 342 scalability and an inability to exploit locality was shown for the flooding approach for event distribution. Flooding 343 is a good approximation of the broadcast approach since most WAN multicast techniques require the use of a 344 series of routers or bridges connecting LAN links. IP multicast [5][1] allows subscriptions to a subrange of 345 possible IP addresses known as class D addresses. Subscriptions to these groups is propagated back through 346 the network routers implementing IP. Pragmatic General Multicast [16] has been proposed as an internetwide 347 multicast protocol with a higher level of service. This protocol is an extension of IP multicast that provides 348 -TC P-like? reliability, and therefore is also reliant on multicastenabled routers. A mechanism for multicast 349 in a network of bridge-connected LANs is proposed in [7]. In this approach, members of a group periodically 350 broadcast to an all-bridge group their membership in a multicast group. Bridges note these messages and update 351 entries in a multicast table, including an expiration time. 352

The content-based subscription systems that have been developed do not yet address wide-area, scaleable 353 354 event distribution, i.e. although they are content-based subscription systems, they are not content-based routing 355 systems. SIENA allows content-based subscriptions to a distributed network of event servers (brokers) [6]. SIENA 356 filters events before forwarding them on to servers or clients. However, a scaleable matching algorithm for use at each server has not been developed. The Elvin system [14] uses an approach similar to that used in SIENA. 357 Publishers are informed of subscriptions so that they may -qu ench? events (not generate events) for which there 358 are no subscribers. In [14], plans are discussed for optimizing Elvin event matching by integrating an algorithm 359 similar to the parallel search tree. This algorithm, presented in [8], converts subscriptions into a deterministic 360 finite automata for matching. However, no plans for optimizations for broker links (such as our optimization 361 through trit annotation) are discussed. Another algorithm for optimizing matching is discussed in [9]. At analysis 362 time, one of the tests aij of each subscription is chosen as the gating test; the remaining tests of the subscription 363 (if any) are residual tests. At matching time, each of the attributes aj in the event being matched is examined. 364 The event value v_j is used to select those subscriptions i whose gating tests include $a_{ij} = v_j$. The residual tests 365 of each selected subscription are then evaluated: if any residual test fails, the subscription is not matched; if 366 367 all residual tests succeed, the subscription is matched. Our parallel search tree performs this type of test for 368 each attribute, not just a single gating test attribute. One outlet for the work presented in this paper could be through Active Networks [17]. Active Networks have been touted as a mechanism for eliminating the strong 369 dependence of route architectures on Internet standards. Active Networks allow the dynamic inclusion of code 370 either at routers or by replacing passive packets with active code. The SwitchWare project [3] follows the former 371 approach and is most appropriate to the type of router customization proposed in this paper. With SwitchWare, 372

373 digitally signed type-checked modules may be loaded into network routers.

374 Our matching and multicasting component could be one such module.

375 15 VI. CONCLUSIONS

In this paper, we have presented a new multicast technique for content-based publish/subscribe systems known 376 as link matching. Although several publish/subscribe systems have begun to support content-based subscription, 377 the novel contribution of link matching is that routing is based on a hop-by-hop partial matching of published 378 events. The link matching approach allows distribution of events to a large number of information consumers 379 distributed across a WAN without placing an undo load on the network. The approach also exploits locality of 380 subscriptions. We evaluate how an implementation of content-based routing protocol performs by showing that a 381 broker network stays up while running the link matching algorithm whereas brokers get overloaded for the same 382 event arrival rate running the flooding algorithm, since brokers have larger numbers of events to process in the 383 flooding case. We also describe a broker implementation that can handle message loads of up to 14000 events per 384 second on a 200 MHz Pentium PC. This shows that content-based routing using link matching supports a more 385 general and flexible form of publish-subscribe while admitting a highly efficient implementation. Future work is 386 concentrating on further validation of our approach to content-based routing. We are currently working to deploy 387 our content-based routing brokers on a large private network. This will allow us to conduct system tests under 388 actual application loads. Sample applications will include some from the financial trading and process control 389 domains. In addition to these system tests, we are also continuing work with our simulator to examine different 390 types of messaging loads. In particular, since many publish/subscribe applications exhibit peak activity periods, 391 we are examining how our protocol performs with bursty message loads. VII. 1 392

¹In some cases, where some destinations reachable through a link downstream on some spanning trees and are not on others, the search may be optimized by splitting the link into two or more "virtual" links.



