# Multicast Topology Inference from End-to-end Measurements[*]

N.G. Duffield[‡]    J. Horowitz[♠]    F. Lo Presti[‡,§]    D. Towsley[§]

[‡]AT&T Labs–Research
180 Park Avenue
Florham Park, NJ 07932, USA
{duffield,lopresti}@research.att.com

[♠]Dept. Math. & Statistics
University of Massachusetts
Amherst, MA 01003, USA
joeh@math.umass.edu

[§]Dept. of Computer Science
University of Massachusetts
Amherst, MA 01003, USA
towsley@cs.umass.edu

### Abstract

This paper describes a class of statistical estimators of multicast tree topology based on end-to-end multicast traffic measurements. This approach allows the determination of the logical multicast topology without assistance from the underlying network nodes. We provide six instances of the class, variously using measurements of loss or delay, or marks set by network elements. We compare their accuracy and computational cost, and recommend the best choice in each of the light and heavy traffic load regimes.

## 1   Introduction and Motivation

The use of multicast shows great promise for determining internal network characteristics based solely on end-to-end measurements. This is because multicast introduces correlation in the end-to-end behavior observed by different receivers within the same multicast session. This correlation can be used to estimate packet loss rates, [1], packet delay distributions, [6], and packet delay variances, [5]. These methods can be used as part of a multicast–capable measurement infrastructure, such as NIMI (National Internet Measurement Infrastructure) [10], for the purpose of monitoring internal network behavior.

All of these multicast-based methods can be applied to end-to-end multicast observations made of packets traversing a fixed, but arbitrary, topology. Knowledge of the multicast topology is required in order to apply the methods. Unfortunately, knowledge of the tree topology is not always available. This motivates the need for algorithms that can identify the topology of the multicast tree. Another motivation is that knowledge of the multicast

---

topology can be of use to multicast applications. For example several reliable multicast protocols (e.g., RMTP [9]) rely on logical hierarchies based on the underlying topology if possible. Other applications attempt to group receivers that share the same network bottleneck, [12].

In this paper, we present a general framework within which to develop algorithms for identifying the multicast topology based on end-to-end observations. These observations can be end-to-end measurements of packet loss or delay, or can be the values observed at receivers of marks set in packets by routers on the end-to-end path. An example of the latter is the recent proposal for Explicit Congestion Notification (ECN) [11] by network elements. Topology inference from observed mark distributions will be useful if end-point congestion control mechanisms, coupled with ECN, make packet loss and delay too infrequent to infer topology from.

Once the topology has been identified, any of the methods mentioned above for identifying internal network behavior can then be applied. The development of an algorithm for topology identification is based on the presence of a packet performance measure that monotonically increases as the packet traverses down the tree and that can be estimated based on observations made at the receivers. We provide additional constraints on the measures such that the resulting algorithm can be shown to be asymptotically consistent, i.e., it identifies the correct topology almost surely as the number of observations goes to infinity. Examples of performance measures that yield such algorithms include loss rate, delay variance, average delay, and link utilization; the corresponding measure for packet marking is the marking probability.

Several algorithms have been proposed for identifying multicast topologies based on loss observations made at receivers. For example, [12] presented an algorithm for identifying a multicast tree when it is a binary tree. [2] established the correctness of this algorithm and introduced several other loss-based algorithms for identifying general trees. They concluded that an algorithm that constructs a binary tree and subsequently prunes branches whose loss rates are close to zero was best suited for topology identification. The framework presented in this paper comes from the recognition that this last approach can be applied to observations of other end-to-end measures such as delays.

Topology discovery is an essential part of several current measurement infrastructure projects, including CAIDA, Felix, IPMA, NIMI and Surveyor; see [3]. We contrast our approach with that of the commonly used diagnostic tools `traceroute` and `mtrace` [7] that discover physical topology. These require cooperation from intervening nodes (in the generation of ICMP messages, or in maintaining counters) and their widespread use raises issues of scaling in topologies with many leaves. The present methods complement these, being able to discover logical multicast topology and its changes without cooperation from the network. Moreover, the use of multicast probes yields good scaling properties for measurement traffic volumes in larger networks; see [2] for further discussion.

The paper is organized as follows. We introduce our framework in Section 2 and provide

conditions under which the resulting algorithms are strongly consistent. Applications to loss and delay measures are presented in Section 3. In Section 4 we analyze the probability of topology misclassification, asymptotically for large numbers of probes. Section 5 reports on a simulation study on the effectiveness of different algorithms obtained through this approach and makes recommendations as to when they perform well. Section 6 concludes the paper.

## 2  Framework for Topology Inference

In this section we set up a general formalism for topology inference, and prove a general result that establishes the consistency of topology inference from end measurements that satisfy certain axioms. The first step is to set up the notation with which to describe multicast topology. Inference of this topology from end-to-end measurements is possible as a result of inferring characteristics, such as loss and delay distributions, of paths within the topology. For each multicast packet, we associate with each link a sample value of the characteristic, e.g., the delay that would be accrued by a packet traversing the link. We call these values **marks**. The term can apply literally, as when we consider marks set in the packet for ECN, or figuratively, as for delay. Each path from the source to a receiver gives rise to a set of marks, i.e., those of the links in the path. We can associate a mark with an entire path by **composition** of the marks of the constituent link, e.g., by adding the link delays to form the path delay. These composite marks represent the observations made at the receivers for the packet. It turns out to be useful to have a compact way to represent the collected measurements at groups of receivers as a single mark. For this purpose we define a further operation on marks, called **aggregation**. The fact that packets were multicast gives rise to certain algebraic relations between the composition and aggregation operations. The key behind topology inference is that the characteristics of paths are both **estimable** (in that they can be determined from measurement at receivers) and **increasing** (in that extending the path increases the corresponding characteristic; e.g. as for mean delay along a path). This property allows us to identify as siblings, those pairs of nodes for which the characteristic on the path to their closet common ancestor is maximized. Iteration of this argument allows us to identify the whole tree. We now set out these structures in detail.

