

Available online at www.sciencedirect.com

Computer Communications xxx (2006) xxx–xxx

computer
 communications

www.elsevier.com/locate/comcom

Explicit routing in multicast overlay networks

Torsten Braun ^{a,*}, Vijay Arya ^b, Thierry Turletti ^b

^a *University of Bern, Neubrückstrasse 10, CH-3012 Bern, Switzerland*

^b *INRIA, 2004 route des Lucioles, B.P. 93, F-06902 Sophia Antipolis*

Received 10 February 2005; received in revised form 16 February 2006; accepted 21 February 2006

8 Abstract

9 Application Level Multicast is a promising approach to overcome the deployment problems of IP level multicast by establishing delivery trees using overlay links among end systems. This paper presents algorithms to support traffic engineering, to improve the reliability of multicast delivery, and to facilitate secure group communications. First, we introduce the so-called backup multicast tree algorithm to compute a set of $n - 1$ backup multicast delivery trees from the default multicast tree. Each backup multicast tree has exactly one link of the default multicast tree that is replaced by a backup link from the set of available links. The algorithm can calculate this set of trees with a complexity of $O(m \log n)$, which is identical with the complexity of well known minimum spanning tree algorithms. The so-called reduced multicast tree algorithm is based on the backup multicast tree algorithm and can calculate a tree from the default multicast tree by removing a particular node and by replacing the links of the removed node. Using the algorithms trees can be calculated individually by each of the nodes but it requires global topology knowledge. We therefore discuss distributed versions of the algorithms.

18 © 2006 Published by Elsevier B.V.

19 *Keywords:* Multicast; Overlay networks; Explicit routing

21 1. Introduction

22 Application level multicast also known as end system
 23 multicast or overlay multicast has become a very popular
 24 research topic during the last years due to the deployment
 25 problems of IP multicast. Typically, application level mul-
 26 ticast approaches apply similar concepts as IP multicast
 27 such as running multicast routing protocols and building
 28 multicast delivery trees, but with the difference that these
 29 operations are performed on application rather than on
 30 network level. Application level multicast avoids multicast
 31 deployment problems in the Internet and can be used to
 32 bypass routes established by underlying routing protocols
 33 that do not consider the current load or congestion level
 34 for the routing decision. Application level multicast is
 35 based on the establishment of overlay networks. Multicast

packets are forwarded between the end systems via such
 overlay networks.

Mechanisms for explicit path selection are not included in most multicast distribution concepts. With explicit path selection, the sender of a multicast packet can explicitly select the distribution path (usually a tree) of a single multicast packet. This allows a sender selecting individual multicast trees for each single packet in order to react on events such as link breaks, node failures, congested links, and group member leaves. We propose that a sender of a multicast packet can select a backup multicast tree instead of the default multicast tree by inserting a fixed size identifier to the multicast packet. A multicast delivery tree is typically established by multicast routing protocols in case of IP multicast and by peer-to-peer protocols in case of application level multicast. Such a multicast delivery tree is then used for the distribution of multicast data. The selected backup multicast tree can then be used to immediately react on link failures without any delay caused by reestablishing a new multicast delivery tree for the new topology. Load balancing can be achieved by using different trees simultaneously and

* Corresponding author. Tel.: +41 31 631 4994; fax: +41 31 631 3261.

E-mail addresses: braun@iam.unibe.ch (T. Braun), Vijay.Arya@sophia.inria.fr (V. Arya), Thierry.Turletti@sophia.inria.fr (T. Turletti).

57 can be applied when a particular link of the default multicast
58 tree becomes congested or for increasing throughput.

59 Another usage of explicit path selection might be prevent-
60 ing particular nodes of a multicast group to receive a multi-
61 cast message. This might be useful in secure group
62 communications. In such cases, group keys must be updated
63 whenever a node joins or leaves [1]. For joining nodes, the
64 new key can be distributed along the old multicast delivery
65 tree and by explicitly sending the new key to the joining mem-
66 ber. However, in case of a leaving node, the old multicast
67 delivery tree can not be used, because the leaving node would
68 receive the multicast message with the new group key. A
69 naïve approach for key distribution is to distribute the group
70 key via point-to-point connections between the root generat-
71 ing the new key and the individual group members, but this
72 approach is not scalable. Other tree based and hierarchical
73 approaches form sub-groups within the group and distribute
74 the new group key along the sub-group trees. Only a small
75 part of the sub-groups must be re-established for a joining
76 member and most of the established multicast distribution
77 trees for the various sub-groups can be used for efficient
78 key distribution [1,2].

79 In Section 2 we review related work on explicit path selec-
80 tion and on application level multicast, in particular con-
81 cepts to improve reliability. Section 3 introduces our
82 concept of explicit path selection in multicast overlay net-
83 works and presents an appropriate signaling protocol. The
84 proposed algorithm for constructing $n - 1$ backup trees
85 out of a single default multicast tree is described in Section
86 4. Also the complexity and performance of this algorithm
87 is evaluated. We will show that the complexity for calculating
88 $n - 1$ backup multicast trees for each of the links of the
89 default multicast tree has a similar complexity as calculating
90 a single minimum spanning tree. The result is confirmed by
91 an implementation of the backup multicast tree algorithm.
92 This algorithm is one of the key contributions of this paper.
93 Section 5 presents an algorithm for constructing a multicast
94 delivery tree, but with the condition that a particular node is
95 removed from the default multicast tree. This algorithm is
96 based on that one introduced in Section 4 and can be used
97 to calculate additional n alternative multicast delivery trees.
98 Section 6 presents our proposed encoding scheme to specify
99 within a multicast packet, which of the alternative multicast
100 delivery trees shall be used for multicasting the packet. The
101 encoding scheme is based on a cardinal representation of
102 trees and is another main contribution of the paper. Section
103 7 presents distributed versions of the algorithms presented in
104 Sections 4 and 5. The signaling protocol introduced in Sec-
105 tion 3 is extended accordingly. Section 8 concludes the paper.

106 2. Related work

107 2.1. Explicit path selection

108 Explicit path selection can be implemented by explicitly
109 specifying the nodes to be traversed, e.g., using the routing
110 header in IPv6 [3] or by describing the multicast group

member addresses in explicit multicast [4]. Since specifying
all the nodes to be traversed does not scale for large multi-
cast groups, it has been proposed for small groups only. As
an alternative, packets can be marked with a unique path
identification (ID) such as a label like in multi-protocol
label switching (MPLS) [5]. The label must be assigned
with a certain path using label distribution protocols,
which adds significant overhead in terms of signaling band-
width and delay. MPLS-like approaches add hard states to
the involved protocol entities.

In contrast to MPLS, the BANANAS concept [6] pro-
vides path IDs for IP level unicast forwarding without intro-
ducing a special signaling protocol. The path ID is derived
from a link state database, which must be known in advance
within a routing domain, and is encoded as a concatenation
of local link IDs of the routers to be traversed. Four bits are
sufficient per router with up to 15 interfaces. In this case, a
128 bit path ID can encode paths with a length of up to 32
hops. In order to eliminate the signaling overhead, BANAN-
AS proposes to calculate alternative paths in a distributed
manner such that each node calculates the same set of paths.
The concept is based on the assumption that a limited set of
possible alternative paths can be calculated in reasonable
time. It is proposed that each node i calculates the k_i shortest
paths to each destination. Since the calculation is performed
at each node independently from each other, a distributed
validation process in order to harmonize the calculations is
required. According to [7], k shortest paths of a graph with
 n vertices and m edges can be calculated with complexity
 $O(m + n \log n + k)$.

2.2. Application level multicast

Several Application level multicast schemes have been
proposed by other researchers. Some of them provide
mechanisms to support load balancing or reliability in case
of error situations such as node failures and broken links.
However, none of them supports the delivery of multicast
packets to exactly $n - 1$ group members, which can be
helpful for efficient key distribution in multicast groups.

Scribe [8] is built on Pastry [9], a generic peer-to-peer
object location and routing substrate. Scribe generates a
tree routed at a rendezvous point, which corresponds to
the node with the closest node ID to the group ID of the
multicast group. A node can join the group by sending a
join message towards the root of the tree. Forwarding
entries are created in forwarder nodes as a result of join
messages. Nodes wishing to leave a group transmit leave
messages, which result in removing the forwarding entries
in the forwarder nodes. In case of a parent node or link
failure, a node must retransmit a join message towards
the tree in order to repair its branch.

Bayeux [10] is an application level multicast system on
top of Tapestry [11] routing. Multicast receivers are orga-
nized in a tree with a single root. Load balancing can be
achieved by replicating root nodes. A member joins by
sending a join message towards the root of the tree.

166 The root replies with a tree message. When receiving a tree
167 message, nodes on the path between the root and the join-
168 ing member add the new member as a node they have to
169 serve. Similarly, leave messages trigger prune messages in
170 case a member leaves a group.

171 NICE [12] arranges end systems into a hierarchy. Each
172 end system is assigned to a certain level in the hierarchy.
173 Members of the same level are grouped into different clus-
174 ters, which are controlled by a cluster leader. The cluster
175 leader is selected such that it has a minimum distance to
176 the other cluster members. The hierarchy is used to define
177 different structures for control message and data delivery.
178 Control messages are required for cluster management. A
179 joining host is mapped to a cluster, such that it has rather
180 close neighbors. Each member of a cluster on level n must
181 be the head of a cluster on level $n - 1$. Clusters exceeding a
182 certain size are split or merged if the cluster size violates
183 some upper or lower bounds.