Figure 1:



Figure 2:



Figure 3:

15 VI. CONCLUSIONS

- [Hanson et al. ()] 'A predicate Matching Algorithm for Database Rule Systems'. Eric N Hanson , Moez Chaabouni
 , Chang-Ho Kim , Yu-Wang Wang . SIGMOD 1990, May 23-25 1990. p. .
- [Tennenhouse et al. (1997)] A Survey of Active Network Research,? IEEE Communications Magazine, J
 Tennenhouse, W D Smith, D Sincoskie, G Wetherall, Minden. January, 1997. 35 p. .
- Birman (1993)] K P Birman . T he process group approach to reliable distributed computing,? pages, Dec. 1993.
 36 p. .

399 [Mishra et al. (1991)] Consul: A Communication Substrate for Fault-Tolerant Distributed Programs, Dept. of

- *computer science*, Shivakant Mishra , Larry L Peterson , Richard D Schlichting . TR 91-32. Nov. 1991. The
 University of Arizona
- 402 [CORBA services: Common Object Service Specification (1998)] CORBA services: Common Object Service
 403 Specification, July 1998. Object Management Group ; Object Management Group (Technical report)
- [Aguilar ()] 'Datagram Routing for Internet Multicasting'. Lorenzo Aguilar . ? ACM Computer Communications
 Review 1984. 14 (2) p. .
- 406 [Deering ()] Stephen E Deering . M ulticast Routing in InterNetworks and Extended LANs, 1988. 18 p. .
- 407 [Carzaniga et al.] Design of a Scalable Event Notification Service: Interface and Architecture, Antonio Carzaniga
 408 , David S Rosenblum , Alexander L Wolf . http://www.cs.colorado.edu/users/carzaniga/siena/
 409 index.html
- [Segall and Arnold (1997)] El vin has left the building: A publish/subscribe notification service with quenching,?
 Proceedings of AUUG97, Bill Segall, David Arnold. September, 1997. Brisbane, Australia.
- [Gough and Smith ()] John Gough , Glenn Smith . Eff icient Recognition of Events in a Distributed System,?
 Proceedings of ACSC-18, (Adelaide, Australia) 1995.
- 414 [Powell Guest Editor (1996)] Group Communication?, David Powell, (Guest Editor (ed.) April 1996. 39 p. .
- [Aguilera et al. ()] 'Matching Events in a Content-Based Subscription System'. Marcos Aguilera , Rob Strom ,
 Daniel Sturman , Mark Astley , Tushar Chandra . http://www.research.ibm.com/gryphon Upcoming
- IBM Technical Report 1998.
- [Oki et al. (1993)] Brian Oki , Manfred Pfluegl , Alex Siegel , Dale Skeen . The Information Bus An Architecture for Extensible Distributed Systems,? pages, Dec. 1993. 27 p. .
- 420 [Speakman et al. (1998)] Tony Speakman , Dino Farinacci , Steven Lin , Alex Tweedly . PGM Reliable Transport
 421 Protocol,? IETF Internet Draft, August 24. 1998.
- [Alexander (1998)] The SwitchWare Active Network Architecture,? IEEE Network Special Issue on Active and
 Controllable Networks, Scott Alexander . July 1998. 12 p. .
- 424 [Uyless Black. TCP/IP Related Protocols ()] Uyless Black. TCP/IP & Related Protocols, 1995. McGraw-Hill.
 425 p. . (Second Edition)
- 426 [Skeen] Vitria's Publish-Subscribe Architecture: Publish-Subscribe Overview, Dale Skeen . http://www. 427 vitria.com/