**Tree Model.**  The physical multicast tree comprises actual network elements (the nodes), and the communication links that join them. When a packet is multicast to a set of receivers, only one copy of the packet traverses each link of the tree; copies are created at the branch points of the tree, one per outgoing link. The logical multicast tree comprises the branch points of the physical tree and the logical links between them. The logical links comprise one or more physical links. Thus each node in the logical tree, except the leaf nodes and, possibly, the root, must have 2 or more children.

3

Let $\mathcal{T} = (V, L)$ denote a logical multicast tree with nodes $V$ and links $L$. We identify the root node $0$ as the source of probes, and $R \subset V$ as the set of leaf nodes (identified as the set of receivers). The set of children of node $j \in V$ is denoted by $d(j)$. Each node, $k$, apart from the root, has a parent $f(k)$ such that $(f(k), k) \in L$. Define recursively the compositions $f^n = f \circ f^{n-1}$ with $f^1 = f$. We will sometimes refer to the link $(f(k), k)$ as simply link $k$. Nodes are said to be siblings if they have the same parent. If $k = f^m(j)$ for some $m \in \mathbb{N}$ we say that $j$ is descended from $k$ (or equivalently that $k$ is an ancestor of $j$) and write the corresponding partial order in $V$ as $j \prec k$. $a(i, j)$ will denote the minimal common ancestor of $i$ and $j$ in the $\preceq$-ordering. For $k \in V$ we let $\mathcal{T}(k) = (V(k), L(k))$ denote the subtree of $\mathcal{T}$ that is rooted at $k$, and set $R(k) = R \cap V(k)$.

**Tree Marks.** The experience of a multicast packet on its passage down the tree is modeled by a random process of **marks** $x = (x_k)_{k \in V}$. Each mark $x_k$ takes a value in a set $\mathcal{X}$ appropriate to the problem of interest. $x_k$ specifies the experience of a packet traversing link $k$. In the setting of packet loss, for example, we take $\mathcal{X} = \{0, 1\}$, where $x_k = 1$ indicates that a packet present at node $f(k)$ is successfully transmitted to node $k$, while $x_k = 0$ indicates that it would be lost.

**Composing Marks Along Paths.** A path is a set of contiguous links, identified by the ordered set of link endpoints $(k_1, \ldots, k_\ell)$ where $k_i = f(k_{i+1})$. We will sometimes use the notation $\boldsymbol{p}(k)$ to denote the path that from the root $0$ to node $k$. The experience of the multicast packet on a path is obtained by composing the marks from each link on the path to form a mark for the path. *Composition* is an associative and commutative binary operation $\otimes$ on $\mathcal{X}$. A path $\boldsymbol{p} = (k_1, \ldots, k_\ell)$ has mark $x_{\boldsymbol{p}}$ formed by successively composing the marks of its constituent links: $x_{\boldsymbol{p}} = x_{k_1} \otimes \ldots \otimes x_{k_\ell}$. We assume that $\mathcal{X}$ contains an identity $\boldsymbol{z}$ such that $\boldsymbol{z} \otimes x = x$ for all $x \in \mathcal{X}$. This $\boldsymbol{z}$ leaves path marks unaltered. In the delay example, composition is by addition and $\boldsymbol{z} = 0$, i.e. zero delay. In the loss example, we can compose link marks using the minimum $x_1 \otimes x_2 = x_1 \wedge x_2$. This models the physical property the loss occurs on a path if it occurs for any link on the path. The identity is $\boldsymbol{z} = 1$.

**Measurements and Marks.** We assume that the individual marks $x_k$ are not directly knowable. Rather, end-to-end measurements comprise the marks $x_{\boldsymbol{p}(k)}$ along paths terminating at leaf nodes $k \in R$. Our task will be to infer the underlying topology from these path marks alone. This information can also be used to characterize the individual links, by inferring the distribution of their link marks.

**Mark Aggregation.** We also equip $\mathcal{X}$ with an aggregation operation that summarizes the experience of packets over a set of possibly intersecting paths. We restrict attention to

binary trees. *Aggregation* is then a binary operation $\oplus$ on $\mathcal{X}$. We express the multicasting property in the axiom that composition $\otimes$ is distributive over aggregation $\oplus$, i.e.,

$$(x_1 \otimes x_2) \oplus (x_1 \otimes x_3) = x_1 \otimes (x_2 \oplus x_3) \tag{1}$$

for any $x_1, x_2, x_3 \in \mathcal{X}$. It expresses the fact that contribution to the path marks from a single probe on two intersecting paths from their common portion is identical. Consider, for example, a four-node logical multicast tree with root $0$ having a single child $1$, the latter node having children $2, 3$ that are leaves. Eq. (1) says that when calculating the aggregate mark for intersecting paths $(1, 2)$ and $(1, 3)$, we can factor out the common mark on the common link $1$. When dealing with loss, we will take aggregation as maximization, i.e. $x_1 \oplus x_2 = x_1 \vee x_2$. When dealing with delay, we will also take aggregation as minimization, i.e. $x_1 \oplus x_2 = x_1 \wedge x_2$. The property (1) holds for all the examples that we consider.