184 A mechanism called Probabilistic Resilient Multicast
185 [13] is based on forwarding data not just once but multiple
186 times to randomly selected nodes. Obviously, this introduc-
187 es some bandwidth overhead. Negative acknowledgements
188 are used to detect losses.

189 Narada [14,15] constructs a multicast distribution tree in
190 two steps: first, a mesh is established out of a set of possible
191 links between two nodes based on continuous performance
192 measurements between two nodes. This mesh is then the
193 basis for the spanning tree construction in the second step
194 by applying a reverse shortest path mechanism. Spanning
195 trees are constructed for each potential sender in order to
196 optimize trees for each source. In case of node or link fail-
197 ures, new overlay links need to be added to the mesh. This
198 results in some delay to repair a failure.

199 The HostCast protocol [16] establishes an overlay data
200 delivery tree and a corresponding control mesh. Both cover
201 all group members. Reliability is achieved by establishing
202 specific secondary links between nodes and their grandpar-
203 ent nodes as well as their uncle nodes. Broken links of the pri-
204 mary data delivery tree can then be replaced by the secondary
205 links. HostCast requires establishing a control mesh. More-
206 over, the number of secondary links is quite large ($>2n$).

207 Network Coding and distribution of specially encoded
208 multicast streams over a redundant multicast graph has
209 been proposed in [17] The concept not only increases the
210 throughput that can be achieved by distributing data to
211 receivers that are reachable via various paths. If different
212 streams have to pass the same links, they can be combined
213 using network coding mechanisms. Each node has two
214 redundant paths to the source, but this does not protect
215 from arbitrary link failures.

216 3. Explicit routing in multicast overlay networks

217 3.1. Overview

218 The concept of explicit routing in multicast overlay net-
219 works as introduced in this paper proposes path IDs for

multicast data delivery. This allows selecting a certain mul- 220
ticast tree for explicit delivery of a multicast packet. We 221
propose to select one multicast delivery tree for application 222
level multicast packet distribution from a set of up to n 223
trees, where up to $n - 1$ backup multicast trees are con- 224
structed from the default multicast delivery tree. Each of 225
the $n - 1$ backup multicast trees has $n - 2$ links in common 226
with the default multicast tree and differs in exactly one 227
link. A sender of a packet can then choose among the n 228
trees to distribute the packet. The chosen tree must be iden- 229
tified by an ID within each multicast message. Since we 230
assume a limited ID space, we have to limit the set of pos- 231
sible trees among which a sender can choose. We present 232
the backup multicast tree algorithm that calculates the 233
 $n - 1$ backup multicast trees belonging to a single default 234
multicast delivery tree. Based on this algorithm we present 235
the reduced multicast tree algorithm computing multicast 236
delivery trees that include all nodes of a group except a sin- 237
gle particular node to be removed from the multicast 238
group. For each node to be removed its links will be 239
replaced, if possible by links calculated for the $n - 1$ back- 240
up multicast trees mentioned above. 241

242 Alternatively, a tree that is as disjoint as possible to the
243 default multicast tree could be calculated. However, such a
244 single disjoint tree can not be used for the construction of
245 trees that exclude particular nodes. Moreover, our algo-
246 rithm for the backup multicast trees minimizes the number
247 of backup links that are required to build the $n - 1$ backup
248 multicast trees.

249 Since we believe that application level multicast will be a
250 promising basis for future multicast services and applica-
251 tions, we develop our concept within this context. For
252 our work, we assume some kind of overlay network such
253 as CAN [18], Chord [19], RON [20] or Tapestry [11] on
254 top of which the application level multicast protocol can
255 run. Our concept is independent of the underlying proto-
256 cols. We only assume that the underlying protocols estab-
257 lish a mesh of links between nodes, not necessarily a full
258 mesh. A certain connectivity is also required to be able to
259 identify backup links that shall replace the links of the
260 default multicast tree in certain conditions.

261 Our concept is independent from specific applications
262 and could support applications such as streaming,
263 audio/video conferencing, games and computer-supported
264 collaborative work. It is based on the calculation of a
265 spanning tree for multicast data delivery. Spanning trees
266 can be used for both any source multicast and source-spe-
267 cific multicast. Scalability is limited by the overhead to
268 distribute topology information and to compute a span-
269 ning tree based on this information. However, by introduc-
270 ing hierarchical structures as proposed in NICE [12],
271 the concept should be able to support large numbers of
272 group members. We therefore think that our proposed
273 mechanisms could be integrated into NICE. However,
274 the concept should also be applicable to other protocols
275 that are generating spanning trees for multicast data
276 delivery, e.g. Narada [14].

277 3.2. Signaling support

278 The mechanism for computing backup multicast trees
279 can be used in three modes:

- 280 1. **Independent mode.** Each node must get a complete view
281 of the multicast overlay topology, in order to calculate
282 the backup multicast trees independently from any other
283 node. This is a similar requirement as in [6], where each
284 node needs the complete knowledge of a domain's topol-
285 ogy. In our case, each multicast overlay node performs
286 exactly the same algorithm and computes exactly the
287 same set of backup multicast trees for a given default
288 multicast tree. The algorithm is described in Section 4.
289 In order to get this complete overlay topology view,
290 the exchange of topology information is required. In this
291 section, we propose a simple signaling protocol that sup-
292 ports the distribution of topology information.
- 293 2. **Distributed mode.** The exchange of complete topology
294 information can be avoided by a distributed version of
295 the proposed mechanism. This allows reducing the
296 amount of exchanged information. It also allows a node
297 to know only the local neighborhood, but it requires a
298 sophisticated signaling protocol, which is tailored to
299 support the backup multicast tree algorithm. This prot-
300 ocol is described in Section 7 in more detail.
- 301 3. **Central mode.** The algorithm can also be used at the root
302 of the multicast delivery tree only. In this case, only the
303 root calculates an appropriate tree for multicast data
304 delivery and specifies the tree using a self-describing
305 specification of the multicast tree in the multicast data
306 packet. This can be done using a cardinal representation
307 as described in Section 6. However, such a cardinal rep-
308 resentation might exceed the space available for a tree
309 description and might be applicable to limited group siz-
310 es only.

312 An important task to enable multicast data distribution
313 is the management of a multicast group and the multicast
314 delivery tree establishment. We propose to use a simple
315 protocol in order to exchange complete topology informa-
316 tion among the multicast overlay network and to support
317 the independent mode (1.) as described at the beginning
318 of this sub-section. The signaling protocol makes use of
319 three signaling messages: join, leave and tree. Those signal-
320 ing messages can be encrypted or signed depending on
321 security requirements.

322 The **join** message is sent by any node that wants to join
323 the multicast group. The join message contains information
324 about the connectivity of the new member to other peers
325 and is forwarded towards the root of the multicast tree.
326 In order to limit the join implosion problem in case of a
327 large group, we can have multiple root nodes that individ-
328 ually serve a certain subset of multicast member nodes. In
329 that case, these root nodes have to organize themselves on
330 a higher level such that each root node gets all multicast
331 packets sent to the multicast group. Naturally, the root

nodes of the various sub-trees perform the same protocol 332
but just one level higher. This two-tier architecture also 333
corresponds to peer-to-peer networks with *super peers*. 334
Super peers are peers with special characteristics such as 335
higher access network bandwidth or higher processing 336
power. They are ideal candidates to serve as root nodes 337
for sub-trees. Such a two-tier architecture as depicted 338
in Fig. 1 can also preserve the scalability of the approach 339
in case of large groups. Note that similar as proposed in 340
[12] more than two levels can be formed. In that case, a 341
node must only know the topology of its own sub-tree. 342

In response to a join message, the root sends a **tree** mes- 343
sage to the group members possibly after the root has 344
checked whether the node is allowed to join the group. 345
The purpose of the tree message is to inform the other 346
group members about newly joined nodes and to update 347
the connectivity information of that peers. We assume that 348
the exchange of tree messages ensures that the peers are 349
always aware of the connectivity within the multicast over- 350
lay network. If the overlay network does not provide infor- 351
mation about the quality of the links, the nodes might 352
measure parameters like round trip times or available 353
bandwidth to other nodes themselves. In case of a super 354
peer based network, each peer only needs to know the con- 355
nectivity of the peers belonging to the same sub-tree served 356
by a super peer. After receiving the tree message each node 357
updates its information about the peer-to-peer network. It 358
also calculates the alternative multicast delivery trees, i.e. 359
the default multicast tree, e.g. using a minimum spanning 360
tree algorithm, and the backup multicast trees using the 361
backup multicast tree algorithm presented in Section 4. 362
The alternative multicast delivery trees are assigned to 363
some unique tree ID as discussed in Section 6. Each multi- 364
cast message must carry this tree ID. A node receiving a 365
message can derive from the tree ID and the knowledge 366
about the topology how to forward the message. 367

If a peer node wants to leave a group, it sends a **leave** 368
message towards the root. In this case, the root updates 369
its member list as well as the overlay network topology 370
and sends a tree message to the group in order to update 371

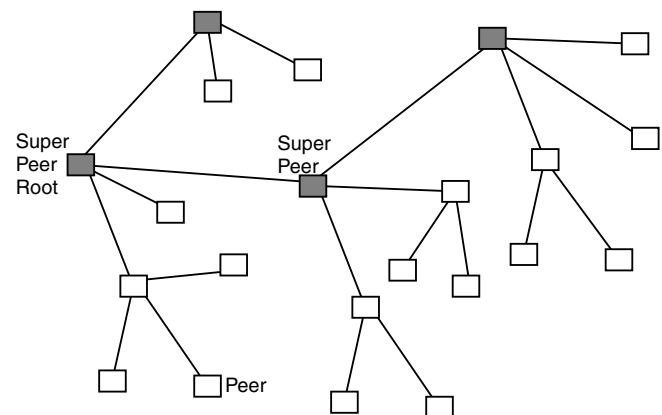


Fig. 1. Two-tier application level multicast architecture.