**Aggregating Receiver Marks.**    It will turn out that identification of the topology does not need the whole joint distribution of the receiver marks, but only of certain combinations that we now describe. Since we deal with binary trees, nodes $k \in V \setminus (R \cup \{0\})$ have two children, which we denote by $h(k)$ and $h^*(k)$. For each $k \in V \setminus \{0\}$ we define the aggregate marks over the paths to all receivers descended from $k$ recursively through
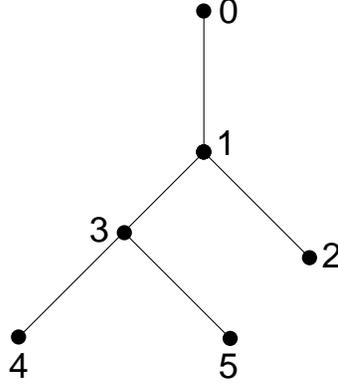
$$\widetilde{x}_k = \left\{ \begin{array}{ll} x_{\boldsymbol{p}(k)} & k \in R \\ \widetilde{x}_{h(k)} \oplus \widetilde{x}_{h^*(k)} & \text{otherwise} \end{array} \right. . \tag{2}$$

When $\oplus$ is associative (i.e. $(m_1 \oplus m_2) \oplus m_3 = m_1 \oplus (m_2 \oplus m_3)$) we can write $\widetilde{x}_k = \oplus_{j \in R(k)} x_{\boldsymbol{p}(j)}$. This is the case for loss, where $\widetilde{x}_k$ is $1$ if the packet reaches some receiver descended from $k$ and $0$ otherwise. However we give an example where $\oplus$ is not associative; see Section 3.2. Note the special case of the aggregate marks, namely, the receiver marks $\widetilde{x}_k = x_{\boldsymbol{p}(k)}$ for $k \in R$.

**Mark Distributions.**    We assume that the marks are independently distributed on different links, according to a mark distribution $\omega = (\omega_k)_{k \in V}$, where $\omega_k$ is the distribution of the mark $x_k$.

   The distribution $\omega_{\boldsymbol{p}}$ of the composite mark of a path $\boldsymbol{p} = (k_1, \ldots, k_\ell)$ is determined by convolution in the usual way: for a measurable subset $B$ of $\mathcal{X}$, $\omega_{\boldsymbol{p}}(B) = (\omega_{k_1} * \ldots * \omega_{k_\ell})(B) := \int_{x_1 \otimes \ldots \otimes x_\ell \in B} \omega_{k_1}(dx_1) \ldots \omega_{k_\ell}(dx_\ell)$. The joint distribution of the aggregate marks $\widetilde{x}_{k_1}, \ldots, \widetilde{x}_{k_\ell}$ will be denoted $\widetilde{\omega}_{k_1, \ldots, k_\ell}$.

**Deterministic Reconstruction of Binary Trees.**    The classification of trees relies on being able to identify certain characteristics of paths that do not terminate at leaves from the characteristics of those that do. Such a characteristic $\phi$ will be termed *estimable*; the precise

$$\Phi_{4,5} = \phi(\omega_{p(3)}) > \phi(\omega_{p(1)}) = \Phi_{2,4} = \Phi_{2,5}$$

Figure 1: RECONSTRUCTION OF A BINARY TREE. $\phi(\omega_{p(3)}) > \phi(\omega_{p(1)})$ if $\phi$ is increasing. The equalities hold if $\phi$ is estimable. Since $\Phi_{4,5} > \Phi_{2,4}, \Phi_{2,5}$, we chose to pair nodes $4$ and $5$ as siblings, rather than with node $2$.

definition is below. We seek estimable characteristics that *increase* as paths lengthen; this allows us to select as siblings those nodes for which the characteristic $\phi$ on the common portion of the path from $0$ is maximized. Packet loss rate and mean packet delay along a path are amongst the examples of such characteristics that we consider.

Let $\mathcal{T} = (V, L)$ be a binary tree. Let $\Omega$ be a set of probability distributions on $\mathcal{X}$ that is closed under convolution, and denote by $\Omega^V$ the set of corresponding product distributions of marks on $\mathcal{T}$. The characteristic $\phi$ is taken to be a weakly continuous functional on the set of measures on $\mathcal{X}$, which takes values in some partially ordered topological space $\mathcal{Q}$. In most examples $\mathcal{Q}$ will be the real numbers. Weak continuity will be useful when we study convergence of the estimators as the number of probes grows. Let $\delta_{\boldsymbol{z}}$ denote the distribution which has unit mass at the identity $\boldsymbol{z}$.

**Definition 1**    *(i)*  $\phi$ *is called* **estimable** *if there exists a function $\Phi$ that, for each $\omega \in \Omega^V$, and $j, k \in V$ with $a(j, k) \neq j, k$, expresses the characteristic of the path $\boldsymbol{p}(a(j, k))$ to the common ancestor $a(j, k)$ in terms of the distributions of aggregate marks $\widetilde{\omega}_{j,k}$ at the descendant nodes of $j, k$ through $\phi(\omega_{\boldsymbol{p}(a(j,k))}) = \Phi(\widetilde{\omega}_{j,k})$.*