372 the group membership and topology information. All
373 group members have to update their group membership
374 and topology information accordingly and have to recalcu-
375 late the alternative multicast delivery trees including their
376 tree IDs.

377 Tree messages are sent in response to join or leave mes-
378 sages, but should also be sent periodically (typically in the
379 range of several seconds or tens of seconds as usual in link
380 state routing protocols) in order to refresh group member-
381 ship and topology information. In case of a refresh the tree
382 messages can use an already existing multicast delivery tree
383 and can be sent by the root with the corresponding tree ID.
384 The same is true in the case of a leave message. The leaving
385 node should take this tree message as an acknowledgement
386 that the root has received the leave message. For the other
387 nodes, the tree message serves as a message to indicate that
388 a node has left the group. Recalculation of the alternative
389 multicast delivery trees and updating the tree IDs should
390 occur after forwarding the tree message to the next node.
391 If the tree message is sent in response to a join message,
392 it can be sent using a previously used tree ID identifying
393 the multicast delivery tree without the new node. In addi-
394 tion, a separate tree message with complete group member-
395 ship and topology information can be unicast to the new
396 peer.

397 The transmission of a multicast message can be per-
398 formed always via the root node in order to allow admis-
399 sion control for group messages. If an incoming message
400 contains an unknown tree ID, the message should be for-
401 warding as a broadcast message to the neighbor nodes of
402 the peer. Broadcast messages should be kept in memory
403 for a certain duration in order to detect and discard dupli-
404 cated broadcast messages.

405 4. Backup multicast trees

406 4.1. Overview

407 The main motivation of explicit multicast routing is to
408 enforce packet forwarding along pre-selected alternative
409 paths. Alternative paths can be used in several situations
410 such as link failures, congestion, and leaving nodes. The
411 corresponding multicast delivery trees need to be calculated
412 in advance and each node must be able to assign a tree ID.

413 We assume that n nodes (vertices) of an overlay network
414 are interconnected by a mesh of m links (edges). For a full
415 mesh with n vertices, the number of edges is $m = n(n - 1)/2$.
416 However, usually application level multicast approaches do
417 not establish full meshes, but only certain links between
418 nodes. For example, Narada [14] proposes to select the best
419 links that fulfill certain quality requirements. Moreover,
420 since overlay networks are established between nodes
421 behind firewalls, we have to assume that not each node
422 can connect arbitrarily to any other node. On the other
423 hand, in most approaches each node tries to establish a
424 certain number of links to other nodes. This also improves
425 the reliability of the overlay network.

426 Out of the finally resulting mesh a huge number of possi-
427 ble multicast delivery trees could be calculated. We
428 assume that multicast delivery trees are spanning trees con-
429 sisting of n vertices and $n - 1$ edges. Since the tree ID is
430 limited to a certain size, we need an algorithm that restricts
431 the number of possible multicast delivery trees. We propose
432 to restrict this number to n and to compute $n - 1$ backup
433 edges that can replace each of the $n - 1$ edges of the default
434 multicast tree in case of link breaks or congestion situa-
435 tions. The basic idea of our approach is compute a backup
436 edge from the set $G-T$ (G : set of edges of the graph, T : set
437 of default multicast tree edges) for each edge of the default
438 multicast tree T . Replacing each of the $n - 1$ edges of T by
439 its corresponding backup edge results in $n - 1$ backup mul-
440 ticast trees. Note that a single edge can serve as backup
441 edge for more than one default multicast tree edges. The
442 default multicast tree and the $n - 1$ backup multicast trees
443 result in n alternative multicast delivery trees that can be
444 used for a delivery over a multicast overlay network.

445 In this section, we present the algorithm for modes 1 and
446 3 as described at the beginning of Section 3.2. We assume
447 that each node knows the default multicast tree. There
448 are several ways how to build such a default multicast tree.
449 It can be built based on multicast routing protocols or by
450 using a minimum spanning tree algorithm. Assuming
451 Prim's minimum spanning tree algorithm [21] with a com-
452 plexity of $O(m \log n)$, one could naively calculate $n - 1$
453 minimum spanning trees for the $n - 1$ graphs $G - e_i$ with
454 $e_i =$ edge i of the minimum spanning tree, $i = 1, \dots, n - 1$.
455 The complexity for calculating $n - 1$ backup multicast
456 trees is $O(mn \log n)$ in this case. Our intention is not to cal-
457 culate backup multicast trees with optimal link weights,
458 but those that can be calculated at low cost. Another goal
459 is to determine a rather small set of edges that can serve as
460 backup edges for the edges of the default multicast tree.
461 Moreover, we select backup paths in such a way that they
462 have minimum overlap with the default path.

463 The idea for the proposed algorithm to calculate the
464 $n - 1$ backup multicast trees for a given default multicast
465 tree is taken from the observation that one broken edge
466 of the default multicast tree can either be repaired with a
467 single replacement of this edge by a backup edge that is
468 not in the default multicast tree or it can not be repaired
469 at all. A backup edge can repair all other edges of the
470 default multicast tree with which the backup edge forms
471 a cycle. For example, in Fig. 2 edges (C, E), (E, I), (I,
472 M), (M, N), (C, H), and (H, L) can be repaired by edge
473 (N, L) or (L, N). In this case, edges (N, L) or (L, N) gener-
474 ate a cycle with the other edges that can be repaired by
475 that edge. Also, (J, L) forms a circle with the edges (A,
476 C), (C, H), (H, L), (A, B), (B, D), (D, F), and (F, J). Again
477 (J, L) or (L, J) can repair all these edges by a single replace-
478 ment. From these examples we also see that both (C, H)
479 and (H, L) can be repaired by either (N, L) or (J, L).
480 Our algorithm selects the backup edge that has the lowest
481 match with the path to be repaired. This is motivated by
482 the goal to circumvent the edge to be replaced as far as pos-

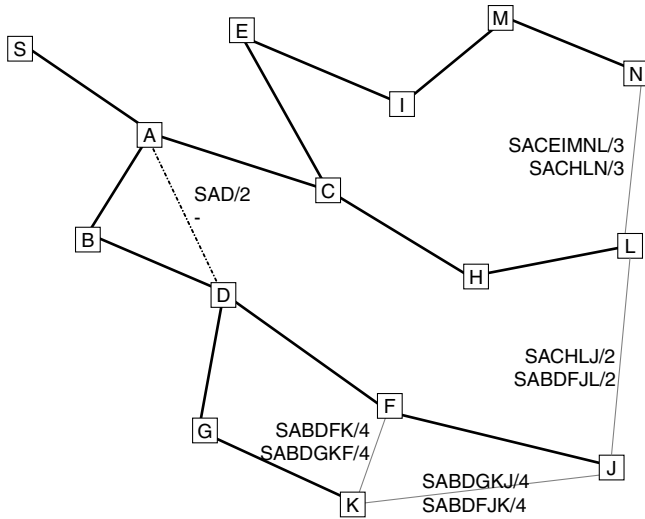


Fig. 2. Default multicast tree with backup edges.

483 sible. In our example, (H, L) will be replaced by (J, L), but
 484 not by (N, L), because (J, L) generates the path SABDFJL
 485 from S to L, while (N, L) generates the path SACEIMNL.
 486 SABDFJL has only the first two nodes in common with the
 487 original path SACHL, while SACEIMNL has three nodes
 488 in common with SACHL. Note that in Fig. 2, edge (S, A)
 489 can not be repaired.

490 4.2. Backup multicast tree algorithm

491 In the following we describe an algorithm for uniquely
 492 determining the backup edge for each single default multi-
 493 cast tree edge using a C-like pseudo code. We assume that
 494 the complete graph is stored in a linked data structure of
 495 vertices with pointers to their edges and neighbor vertices
 496 as it would result from a minimum spanning tree calcula-
 497 tion. We also assume that the default multicast tree edges
 498 are labeled as such and that each vertex has a vertex ID.

499 In the first part, we calculate the path from the root to
 500 each vertex in T and store this path with each vertex
 501 (path_from_root). We also store the distance (distan-
 502 ce_to_root) of each vertex to the root. We begin with
 503 the root as first vertex and add all vertices that can be
 504 directly reached from the root to set V. The path from
 505 the root as well as the distance is stored at all the vertices
 506 that can be directly reached from the root. After processing
 507 all direct neighbors of the root, we select all the vertices of
 508 set V after each other and perform the same operations as
 509 for the root. This is repeated until all nodes of the tree have
 510 been processed.