  *(ii)*  $\phi$ *is called* **increasing** *if extending a path increases the characteristic $\phi$, i.e., $\phi(\delta_{\boldsymbol{z}}) = 0 \in \mathcal{Q}$, and for all $\omega \in \Omega$, $\phi(\omega_{\boldsymbol{p}}) < \phi(\omega_{\boldsymbol{q}})$ when $\boldsymbol{p}$ is a proper subpath of $\boldsymbol{q}$. The condition $\phi(\delta_{\boldsymbol{z}}) = 0$ says $\phi$ is not increased by a link that does not change the marks of paths traversing it.*

6

1. *Input*: The set of receivers $R = \{i_1, \ldots, i_r\}$ and
   the leaf mark distributions $\widetilde{\omega}_{i_1, \ldots, i_r}$.
2. $R' := R; V' := R'; L' = \emptyset$ ;
3. **while** $|R'| > 1$ **do**
4.      **select** $U = \{j, k\} \subseteq R'$ with maximal $\Phi_{j,k}$;
5.      $V' := V' \cup \{U\}$;
6.      $L' = L' \cup \{(U, \ell) : \ell \in U\}$;
7.      $R' := (R' \setminus U) \cup \{U\}$;
8. **enddo**
9. **if** $\Phi_{j,k} > 0$ **do**
10.      $V' := V' \cup \{0\}$ ; $L' = L' \cup \{(0, R')\}$ ;
11. **enddo**
12. *Output*: tree $(V', L')$ ;

Figure 2: Deterministic Binary Tree Classification Algorithm (DBT).

Given an estimable, increasing $\phi$, and a distribution $\omega$, we use the shorthand $\Phi_{j,k}$ for $\Phi(\widetilde{\omega}_{j,k})$. The topology can be reconstructed using the $\Phi_{j,k}$ as follows. The key observation from (ii) is that $\Phi_{j,k} > \Phi_{j',k'}$ whenever $a(j, k) \prec a(j', k')$. Thus $\Phi_{j,k}$ is maximized when $j, k$ are siblings in $R$. If not, then one of the receivers, say $j$, would have a sibling $k'$ for which $\Phi_{j,k'} > \Phi_{j,k}$. Thus the siblings can be identified on the basis of leaf distributions alone. This is illustrated for a three-leaf tree in Figure 1. Substituting a composite node that represents their parent and iterating, should then reconstruct the binary tree. This approach is formalized in the Deterministic Binary Tree Classification Algorithm (DBT); see Figure 2.

DBT operates as follows. $R'$ denotes the current set of nodes from which a pair of siblings will be chosen, initially equal to the receiver set $R$. We first find the pair $U = \{j, k\}$ that maximizes $\Phi_{j,k}$ (line 4). This identifies the members of $U$ as siblings, and the set $U$ is used to represent their parent. Correspondingly, we add $U = \{j, k\}$ to the list $V'$ of nodes (line 5), we add $(U, j), (U, k)$ to the list $L'$ of links (line 6), and replace $j$ and $k$ by $U$ in the set $R'$ of nodes available for pairing in the next stage (line 7). This process is repeated until all sibling pairs have been identified (loop from line 3). Finally, we test in line 9 whether the last node grouped should be taken as the root node. If the last node identified were not the root node, equality in the test would contradict the increasing property of $\phi$. Otherwise, we adjoin a root node, and a link joining it to its single descendant (line 10).

We say that DBT reconstructs the binary logical multicast tree $(V, L)$ if given the receiver set $R = \{i_1, \ldots, i_r\}$ and the leaf mark distributions $\widetilde{\omega}_{i_1, \ldots, i_r}$, it produces $(V, L)$ as its output. Clearly this happens if and only if before each iteration of the while loop 3 in Figure 2, $(V', L')$ can be decomposed in terms of disjoint subtrees $V' = \sum_{k \in R'} V(k)$ and

$L' = \sum_{k \in R'} L(k)$. These subtrees may just be trivial ones $\mathcal{T}(k) = (\{k\}, \emptyset)$ comprising a root node $k$. We note also that these trees cover $R$, i.e. $R = \cup_{k \in R'} R(k)$. These properties hold before the first while loop, and hold subsequently since each loop of a successful reconstruction amalgamates binary subtrees rooted at siblings.

**Theorem 1** *Let $\mathcal{T}$ be a binary tree, equipped with an estimable and increasing function $\phi$. Then* DBT *reconstructs $\mathcal{T}$.*

**Proof of Theorem 1:** Suppose the algorithm does not reconstruct the tree. Then there must be an iteration of the while loop for which $j$ and $k$ in line 4 of Figure 2 are not siblings. Consider $R', V'$ at the start of the first loop that this occurs. Let $\ell$ be the sibling of $j$. $\ell \notin R'$ since $a(j, \ell) \prec a(j, k)$ implies $\Phi_{j,\ell} > \Phi_{j,k}$, contradicting the maximality of $\Phi_{j,k}$. Since the subtrees comprising $(V', L')$ are disjoint, no ancestor of $j$ (or hence of $\ell$) can lie in $R'$. Since the tree is binary, $\ell$ must have at least two descendents $t_1, t_2$ in $R'$ since otherwise $\cup_{r \in R'} R(r)$ would not cover $R$. Since $a(t_1, t_2) \prec \ell$, then $\Phi_{t_1, t_2} > \phi(\omega_{\boldsymbol{p}(\ell)}) > \phi(\omega_{\boldsymbol{p}(a(j,k))}) = \Phi_{j,k}$, contradicting the maximality of $\Phi_{j,k}$. ∎

**Reconstruction of Binary Trees from Measurements.** Now we switch to the context that a stream of probes is dispatched from the source, each giving rise to an independent realization of the mark process. Let $x^{(i)}$ denote the marks of the $i^{\text{th}}$ such realization. Each realization gives rise to a set of measurements $\{x^{(i)}_{\boldsymbol{p}(k)} : k \in R\}$ at the leaves. Suppose that some subtree $(V(k), L(k))$ of the tree is already identified. Then we can aggregate the measured leaf marks analogously to (2), defining $\widetilde{x}^{(i)}_k = x^{(i)}_{\boldsymbol{p}(k)}$ for $k \in R$, and forming $\widetilde{x}^{(i)}_k = \widetilde{x}^{(i)}_{h(k)} \oplus \widetilde{x}^{(i)}_{h^*(k)}$ by recursion for $k \notin R$.