511 In the second part, we calculate for each edge in G–T the
 512 resulting path from the root via the first vertex to the sec-
 513 ond vertex of the edge (backup_path). That path is called
 514 backup path for the second vertex, since it allows reaching
 515 a vertex via an alternative path other than the default path
 516 along the default multicast tree. The backup paths can be
 517 used to reach vertices in case of edge failures. For example,

edge (L, N) stores backup path SACHLN (S → L, N), 518
 while (N, L) stores backup path SACEIMNL (S → N, 519
 L). Edge (A, D) stores backup path SAD (S → A, D), while 520
 edge (D, A) stores SABDA (S → D, A), which will later be 521
 detected as not valid, because A occurs twice in it. This 522
 means that a vertex can be reached via several paths: First 523
 via the path along the default multicast tree and in addition 524
 via several backup paths. For example, node L can be 525
 reached via the default multicast tree, but also via a backup 526
 path via J and another backup path via N. For all backup 527
 paths we also determine at which node the backup path 528
 and the path along the default multicast tree begin to differ. 529
 For example, node L can be reached via the backup path 530
 SACEIMNL (S → N, L) and by the path along the default 531
 multicast tree SACHL (S → L). These two paths begin to 532
 differ in the fourth vertex (E vs. H), but the first three ver- 533
 tices (SAC) are identical. Therefore, we store the value of 3 534
 (also called common path length) with the backup path 535
 (SACEIMNL/3) at (N, L). Fig. 2 shows the result of the 536
 first two parts of the algorithm. 537

538 In the third part, we copy the backup paths from the 539
 edges in G–T to the leaf edges of the default multicast 540
 delivery tree. After the copy operation we extend the back- 541
 up path by that node of the edge that is closer to the root 542
 of the tree. If a leaf edge connects to two edges of G–T, we 543
 only keep one backup path, in particular that path with 544
 the lowest common path length. For example, edge (H, 545
 L) connects to two edges in G–T: (N, L) and (J, L). We 546
 only keep backup path SABDFJLH/2 from (J, L), but do 547
 not keep SACEIMNLH/3 from (N, L) due to the lower 548
 common path length at (J, L)h. This approach selects 549
 among several alternative backup paths that one with the 550
 lowest number of common nodes shared with the path 551
 from the root to a node along the default multicast deliv- 552
 ery tree. We consider that backup path as the best choice. The 553
 inner while loop of part 3 is executed twice, once for all 554
 edges in T and once for all edges in G–T. Only the “best” 555
 backup path is copied from edge f to the considered edge e. 556
 After processing the inner while loop twice, the considered 557
 edge becomes labeled. If the upstream edge of e(g) then 558
 only has labeled downstream edges, g is added to set E. 559
 The set E contains all edges that are ready to be processed 560
 by the inner while loop. 561

Algorithm.

```

562 vertex_set N, V;
563 vertex root, n, v, x;
564 edge_set E, F, G, T;
565 edge e, f, g, h;
566 V := {root}; root.path_from_root := (root);
567 root.distance_to_root := 0;
568 while (V != ∅) /* part 1: O(n) */
569 {
570     v := first_element(V); V := V - v;
571     N := {all vertices x | edge (v, x) ∈ T};
572     while (N != ∅) {
  
```

```

573     n := first_element(N); N := N - n;
574     n.path_from_root := (v.path_from_
575     root, n);
576     n.distance_to_root := v.dis-
577     tance_to_root + 1;
578     V := V + n;
579 }
580 }
581 }
582 }
583 /* here, for all vertices x: x.path_from_
584 root */
585 /* describes the path from the root to x, */
586 /* x.distance_from_root describes the num-
587 ber of */
588 /* hops from root to x */
589 F := G-T;
590 while (F != ∅) { /* part 2: O(m log n) */
591     f := first_element(F); F := F - f;
592     f.backup_path := (f.x.path_from_root,
593     f.y);
594     x := (lowest common parent of f.x and
595     f.y);
596     /* we apply binary search: complexity
597     O(log n) */
598     f.backup_path_common := x.distance_
599     to_root + 1;
600 }
601 }
602 /* e.x, e.y are the vertices of e (e.x ->
603 e.y) */
604 /* at this point each edge e from G-T stores
605 the */
606 /* path from the root to e.x plus edge e.y */
607 E := {edges of T with leaf vertices};
608 /* leaf vertex: vertex with degree 1 */
609 while (E != ∅) { /* part 3: O(m) */
610     e := first_element(E); E := E - e;
611     e.backup_path = ();
612     e.backup_path_common = MAXINT;
613     for (j := 2; j > 0; j-){
614         if (j == 2)
615             F := {all edges f | f ∈ T && f.x == e.y};
616         else
617             F := {all edges f | f ∈ G-T && f.y ==
618             e.y};
619         while (F != ∅){
620             f := first_element(F); F := F - f;
621             if ((e.y is not twice in f.back-
622             up_path) &&
623                 (f.backup_path_common
624                 < e.backup_path_common))
625                 e.backup_path := (f.backup_path,
626                 e.x);
627                 e.backup_path_common := f.backup_
628                 path_common;
629         }
630     }
631 }
632 }
633 }
634 }
635 }
636 }
637 }
638 e.labelled := TRUE;
639 if (all edges g ∈ T with g.x == e.x

```

```

640 are labeled) {
641     h := (edge ∈ T with h.y == e.x);
642     E := E + h;
643 }
644 }
645 }
646 }
647 }

```

Fig. 3 shows the state after copying the backup paths from the edges in $G-T$ to the leaf edges. We consider all the leaf edges as labeled. After that the copy operations are performed on the other edges of T , which are not yet labeled.

Fig. 4 shows the final result of the algorithm. After processing leaf edges in the first round, edges (D, F) , (D, G) , (C, H) , and (I, M) are processed in the second round. Basically, the backup paths are copied towards the root of the tree. For edge (D, F) the backup path copied from edge (K, F) , i.e. $SABDGKFD/4$ is copied, but becomes removed, because the other backup path $(SACHLJFD/2)$ from edge (F, J) has a lower common path length (2 vs. 4). In a later round, edge (A, C) is considered. From the backup paths copied from edges (E, C) and (H, C) only the one from (H, C) is appropriate $(SABDFJLHCA/2)$. The other one from edge (C, E) $(SACHLNMIECA/3)$ contains C twice and must therefore be removed. When edge (S, A) is considered, all backup paths from edges (A, C) and (A, B) contain A twice. Therefore, (S, A) can not be repaired and is marked as not repairable.

4.3. Complexity analysis

4.3.1. General case

The complexity of the first part is $O(n)$. We basically calculate for each edge $(x, y) \in T$ the path from the root of the tree to y and the distance of y to the root.

The while loop in the second part is executed for each edge in $G-T$, i.e. with $O(m)$. The determination of the lowest common parent can be performed in $O(\log n)$

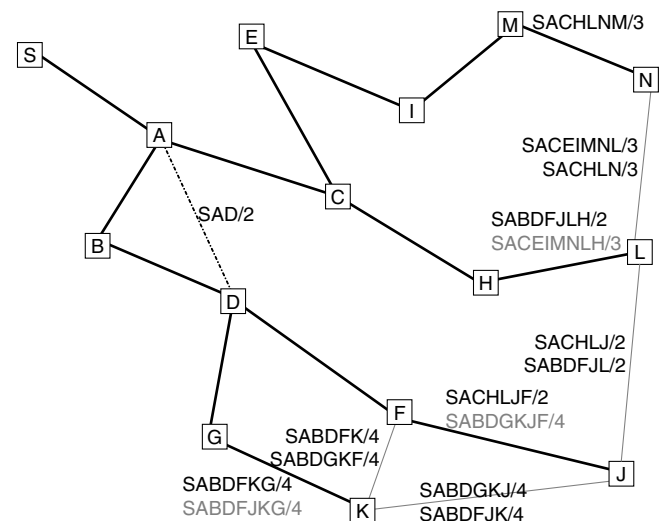


Fig. 3. Result of first round of backup multicast tree construction.

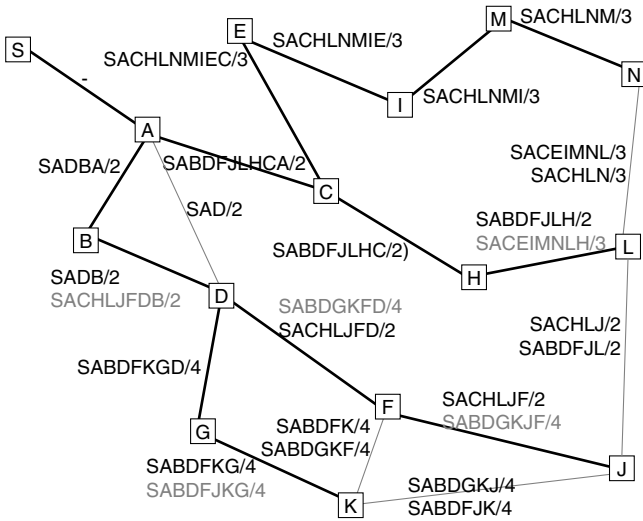


Fig. 4. Final result of backup multicast tree construction.

677 steps, if we apply binary search. For example, when we
 678 calculate the lowest common parent of vertices F and
 679 K, we have to consider the paths from the root to F
 680 and K, respectively. These paths are SABDF and SAB-
 681 DGK. The lowest common parent is D and D's distance
 682 to the root is 3. In our example, where the lowest com-
 683 mon parent for edge (F, K) is calculated, we could have
 684 selected position 3 (B), then position 5 (G vs. F), and
 685 finally position 4 (D). Applying binary search, searching
 686 the lowest common parent of the two nodes of an edge
 687 has a complexity of $O(\log D)$, with $D =$ the depth of the
 688 default multicast tree. Since $D \leq n$, the total complexity
 689 $O(m \log n)$.

690 In the third part of the algorithm the inner loop is exe-
 691 cuted once for each edge in G. This results in a complexity
 692 of $O(m)$. This also means that the overall complexity of the
 693 backup multicast tree algorithm is $O(m \log n)$.

694 4.3.2. Binary tree with full mesh

695 We now assume that the default multicast tree is a
 696 complete binary tree and the graph is a full mesh, i.e.
 697 each node is connected to any other. For the complexity
 698 analysis we assume that we perform the determination
 699 of the lowest common parent by sequential search in
 700 the paths starting at the root of the tree. We now take
 701 one edge (x, y) from G and perform the comparison
 702 described above for the two paths $S \rightarrow x$ and $S \rightarrow y$.
 703 The probability that the two paths along the binary tree
 704 to the two nodes x and y differ at the second node is at
 705 least $1/2$. This means that for $m/2$ edges (x, y) of G a
 706 single basic comparison operation (comparison whether
 707 two nodes are different) is sufficient to find a difference
 708 in the two paths $S \rightarrow x$ and $S \rightarrow y$. The probability that
 709 the two paths along the binary tree to the two nodes x
 710 and y are equal at the second node but differ at the
 711 third node is at least $1/4$. The probability that the
 712 two paths along the binary tree to the two nodes x

and y are equal at the i^{th} node but differ at the
 $i + 1^{\text{th}}$ node is at least $(1/2)^i$.

In case of a complete binary tree for the default mul-
 ticast tree and a full mesh for the complete graph an
 upper limit for the total number of path comparisons
 to find a difference in the two paths $S \rightarrow x$ and $S \rightarrow y$
 for all m edges (x, y) of the mesh is given by the follow-
 ing formula:

$$\frac{1}{2} \cdot m \cdot 1 + \frac{1}{4} \cdot m \cdot 2 + \frac{1}{8} \cdot m \cdot 3 + \frac{1}{16} \cdot m \cdot 4 + \dots + \frac{1}{2^{D-1}} \cdot m \cdot (D-1) = m \sum_{i=1}^{D-1} \frac{i}{2^i}$$

$$\lim_{D \rightarrow \infty} m \sum_{i=1}^{D-1} \frac{i}{2^i} = m \sum_{i=1}^{\infty} \frac{i}{2^i} = m \sum_{j=1}^{\infty} \sum_{i=j}^{\infty} \frac{1}{2^i} = m \sum_{j=1}^{\infty} \frac{1}{2^{j-1}} = m \sum_{j=0}^{\infty} \frac{1}{2^j} = m \cdot \frac{1}{1-\frac{1}{2}} = 2m.$$

This means that in a full mesh with a complete binary
 tree as default multicast tree the total number of compar-
 isons to find the lowest common parent of the two
 nodes of any edge is limited by $2m$. The total com-
 plexity of part 3 becomes $O(m)$. In that case, calculat-
 ing $n - 1$ backup links can even be performed with a
 lower complexity than computing the minimum span-
 ning tree.

4.4. Performance measurements

Fig. 5 shows the performance of an implementation
 of the backup multicast tree algorithm using gcc on
 an Intel™ Xeon 3.06 GHz CPU with 512 KB cache
 and 1 GB main memory. Random topologies with up
 to 500 nodes have been generated. For each node, links
 to 20% of other nodes randomly selected are generated
 $(m = 0.2 n^2)$. A topology with $n = 500$ nodes has
 $m = 500 * 500 * 0.2 = 50,000$ links. Calculating all back-
 up links for the $n - 1$ default multicast tree links takes
 less than 40 ms for a 500 node topology. The graph
 shows the measured time using 95% confidence intervals
 compared with the function $f(m, n) = 0.3 (m \log n)$. We
 see that the curve for $f(m, n)$ is growing faster than the
 curve representing the computing time of the backup
 spanning tree algorithm. Fig. 6 shows the performance
 when each node establishes only 4 links to other nodes
 $(m = 4n)$. In addition, the function $f(m, n) = 0.4 (m \log n)$
 is plotted. Again this indicates that our determined
 complexity of $O(m \log n)$ is correct.

Although we did not optimize the implementation, we
 compare the performance numbers with related work.
 Since we are not aware of a similar algorithm, we have
 to compare the performance with shortest path computa-
 tion algorithms. In [22], a fast shortest path algorithm
 has been implemented and the processing time per edge
 has been presented. Values $> 5 \mu\text{s}$ per edge are reported,
 however using less powerful hardware than for our eval-
 uation. Our implementation is running on a computer
 with a four times faster CPU and with four times more
 memory. It needs below $1 \mu\text{s}$ per edge. Note that a short-
 est path algorithm would have to run n times to get n
 backup trees.

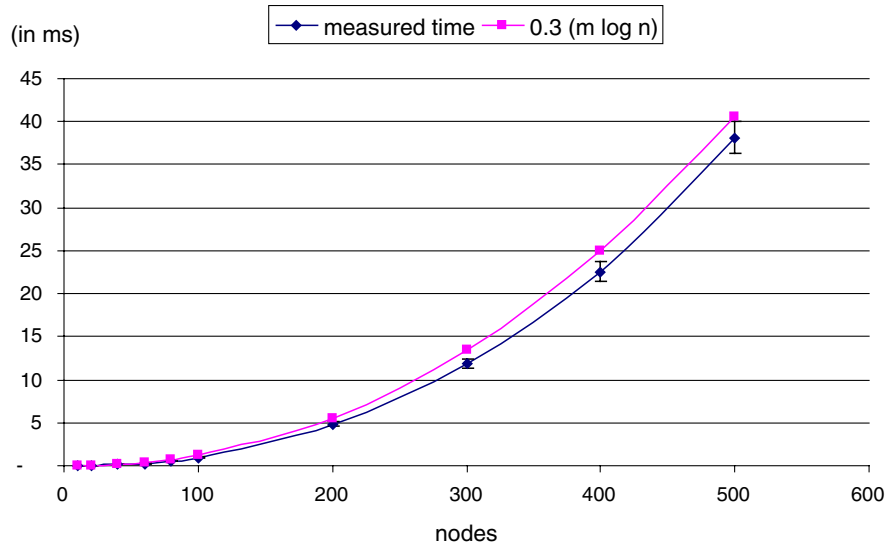


Fig. 5. Backup multicast tree performance (each node has connections to 20% of the other nodes).

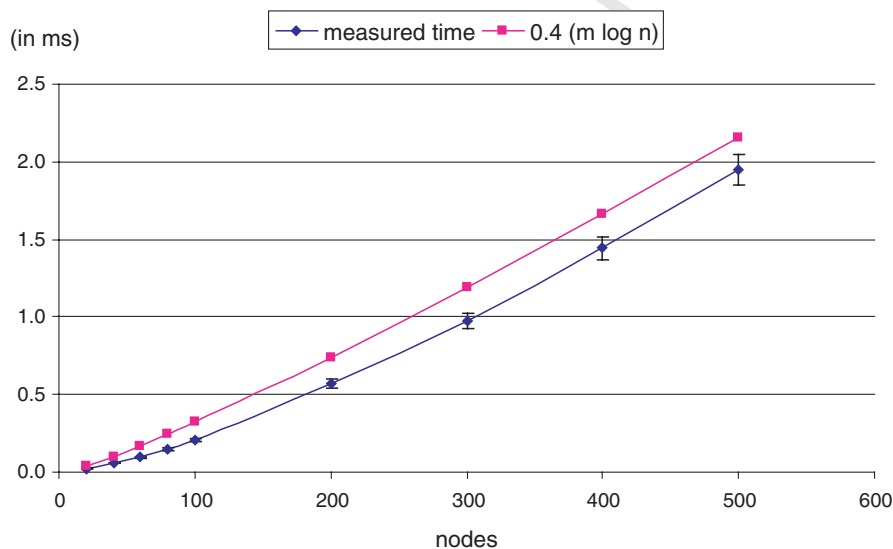


Fig. 6. Backup multicast tree performance (each node has four connections to other nodes).

764 5. Removing nodes from a multicast tree

765 5.1. Overview

766 Backup multicast trees can also support situations,
 767 where nodes leave the multicast group and new group keys
 768 need to be distributed among the remaining group mem-
 769 bers efficiently, but such that the leaving node does not
 770 receive the key. Our goal is to construct from the default
 771 multicast tree a new tree that covers all nodes except the
 772 leaving node. This new tree is also called *reduced multicast*
 773 *tree* hereafter. A reduced multicast tree can be derived from
 774 a single backup multicast tree, only if a node leaving the
 775 tree is not a branching point. In that case and assuming
 776 that x is the leaving node, w is the upstream node (the next
 777 node from x to the root), and y is the downstream node

(the directly connected child of x), a backup multicast tree 778
 779 should be constructed by replacing edge (x, y) by the cor-
 780 responding backup edge and by eliminating edge (w, x) .
 781 Given the example of Fig. 4 and assuming that node F
 782 leaves the group, we can construct a reduced multicast tree
 783 by replacing edge (F, J) by backup edge (L, J) and by
 784 removing edge (D, F) .