Let $\widetilde{\omega}^{(n)}_k = n^{-1} \sum_{i=1}^n \delta_{\widetilde{x}^{(i)}_k}$ denote the empirical distribution of $\widetilde{x}^{(n)}_k$ ; here $\delta_y$ denotes the unit mass at $y \in \mathcal{X}$. $\widetilde{\omega}^{(n)}_k$ can be thought of as an discrete approximation to the true distribution $\widetilde{\omega}_k$, based on the measurements from $n$ probes. Correspondingly, we estimate $\mathcal{T}$ by the topology $\mathcal{T}^{(n)}$ obtained by using the $\widetilde{\omega}^{(n)}$ in place of $\widetilde{\omega}$ in the DBT algorithm. Specifically, we use $\Phi^{(n)}_{j,k} := \Phi(\widetilde{\omega}^{(n)}_{j,k})$ in place of $\Phi_{j,k}$ in line 3 of Figure 2. We call the resulting algorithm the Binary Tree Classification Algorithm (BT).

**Theorem 2** *Under the conditions of Theorem 1, with probability* 1, *$\mathcal{T}^{(n)} = \mathcal{T}$ for sufficiently large $n$. Hence $\mathcal{T}^{(n)}$ is a consistent estimator of $\mathcal{T}$ and $\lim_{n \to \infty} \mathsf{P}_\omega[\mathcal{T}^{(n)} \neq \mathcal{T}] = 0$.*

**Proof of Theorem 2:** By the law of large numbers, $\omega^{(n)}$ converges weakly to $\omega$, almost surely. Since $\phi$ is weakly continuous each $\Phi^{(n)}_{j,k}$ converges almost surely to $\Phi_{j,k}$. Then, almost surely, for all sufficiently large $n$, the relative ordering of the $\Phi^{(n)}_{j,k}$ is the same as that of the $\Phi_{j,k}$ for pairs $j, k$ for which the $\Phi_{j,k}$ are distinct. Hence BT reconstructs the tree

8

in the same manner as $\mathsf{DBT}$, except possibly varying the order of grouping amongst sets of sibling pairs $(j, k)$ with identical $\Phi_{j,k}$. The last two statements then follow by standard results. ∎

**Characterizing Link Behavior.**    In many of the examples of the next section $\phi$ is additive over links. i.e. $\phi(\omega_1 * \omega_2) = \phi(\omega_1) + \phi(\omega_2)$. Then we can ascribe a descriptor $\phi_k$, such as a packet loss rate or a mean delay, to each link $(f(k), k)$ through $\phi_k = \phi(\omega_{\boldsymbol{p}(k)}) - \phi(\omega_{\boldsymbol{p}(f(k))})$. (This may conveniently be done during the execution of $\mathsf{DBT}$ or $\mathsf{BT}$).

**Extension to General Trees.**    Inference of general trees is accomplished as follows. For simplicity assume that $\phi$ is additive. Then application of $\mathsf{DBT}$ to an arbitrary tree results in a binary tree that contains fictitious links $k$ such that $\phi_k = 0$. The tree can then be pruned by removing any such link and identifying its endpoints. It can be shown that this procedure yields the true general tree. In $\mathsf{BT}$ it is necessary to apply a threshold $\varepsilon > 0$ and prune all links $k$ with $\phi_k \leq \varepsilon$. This is because for finitely many probes, statistical fluctuations lead the characteristic $\phi$ of the fictitious links to differ slightly from zero. It can be shown that for sufficiently many probes, this approach reconstructs any general tree for which all $\phi_k > \varepsilon$.

## 3   Instances of Topology Inference

In this section we specify instances of the framework described above, specifying the setting (the marks $\mathcal{X}$, etc) and the the forms of the functions $\phi$, $\Phi$ and $\Phi^{(n)}$. Theorem 1 and 2 then apply immediately in each case.

### 3.1   Loss-Based Inference

For the case of loss we take $\mathcal{X} = \{0, 1\}$, where 0 indicates packet loss and 1 transmission. Composition is by taking minima $x_1 \otimes x_2 = x_1 \wedge x_2$ and the identity is $\boldsymbol{z} = 1$; a packet is transmitted on a path if it would be transmitted on all links of that path. Aggregation is by taking maxima $x_1 \oplus x_2 = x_1 \vee x_2$; hence $\widetilde{x}_k = 1$ if a packet reaches any receiver descended from $k$. It can be shown [1] that the path loss probabilities are related to aggregate loss probabilities at descendant nodes through

$$\mathsf{P}_\omega[x_{\boldsymbol{p}(a(j,k))} = 1] = \frac{\mathsf{P}_\omega[\widetilde{x}_j = 1]\mathsf{P}_\omega[\widetilde{x}_k = 1]}{\mathsf{P}_\omega[\widetilde{x}_j = \widetilde{x}_k = 1]}. \tag{3}$$

The generic distribution on $\mathcal{X}$ takes the form $(1-\alpha)\delta_0 + \alpha\delta_1$ for some $\alpha \in [0, 1]$. We define $\phi$ to act by extracting the negative log-weight of the state 1, thus: $\phi((1-\alpha)\delta_0 + \alpha\delta_1) =$

9

$-\log(\alpha)$. With $\boldsymbol{z} = 1$ then $\phi(\delta_{\boldsymbol{z}}) = 0$, since $\delta_{\boldsymbol{z}}$ is the generic measure parameterized by $\alpha = 1$.