785 However, a single backup multicast tree is not able to
 786 support leaving nodes that are branching points of the tree.
 787 In the following we describe a general mechanism to con-
 788 struct a so-called reduced multicast tree. This tree can be
 789 used to reach all nodes of the default multicast tree, but
 790 not a single node that shall be removed from the tree. To
 791 remove a node from a default multicast tree, we have to
 792 remove the upstream edge of that node and to replace
 793 the downstream edges of that node by other links.

794 The replacement is required to re-connect the vertices fur- 849
 795 ther down the tree. The replacement of different down- 850
 796 stream edges can be supported by the reduced multicast 851
 797 tree algorithm presented in Section 5.2. A downstream 852
 798 edge can be replaced by its backup edge, if the following 853
 799 condition is fulfilled: The path $S \rightarrow y$ via the backup edge 854
 800 of (x, y) does not lead via x . It is important to note that 855
 801 if there exists such a backup edge that can replace the 856
 802 downstream edge, the algorithm presented in Section 4.2 857
 803 will find such a backup edge. 858

804 If we consider again the scenario in Fig. 4 and assume 859
 805 that node C leaves the group and needs to be removed, 860
 806 the backup edge for downstream edge (C, H) is edge (J, 861
 807 L). The resulting path from S to H is SABDFJLH. This 862
 808 path meets the condition above and does not lead via 863
 809 the leaving node C. However, the backup edge for edge 864
 810 (C, E) is edge (L, N) and the resulting path from S to 865
 811 E is SACHLNMIE. This path does not meet our condi- 866
 812 tion and leads via node C. Therefore, backup edge (L, 867
 813 N) is not appropriate to replace downstream edge (C, 868
 814 E). We have to emphasize here that if any of the vertices 869
 815 along this sub-tree (E, I, M, N) would have had an edge 870
 816 to another vertex not in the sub-tree of C, this edge would 871
 817 have been found as a backup edge for downstream edge 872
 818 (C, E) by the algorithm presented in Section 4. This 873
 819 means that if this sub-tree (E, I, M, N) can be connected 874
 820 to the default multicast tree without going via C, there 875
 821 must either exist a connection via the other sub-tree 876
 822 beginning at C or it cannot be connected at all to the 877
 823 default multicast tree. This observation leads to the fol- 878
 824 lowing algorithm for removing a node x from a default 879
 825 multicast tree and for constructing a reduced multicast 880
 826 tree. 881

827 5.2. Reduced multicast tree algorithm

828 In the first part of the algorithm given in C-like pseudo 882
 829 code below, we process all edges of T. If the backup path 883
 830 for edge (x, y) leads via vertex x , vertex y is colored red, 884
 831 otherwise it is colored green. If a vertex y is colored green, 885
 832 we can replace edge (x, y) by its backup edge for the con- 886
 833 struction of the reduced multicast tree. 887

834 In the second part, we process each edge of G–T. All 888
 835 these edges might be needed to connect the sub-trees of 889
 836 the red vertices to the reduced multicast tree. For each edge 890
 837 (x, y) of T we create an edge set. Then, we map each edge of 891
 838 G–T to one of these edge sets. An edge (w, z) is mapped to 892
 839 the edge set $D_{x,y}$, if x is the lowest common parent of w 893
 840 and z in T and if w is a child of y in T. The edges in an edge 894
 841 set $D_{x,y}$ are candidates to connect the sub-tree below y to 895
 842 the reduced multicast tree. 896

843 The third part processes all edges mapped to one of the 897
 844 edge sets in the second part. First, we select a node x to be 898
 845 removed and store all direct children of x in sets GREEN 899
 846 or RED depending on their color from the first part. We 900
 847 take one green node y after another and check whether 901
 848 an edge of its set $D_{x,y}$ connects to a sub-tree of a vertex 902

in set RED. If so, the red vertex becomes green and the 849
 edge is added to set I_x (set of interconnection edges for 850
 x). At the end of the algorithm, a reduced multicast tree 851
 can be constructed for all nodes x with only green direct 852
 children y . For constructing the reduced multicast tree, 853
 we take the default multicast tree and remove all edges that 854
 include x . We add all edges of sets B_x and I_x . This results 855
 again in a spanning tree, which includes all vertices of the 856
 default multicast tree except x . 857

If a node x has one or more directly connected red chil- 858
 dren, it is not possible to construct a reduced multicast tree. 859
 In this case, several group members become even discon- 860
 nected from the multicast group. It is not possible to build 861
 a new default multicast tree that includes those vertices 862
 with the current set of edges. New edges (overlay links) 863
 need to be established in this case. 864

865 Algorithm.

```

vertex x, y, a, b; 866
vertex_set RED, GREEN, Y; 867
edge e, f; 868
edge_set E, F, G, T,  $\forall(x, y) \in T: D_{x,y}, \forall x \in T: 869
I_x, B_x; 870
E := T; /* part 1: O(n) */ 871
while (E !=  $\emptyset$ ) { 872
  e := first_element(E); E := E - e; 873
  x := e.x; 874
  if (x is in e.backup_path) { 875
    e.y.color := red; 876
  } else { 877
    e.y.color := green; 878
    B_x := B_x + backup_edge(e); 879
  } 880
} 881
E := G-T; /* part 2: O(m) */ 882
order(E); /* order all edges according to */ 883
/* some predefined criteria */ 884
while (E !=  $\emptyset$ ) { 885
  e := first_element(E); E := E - e; 886
  x := (lowest common parent of e.x and 887
  e.y); 888
  a := (next node from x towards e.x); 889
  D_{x,a} := D_{x,a} + e; 890
} 891
V := {all vertices of T}; /* part 3: O(m) */ 892
while (V !=  $\emptyset$ ) { 893
  x := first_element(V); V := V - x; 894
  Y := {all vertices y | 895
   $\exists$  edge e in T with e.x == x, e.y == y} 896
  GREEN := {all green vertices of Y} 897
  RED := Y - GREEN; 898
  while (GREEN !=  $\emptyset$ ) { 899
    y := first_element(GREEN); GREEN := 900
    GREEN - y; 901
    F := D_{x,y}; 902
    while (F !=  $\emptyset$ ) { 903$ 
```

```

909     f := first_element(F); F := F - f;
910     b := (next node from x towards f.y);
911     if (b.color == red) {
912         Ix := Ix + f;
913         b.color := green;
914         GREEN := GREEN + b;
915         RED := RED - b;
916     }
917 }
918 }
919 }
920 }
921 }
922 }
923 }
924 }

```

925 Fig. 7 illustrates the reduced multicast tree construction
 926 process with a given default multicast tree consisting of the
 927 solid line edges. X is the node to be removed. The sub-trees
 928 of vertices Y₁ and Y₅ have backup edges that connect their
 929 sub-trees to the default multicast tree without going via X.
 930 The backup edges are added to B_X. Therefore, vertices Y₁
 931 and Y₅ are colored green initially, all others (Y₂, Y₃, Y₄)
 932 become red. In the second part, the various links of G–T
 933 are mapped to edge sets. For example, (W₁, Z₂) and (Z₁,
 934 Z₂) are mapped to edge set D_(X, Y₁). One of these two edges
 935 is processed first (let us assume (W₁, Z₂)) and it is discovered
 936 that this edge connects to the sub-tree of the red vertex
 937 Y₂. Y₂ becomes green and (W₁, Z₂) is added to set I_X. Later
 938 edges (W₂, Z₃) and (Z₅, W₄) are also added to I_X. In order
 939 to get a reduced multicast tree for node X we have to
 940 remove all edges with X as a vertex from the default multic-
 941 cast tree and add the edge sets I_X and B_X.

942 Note that if no reduced multicast tree for node X can be
 943 constructed by the given algorithm the nodes downstream
 944 of a leaving node X will be disconnected from the graph
 945 and no tree that includes all of them exists. In this case,
 946 the underlying peer-to-peer network must solve the connec-

947 tivity problem. This problem can be avoided if each node
 948 establishes a certain amount of links to other overlay
 949 nodes. It might also be helpful, if not only links to close
 950 peers are established. Otherwise, network partitioning
 951 might occur with a higher probability in case of link
 952 failures.

953 Instead of distributing a multicast message via a single
 954 reduced multicast tree, one could also use several trees, if
 955 a single reduced multicast tree cannot be constructed.

956 In the case of multiple leaves, we have to serialize the
 957 leaves and construct the reduced multicast trees accord-
 958 ingly. In this case, for a hierarchical scenario as depicted in
 959 Fig. 7 the algorithm can be performed concurrently in each
 960 sub-tree.

961 5.3. Complexity analysis

962 The loop in part 1 of the reduced multicast tree algo-
 963 rithm is performed once for each edge in T. This results
 964 in a complexity of $O(n)$. The second part has a complexity
 965 of $O(m)$, because the loop is performed for each edge in G–
 966 T. Also part 3 has a complexity of $O(m)$. The inner loop is
 967 performed once per edge in G–T. Each edge is in one of the
 968 sets D_(x,y). This means that an overall complexity of $O(m)$
 969 is required to construct n reduced multicast trees, if we
 970 assume that all backup edges are known in advance.