Writing $\Phi_{j,k} = \log \mathsf{P}_\omega[\widetilde{x}_j \cdot \widetilde{x}_k = 1] - \log \mathsf{P}_\omega[\widetilde{x}_j = 1] - \log \mathsf{P}_\omega[\widetilde{x}_k = 1]$, we obtain the required relation between $\Phi$ and $\phi$ in Definition 1.

The characteristic $\phi$ is additive over links, and the link characteristic is $\phi_k = -\log \mathsf{P}_\omega[x_k = 1]$, i.e., the negative log probability of successful transmission over link $k$. Thus $\phi$ is increasing provided the link loss probabilities are strictly positive. Inference from measurements uses $\Phi_{j,k}^{(n)} = \log n + \log(\sum_{i=1}^n \widetilde{x}_j^{(i)} \widetilde{x}_k^{(i)}) - \log(\sum_{i=1}^n \widetilde{x}_j^{(i)}) - \log(\sum_{i=1}^n \widetilde{x}_k^{(i)})$ where we have expressed functions of the empirical measures $\widetilde{\omega}^{(n)}$ in terms of the corresponding empirical means.

## 3.2 Delay Covariance-Based Inference

In this case the increasing function $\phi$ is the variance of the cumulative delay from the root to a given node. In the formalism, $\mathcal{X} = \mathbb{R}_+$, with $x_k$ the delay encountered on link $k$. (The formalism extends to loss by using $x_k = \infty$ to denote loss; we treat this elsewhere [5]). Composition adds delays along a path: $x_1 \otimes x_2 = x_1 + x_2$. The identity is $\boldsymbol{z} = 0$. Aggregation takes the mean of two delays $x_1 \oplus x_2 = (x_1 + x_2)/2$; we note this aggregation is not associative. With $\mathcal{Q} = \mathbb{R}_+$ and $\phi(\omega) = \mathsf{Var}_\omega(x)$ we take

$$\phi(\omega_{\boldsymbol{p}(a(j,k))}) = \mathsf{Var}_\omega(x_{\boldsymbol{p}(a(j,k))}) = \mathsf{Cov}_\omega(\widetilde{x}_j, \widetilde{x}_k) = \Phi(\widetilde{\omega}_{j,k}). \tag{4}$$

The middle equality holds since, by the independence assumption, the only non-zero contribution to $\mathsf{Cov}_\omega(\widetilde{x}_j, \widetilde{x}_k)$ is due to delays on the common portion of the paths to $j$ and $k$. By the independence assumption $\phi$ is additive and so $\phi_k = \mathsf{Var}_\omega(x_k)$, the delay variance of link $k$. $\phi$ is increasing provided delays are not constant. $\Phi_{j,k}^{(n)}$ is the sample covariance, i.e., $\Phi_{j,k}^{(n)} = \frac{1}{n}\left(\sum_{i=1}^n \widetilde{x}_j^{(i)} \widetilde{x}_k^{(i)} - \frac{1}{n}\sum_{i=1}^n \widetilde{x}_j^{(i)} \cdot \sum_{i=1}^n \widetilde{x}_k^{(i)}\right)$.

## 3.3 Delay Distribution-Based Inference

With inference supplying the full distribution of the cumulative delay from the root to a given node, there are several choices of the increasing function $\phi$ available: the complementary cumulative distribution function (ccdf), the delay moments, and the delay variance.

In the formalism, $\mathcal{X} = \{q, 2q, \ldots, dq, \infty\}$, $q > 0, d \in \mathbb{N}$, with $x_k$ the delay on link $k$, discretized in bins of width $q$. $dq$ is a threshold delay above which packets are considered lost, $x_k$ taking the value $\infty$. Composition adds delays along a path: $x_1 \otimes x_2 = x_1 + x_2$. The identity is $\boldsymbol{z} = 0$. Aggregation takes minimum delay between paths $x_1 \oplus x_2 = x_1 \wedge x_2$; hence $\widetilde{x}_k = y$ if the minimum delay from the source to some receiver descended from $k$ is $y$. In [6] it was shown how a generalization of the approach for loss inference in Section 3.1 above can be used to express the discretized distribution $\omega_{\boldsymbol{p}(a(j,k))}$ of the delay

from the root to an interior node, in terms of the distribution $\widetilde{\omega}_{j,k}$ aggregate delays to leaf nodes descended through offspring $j, k$. More precisely, denote $A_k(i) = \mathsf{P}[x_{\boldsymbol{p}(k)} = iq]$, $\gamma_k(i) = \mathsf{P}[\widetilde{x}_{\boldsymbol{p}(k)} \leq iq]$, and $\gamma_{j,k}(i) = \mathsf{P}[\widetilde{x}_{\boldsymbol{p}(j)} \oplus \widetilde{x}_{\boldsymbol{p}(k)} \leq iq]$, $iq \in \mathcal{X}$. Then:

$$A_{a(j,k)}(0) = \frac{\gamma_j(0)\gamma_k(0)}{\gamma_j(0) + \gamma_k(0) - \gamma_{j,k}(0)} \tag{5}$$

and $A_{a(j,k)}(i)$, $i = 1, \ldots, d$, is recursively computed as the smallest solution of the following quadratic equation:

$$
\begin{aligned}
&\gamma_{j,k}(i) - A_{a(j,k)}(0) + A_{a(j,k)}(0) \cdot \prod_{\ell \in \{j,k\}} \left[ 1 - \frac{\gamma_\ell(i) - \sum_{m=1}^{i-1} \beta_\ell(i-m) A_{a(j,k)}(m) - \beta_\ell(0) A_{a(j,k)}(i)}{A_{a(j,k)}(0)} \right] \\
&+ \sum_{m=1}^{i-1} A_{a(j,k)}(m) \left\{ \prod_{\ell \in \{j,k\}} [1 - \beta_\ell(i-m)] - 1 \right\} \\
&+ A_{a(j,k)}(i) \left\{ \prod_{\ell \in \{j,k\}} [1 - \beta_\ell(0)] - 1 \right\} = 0
\end{aligned}
\tag{6}
$$

where $\beta_\ell(i) = \frac{\gamma_\ell(i) - \sum_{m=1}^{i} A_{a(j,k)}(m) \beta_\ell(i-m)}{A_{a(j,k)}(0)}$, $\ell \in \{j, k\}$.

In the first instance we can take $\mathcal{Q}$ as the set of ccdf's, equipped with pointwise partial order, arising from the delay distributions. Excluding from $\Omega$ trivial distributions in which all link delays are zero, then, since link delays are non negative, the map $\phi$ taking distributions to ccdf's is increasing in the sense of Definition 1(i).

In order to avoid comparing entire distributions, we can instead compare summary statistics. Since link delays are non-negative then any function of the form $\phi(\omega) = \mathsf{E}_\omega[h(x) \mid x < \infty]$ is estimable and increasing when $h$ is an increasing function, e.g., $h(x) = x^p, p > 0$. (Here $x$ represents a generic mark with distribution $\omega$.) A special case is the **delay average** estimator, obtained when $p = 1$. This is additive since the mean of the sum of two random variables is the sum of their means. Another estimator is the **delay variance** estimator: $\phi(\omega) = \mathsf{Var}_\omega[x | x < \infty]$. This is additive due to the independence of link delays.

For the delay average and variance classifiers, use $\Phi_{j,k} = \mathsf{E}_\omega[x_{\boldsymbol{p}(a(j,k))} | x_{\boldsymbol{p}(a(j,k))} < \infty]$ $= \sum_{i=0}^{d} iq A_{a(j,k)}(i) / \sum_{i=0}^{d} A_{a(j,k)}(i)$ and $\Phi_{j,k} = \mathsf{Var}_\omega[x_{\boldsymbol{p}(a(j,k))} \mid x_{\boldsymbol{p}(a(j,k))} < \infty] = \sum_{i=0}^{d}(iq)^2 A_{a(j,k)}(i) / \sum_{i=0}^{d} A_{a(j,k)}(i) - \mathsf{E}_\omega^2[x_{\boldsymbol{p}a(j,k)} \mid x_{\boldsymbol{p}(a(j,k))} < \infty]$, respectively, where we condition on the delay being finite. The corresponding $\Phi_{j,k}^{(n)}$ are computed using the estimated distribution $A_{a(j,k)}^{(n)}(i)$, computed through (5) and (6) using the estimates $\gamma_\ell(i)^{(n)} := n^{-1} \sum_{m=1}^{n} \mathbf{1}_{\{x_{\boldsymbol{p}(\ell)}^{(m)} \leq iq\}}$ and $\gamma_{j,k}(i)^{(n)} := n^{-1} \sum_{m=1}^{n} \mathbf{1}_{\{x_{\boldsymbol{p}(j)}^{(m)} \oplus x_{\boldsymbol{p}(k)}^{(m)} \leq iq\}}$ in place of $\gamma_\ell(i)$, $\ell \in \{j, k\}$, and $\gamma_{j,k}(i)$, $iq \in \mathcal{X}$.

## 3.4   Inference from Network Marking

We assume that a single bit mark can be set in a packet by network elements. Explicit Congestion Notification [11] is the example we have in mind, the mark being set to indicate

11

congestion on the path. The mark is not set at the source, but once set cannot be unset. This assumes that a mark on a packet incoming at a node is also set on all outgoing packets.

First assume that there is no packet loss. In the formalism we set $\mathcal{X} = \{0, 1\}$; the mark is left unchanged on link $k$ when $x_k = 1$, a mark is set (whether or not already set) when $x_k = 0$. With these conventions, the scheme is formally identical to that for packet loss in Section 3.1, when "packet lost" in the former case is identified with the event "mark set" in the present case. $e^{-\phi_k}$ is the probability that link $k$ would attempt to set the mark.

The two state model can be extended to packet loss by extending the 0 state to include packet loss as well as packet marked. One potential drawback of this approach is that packet marking and loss probabilities cannot be distinguished in the characteristic $\phi_k$. This can be achieved by extending to the following special case of the delay distribution formalism of Section 3.3. We parameterize the delay state space by $q = d = 1$, giving rise to the state space $\mathcal{X} = \{0, 1, \infty\}$, where the states $0$, $1$ and $\infty$ of $x_k$ are identified with the events "packet transmitted with mark unchanged", "packet transmitted with mark set (whether or not already set)" and "packet lost" on link $k$ respectively.