971 6. Tree IDs for multicast delivery trees

972 In the previous sections we have presented concepts to
 973 calculate backup and reduced multicast trees from a given
 974 default multicast tree. In particular, $n - 1$ backup multicast
 975 trees and $n - 1$ reduced multicast trees can be calculated in
 976 order to support failures or cases with leaving nodes. The

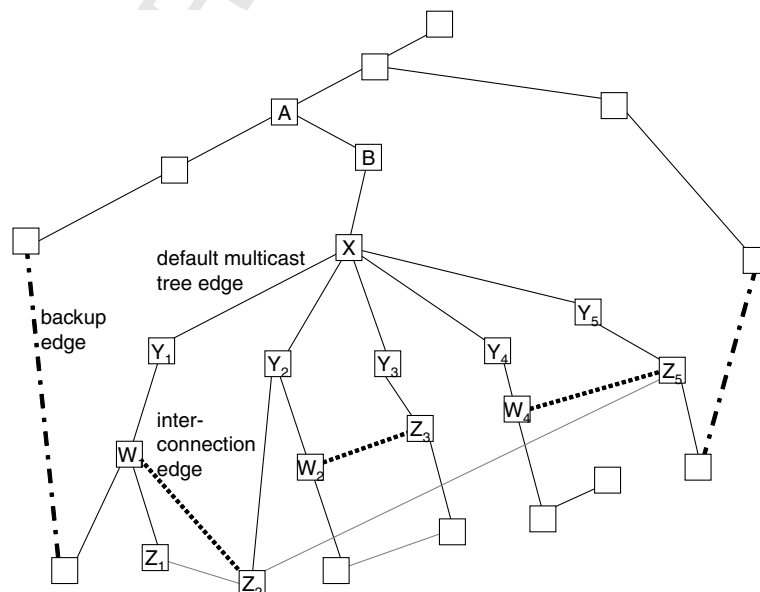


Fig. 7. Reduced multicast tree construction.

idea is that the sender (the root of the tree in case of source-specific multicast) selects an appropriate multicast delivery tree among several alternatives to distribute the multicast message. We propose to calculate for each alternative multicast delivery tree a unique identification called tree ID which allows a forwarding node to discover how a message must be forwarded. Such an ID should have the following characteristics:

- The tree IDs of two consecutive default multicast trees (before and after a node joins or leaves) should differ, because some messages might be delayed and messages that shall be forwarded along different trees might be present simultaneously. Also the protocol operations presented in Section 3.2 require that a tree ID remains valid for some time after a tree has been changed.
- The tree IDs of all backup and reduced multicast trees should be different in order to be able to distinguish which trees associated with a default multicast tree shall be used to distribute a multicast message.
- The size of such a tree ID should be limited in order to scale for large groups.

Based on this discussion, we propose to use three fields for specifying the selected multicast distribution tree:

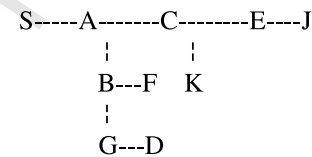
1. a default multicast tree ID for specifying the currently used default multicast tree used for the multicast group,
2. a type ID specifying whether the default multicast, one of its $n - 1$ backup multicast trees or one of its $n - 1$ reduced multicast trees shall be used,
3. a node ID specifying in case of a reduced multicast tree which node shall be excluded from the default multicast tree and in case of a backup multicast tree which upstream link of the specified node shall be replaced by its backup link.

For the calculation of the default multicast tree ID we propose a combination of a cardinal representation and MD5 [23]. A link-based scheme as used in [6] is not appropriate for our case, because links among peers will change more frequently than router links. Moreover, in case of large groups a tree ID consisting of all traversed link IDs might become too large.

The default multicast tree ID is based on a cardinal representation [24,25], which encodes the structure of the tree and its IDs separately. The structure is represented using a balanced parenthesis representation obtained by a *pre-order* traversal wherein a "(" is output when we enter a node and a ")" when we leave a node. This is then combined with the pre-order traversal of the node IDs. This representation takes $(n \log n + 2n)$ bits to encode a tree of n nodes. Further, it can be used to represent arbitrary sub-trees since the nodes are simply listed in pre-order. At each forwarding node, a single traversal of the encoding is sufficient to find the children and construct the encodings of the sub-trees. Thus at each node the number of bits in the encoding

reduce by a significant amount. For the tree given in the example below the following path ID is calculated at root S: (((((()())()((()))))SABGDFCKEJ. Since this path ID is variable in length and can easily become very long in large groups, we have to map it to a constant length identifier. We propose to apply a hashing mechanism on the cardinal representation and use the resulting hash value as path ID. In cases where the cardinal presentation is short enough, hashing could be avoided. Also in this case, the pre-computation of $n - 1$ backup trees or $n - 1$ reduced trees would not be required, since the sender of the multicast packet is able to completely specify the path to be taken by the multicast packet. In this case, our presented algorithms can be used to efficiently calculate alternative routes for single link breaks and node leaves. However, cardinal presentations are not limited in size and may exceed a given maximum size value [25].

Example:



7. Distributed algorithms

For the calculation of the backup and reduced multicast trees, we assumed that each node knows the default multicast tree and the complete mesh topology. Based on this knowledge, the backup and reduced multicast trees can be computed by each node independently according to the algorithms presented in Sections 4 and 5. An issue to be investigated in this section is whether the algorithms can be performed without that each node knows the complete topology. In the following sub-sections we show that this is possible, but under the constraint that additional signaling between the nodes is introduced.

For both sub-sections we assume a more light-weight basic signaling mechanism than described in Section 3.2 based on tree, join, and leave messages. A joining node sends a join message towards the root of the multicast delivery tree. The root in turn confirms the inclusion of the joining node by a tree message and inserts the path from the root to the joining node along the default multicast tree. This way, each node can learn the path from the root node to itself.

7.1. Backup multicast tree

For the distributed algorithm of the backup multicast tree algorithm, all the nodes need to know to which other nodes they connect to and which of these links are used for the default multicast tree. The algorithm makes use of two additional signaling messages in addition to tree, join, and leave messages:

- 1080 • Backup Path Establishment (BPE)
- 1081 • Backup Path Termination (BPT)

1082
1083 The protocol begins with the transmission of BPE mes-
1084 sages over links that are not part of the default multicast
1085 tree. Initially, the node originating a BPE message will
1086 put the IDs of the nodes along the default multicast tree
1087 path from the root to itself into the BPE message. We
1088 assume that each node has learnt that path before, e.g.
1089 by a tree message.

1090 A leaf node receiving a BPE message will then select
1091 from all received BPE messages that one that includes the
1092 best backup path as defined for the backup multicast tree
1093 algorithm in Section 4. The node will append its own ID
1094 to the backup path in the BPE message and forward it to
1095 its upstream node towards the root of the default multicast
1096 tree. A forwarded BPE indicates that a node is part of a
1097 backup path in order to replace a link. For all BPE messag-
1098 es that are not selected by a leaf node, a BPT message is
1099 sent back towards the originator of the BPE message.
1100 The node creating a BPT message is also called terminating
1101 node hereafter. BPT messages are forwarded in the reverse
1102 direction as the corresponding BPE messages until the final
1103 destination (the originating node of a BPE message) has
1104 been reached. BPT messages contain the recorded path
1105 information from the root via the node originating node
1106 of the BPE message to the terminating node. If a node
1107 receives a BPT message it can learn from the included
1108 backup path, for which links the backup path (and the
1109 backup link) have been selected. In particular, these are
1110 all the links between the nodes originating and terminating
1111 the BPE message.

1112 A node, which is not a leaf node in the default multicast
1113 tree, will receive BPE messages from links that are not part
1114 of the default multicast tree, but also from its downstream
1115 nodes via default multicast tree links. Processing of BPE
1116 messages is performed in exactly the same way as in the
1117 case of leaf nodes. Forwarding of BPE messages to
1118 upstream nodes shall be performed periodically and a node
1119 has to consider BPE messages received from all other links.

1120 For example, we consider Fig. 4. Node L issues a BPE
1121 message with path SACHL to node N. N forwards the
1122 BPE message via M, I, and E to C. C terminates this
1123 BPE message and returns a BPT message including the
1124 backup path SACHLNMIEC. Backup link (L, N) can
1125 serve as backup link for the default multicast tree links
1126 between C and L. These are (C, E), (E, I), (I, M), and
1127 (M, N). Later, any node between the root and C might
1128 detect or be informed (via additional failure notifications)
1129 that there exists a problem on link (C, E). Such a node
1130 might simply insert a tree ID into the multicast packet indi-
1131 cating that the packet shall be forwarded via the backup
1132 multicast tree instead of the default multicast tree. This
1133 means that the packet shall be forwarded via the backup
1134 link (L, N) instead of link (C, E). When node C receives
1135 such a message, it does not forward the packet via link
1136 (C, E). The packet arrives also at node L, which has

1137 learned from the exchange of BPE and BPT messages that
1138 link (L, N) serves as backup link for link (C, E). Therefore,
1139 L forwards the multicast packet via the backup link to N.
1140 Then, the packet travels to E along the branch of the
1141 default multicast delivery tree, but in the opposite
1142 direction.

1143 The signaling overhead associated with this procedure
1144 depends on the number of edges (m). There is not more
1145 than one BPE message on each of the m links. Since the
1146 BPT messages travel exactly on the reverse path back to
1147 the originators of BPE messages, also not more than m
1148 BPT messages are generated. Note that in a dynamic envi-
1149 ronment these messages should be distributed periodically.
1150 In that case one BPE and one BPT message occur on each
1151 link per interval.

7.2. Reduced multicast tree

1152
1153 With the signaling protocol described in Section 7.1 each
1154 node learns the backup links and the backup paths for all
1155 links, to which it is directly connected. With this informa-
1156 tion a node $Y_i (i = 1, \dots, 5)$ in Fig. 7 can derive, whether
1157 there exists a backup path for link (X, Y_i) that does not
1158 lead via node X. If such a backup path exists node Y_i is
1159 a green node and can be connected to the multicast tree
1160 via the backup link for (X, Y_i). If this is not the case, node
1161 Y_i is a red node and has to search for a connection to one
1162 of the other sub-trees of X that are represented by the
1163 nodes Y_i .

1164 This search can be supported by another signaling proto-
1165 col extension. Each node Y_i that can not be connected
1166 via backup links to the multicast tree has to send an *explore*
1167 message towards the children along the default multicast
1168 tree. The message contains the link (X, Y_i) as a parameter.
1169 The message is flooded on the complete sub-tree of Y_i , and
1170 each node on this sub-tree that has a link from (G–T) for-
1171 wards the message to its neighbor. The neighbor receives
1172 this message and discards it, if X is not on the path from
1173 the root to itself. If X is on the path from the root to itself,
1174 the message is forwarded towards X and will be received by
1175 a child of X that is directly connected to X. This direct
1176 child is one of the downstream nodes, for example Y_j . If
1177 Y_j is a green node, it returns an *interconnect* message to
1178 Y_i in the reverse direction than the explore message. The
1179 message passes one link that is not in T. We call this link
1180 interconnection link. The nodes of the interconnection link
1181 learn that this link is required to reach node Y_j , if node X
1182 becomes removed and will later forward multicast packets
1183 that contain a tree ID for the reduced multicast tree for
1184 node X. After node Y_j has received the interconnect mes-
1185 sages, it becomes a green node and can also return intercon-
1186 nect messages for incoming explore messages from other
1187 red nodes.

1188 In our example given in Fig. 7, Y_1 and Y_5 are green
1189 nodes, the other nodes Y_2 – Y_4 are red. Y_2 sends an explore
1190 message, because it has learned that there is no backup link
1191 for (X, Y_2) that does not lead via X. The explore message is

forwarded by W_2 via Z_3 to Y_3 , which is a red node too and does not respond the explore message. However, the explore message is forwarded by Z_2 via Z_1 and W_1 to Y_1 . Y_1 responds with an interconnect message that travels back to Y_2 . W_1 learns that (W_1, Z_2) will be required to reach Y_2 if node X becomes removed. Y_2 receives the interconnect message, becomes red and may later respond to explore messages from other red nodes such as Y_3 with interconnect messages.

The signaling mechanism to support reduced multicast trees requires an exchange explore/interconnect messages for each red node of the tree. The overhead can be reduced by piggy-backing the various explore and interconnect messages that are traveling up and down the sub-tree below the node to be removed.

8. Conclusions

In this paper, we have proposed a concept for explicit routing in multicast overlay networks. In particular, it allows specifying a particular distribution tree for the transmission of multicast data. We proposed to calculate $n - 1$ backup multicast trees for a given default multicast tree (e.g., a minimum spanning tree) that interconnects the n nodes belonging to a group. The complexity for calculating these backup multicast trees is slightly above the complexity for calculating a minimum spanning tree but can be even lower with certain types of graphs. The performance measurements of the backup multicast tree algorithm implementation confirm the determined complexity. In addition, an algorithm has been presented that allows calculating a reduced multicast tree by discarding a particular node, e.g. a node leaving the multicast tree, from the default multicast tree. Unique IDs for the alternative (default, backup, or reduced) multicast trees are required in order to specify which of the trees shall be used for multicast message forwarding. An appropriate encoding scheme has been developed and discussed. We also described a distributed version of the algorithms avoiding that each node needs to be aware of the full topology. This requires introducing a light-weight signaling protocol.

Acknowledgements

We thank Prof. Shivkumar Kalyanaraman, who introduced us to the problem of explicit multicast for overlay networks. The second author has worked with him to design encoding schemes for multicast trees [25].

References

- [1] C.W. Kong, M. Gouda, S. Lam, Secure group communications using key graphs, *IEEE/ACM Transactions on Networking* 8 (1) (2000) 16–30.
- [2] S. Banerjee, B. Bhattacharjee, Scalable secure group communication over IP multicast, *IEEE Journal on Selected Areas in Communications* 20 (8) (2002) 1511–1527.

- [3] S. Deering, B. Hinden, Internet Protocol, Version 6 (IPv6) Specification, RFC 2460, December 1998.
- [4] R. Boivie, N. Feldman, Y. Imai, W. Livens, D. Ooms, O. Paridaen, Explicit Multicast (Xcast) Basic Specification, Internet Draft, work in progress, August 2003.
- [5] E. Rosen, A. Viswanathan, R. Callon, Multiprotocol Label Switching Architecture, RFC 3031, January 2001.
- [6] H.T. Kaur, S. Kalyanaraman, A. Weiss, S. Kanwar, A. Gandhi, BANANAS: an evolutionary framework for explicit and multipath routing in the internet, in: *ACM SIGCOMM 2003 Workshop on Future Directions on Network Architectures (FDNA)*, Karlsruhe/Germany, August 25–27, 2003.
- [7] D. Eppstein, Finding the k Shortest Paths, 35th Symposium on Foundations of Computer Science, IEEE, November 1994.
- [8] M. Castro, P. Druschel, A.M. Kermarrec, A. Rowstron, Scribe: a large-scale and decentralized application-level multicast infrastructure, *IEEE Journal on Selected Areas in Communications* 20 (8) (2002).
- [9] A. Rowstron, P. Druschel, Pastry: Scalable, Distributed Object Location and Routing for Large-scale Peer-to-peer Systems, *IFIP/ACM Middleware 2001*, Heidelberg/Germany, November 2001.
- [10] S. Zhuang, B. Zhao, A. Joseph, R. Katz, J. Kubiatiowicz, Bayeux: an architecture for scalable and fault-tolerant wide-area data dissemination, in: *11th International Workshop on Network and Operating System Support for Digital Audio and Video*, Port Jefferson/USA, June 2001.
- [11] B. Zhao, L. Huang, J. Stribling, S. Rhea, A.D. Joseph, J.D. Kubiatiowicz, Tapestry: a resilient global-scale overlay for service deployment, *IEEE Journal on Selected Areas in Communications* 22 (1) (2004).
- [12] S. Banerjee, B. Bhattacharjee, Ch. Kommareddy, Scalable application layer multicast, in: *ACM SIGCOMM 2002*, Pittsburgh/USA, August 19–23, 2002.
- [13] S. Banerjee, S. Lee, B. Bhattacharjee, A. Srinivasan, Resilient multicast using overlays international conference on measurements and modelling of computer systems, in: *ACM SIGMETRICS 2003*, San Diego/USA, June 9–14, 2003.
- [14] Y. Chu, S. Rao, S. Seshan, H. Zhang, Enabling conferencing Applications on the Internet using an overlay multicast architecture, in: *ACM SIGCOMM 2001*, San Diego, August 27–31, 2001, pp. 55–67.
- [15] Y. Chu, S. Rao, S. Seshan, H. Zhang, A case for end system multicast, *IEEE Journal on Selected Areas in Communications* 20 (8) (2002).
- [16] Zhi Li, P. Mohapatra, Hostcast: a new overlay multicasting protocol, in: *IEEE International Conference on Communications (ICC) 2003*, Anchorage/USA, May 11–15, 2003.
- [17] Y. Zhu, B. Li, J. Guo, Multicast with network coding in application-layer overlay networks, *Journal on Selected Areas in Communications* 22 (1) (2004) 1–13.
- [18] S. Ratnasamy, P. Francis, M. Handley, R. Karp, S. Shenker, A Scalable content addressable network, in: *ACM SIGCOMM 2001*, San Diego/USA, August 27–31, 2001.
- [19] I. Stoica, R. Morris, D. Liben-Nowell, D.R. Karger, M.F. Kaeshoek, F. Dabek, H. Balakrishnan, Chord: a scalable peer-to-peer lookup protocol for internet applications, *IEEE/ACM Transaction on Networking* 11 (1) (2003).
- [20] D. Andersen, H. Balakrishnan, M. Kaashoek, R. Morris, The case for resilient overlay networks, in: *Proceedings of the Eighth Annual Workshop on Hot Topics in Operating Systems (HotOS-VIII)*, Schloss Elmau, Germany, May 2001.
- [21] R. Prim, Shortest connection networks and some generalizations, *Bell System Technical Journal* 36 (1957) 1389–1401.
- [22] S. Pettie, V. Ramachandran, S. Sridhar, Experimental evaluation of a new shortest path algorithm, in: *Proceedings of the Fourth Workshop on Algorithm Engineering and Experiments*, LNCS vol. 2409, pp. 126–142, 2002.
- [23] R. Rivest, The MD5 Message-Digest Algorithm, RFC 1321, April 1992.

- 1312 [24] D. Benoit, E. Demaine, J. Munro, V. Raman, Representing trees of
1313 higher degree, in: Proceedings of the Sixth International Workshop
1314 on Algorithms and Data Structures (WADS'99), pp. 169–180, 1999.
- [25] V. Arya, T. Tulletti, S. Kalyanaraman, Encodings of multicast trees, 1315
in: IFIP Networking Conference, Waterloo, Canada, May 2–6, 2005, 1316
pp. 992–1004. 1317
1318

UNCORRECTED PROOF