In the three-state model, each pairs of nodes $j, k$ gives rise to a characteristic $\Phi_{j,k}$ as a ccdf on $\mathcal{X}$ of the form $q = (q_0, q_1)$ with $q_0 > q_1$ where $q_0 = \mathsf{P}[x_{\boldsymbol{p}(a(j,k))} > 0]$ and $q_1 = \mathsf{P}[x_{\boldsymbol{p}(a(j,k))} > 1]$. Nodes $j, k$ are selected for pairing if $\Phi_{j,k}$ is maximal in the pointwise partial order.

The three state model assumes that loss and marking occur independently at different nodes. This excludes consideration of dependence of between loss and marking that could occur with spatially extended congestion, and also those dropping schemes in which marked packets are dropped selectively in downstream congestion.

### 3.5   Utilization-Based Inference

In this case the increasing function $\phi$ is (a function of) the probability of encountering minimal delay at all links in a path. This case can be regarded as a degenerate case of the delay distribution inference in which $\mathcal{X} = \{0, 1\}$ where $x_k = 0$ indicates no delay on link $k$, while $x_k = 1$ corresponds to any non-zero queueing delay. Hence the fraction of packets that experience $x_k = 1$ is a direct measure of the utilization of link $k$. Since $x_{\boldsymbol{p}(k)} = 0$ iff $x_j = 0$ for all $j$ in the path $\boldsymbol{p}(k)$, the setup maps exactly onto the loss inference described in Section 3.1, except with the roles of 0 and 1 interchanged.

## 4   Misclassification Analysis

In this section, we analyze the probabilities of misclassification for the various instances of BT, and estimate their convergence rates. For simplicity we consider here only the case that the characteristic $\phi$ takes values in $\mathbb{R}$.

Denote by $E_i$ the event that BT has correctly reconstructed the subtree rooted at node $i$, $i \in V$. Since the algorithm proceeds iteratively up the tree, $E_i$ requires first that both the subtrees rooted at its child nodes have been correctly reconstructed, then that its child nodes have been paired together. Therefore, for $i \in V \setminus R$ we can write

$$E_i \supseteq E_{h(i)} \cap E_{h^*(i)} \cap \left( \cap_{(j,k,l) \in \mathcal{S}(i)} Q(j,k,l) \right) \tag{7}$$

where $S(i) = \{(h(i), h^*(i), l), (h^*(i), h(i), l) | i, l \neq a(i,l)\}$ and $Q(j,k,l)$ is the event that

$$D^{(n)}(j,k,l) := \Phi_{j,k}^{(n)} - \Phi_{j,l}^{(n)} > 0 \tag{8}$$

holds. In $\cap_{(j,k,l) \in \mathcal{S}(i)} Q(j,k,l)$, $h(i)$ and $h^*(i)$ are grouped together to form node $i$ for all possible ways to reconstruct the tree. Denote by $E$ the event that the tree is correctly classified. From (7) we immediately have that $E \supseteq \cap_{i \in V \setminus R} \cap_{(j,k,l) \in \mathcal{S}(i)} Q(j,k,l)$. This provides the following upper bound for the misclassification probability, denoted by

$$P^f := \mathsf{P}[E^c] \leq \sum_{i \in V \setminus R} \sum_{(j,k,l) \in \mathcal{S}(i)} \mathsf{P}[Q^c(j,k,l)] \tag{9}$$

**Normal Approximations**   It can be shown that $\sqrt{n}(\Phi_{j,k}^{(n)} - \Phi_{j,k})$ has an asymptotically Gaussian distribution as the number of probes $n \to \infty$ in all instances described in Section 3. See [2, 5] for the loss and delay covariance case, for the other estimators it follows from Theorem 3 in [6].

**Theorem 3** *Under the conditions of Theorem 1, for each $i \in V \setminus R$, $\sqrt{n} \cdot (D^{(n)}(j,k,l) - D(j,k,l))$, $(j,k,l) \in \mathcal{S}(i)$, where $D(j,k,l) = \Phi_{j,k} - \Phi_{j,l}$, converges in distribution, as the number of probes $n \to \infty$, to a Gaussian random variable with mean 0 and variance $\sigma_D^2(j,k,l) = \lim_{n \to \infty} n \cdot \left( \mathsf{Var}(\Phi_{j,k}^{(n)}) + \mathsf{Var}(\Phi_{j,l}^{(n)}) - 2\mathsf{Cov}(\Phi_{j,k}^{(n)}, \Phi_{j,l}^{(n)}) \right)$.*

Theorem 3 suggests that for $(j,k,l) \in \mathcal{S}(i)$ we can approximate $\mathsf{P}[Q^c(j,k,l)]$ by $\Psi\left(-\sqrt{n} \cdot \frac{D(j,k,l)}{\sigma_D(j,k,l)}\right)$, where $\Psi$ is the cdf of the standard normal distribution. For large $n$, we have the following logarithmic asymptotic

$$n^{-1} \log \mathsf{P}[Q^c(j,k,l)] \sim -D^2(j,k,l)/2\sigma_D^2(j,k,l) \tag{10}$$

Since the largest term over $\cup_{i \in V \setminus R} \mathcal{S}(i)$ in (9) should dominate for large $n$, we expect the curve $\log P^f$ vs. $n$ to be asymptotically linear with negative slope

$$\frac{1}{2} \inf_{i \in V \setminus R} \inf_{(j,k,l) \in \mathcal{S}(i)} \frac{D^2(j,k,l)}{\sigma_D^2(j,k,l)} \tag{11}$$
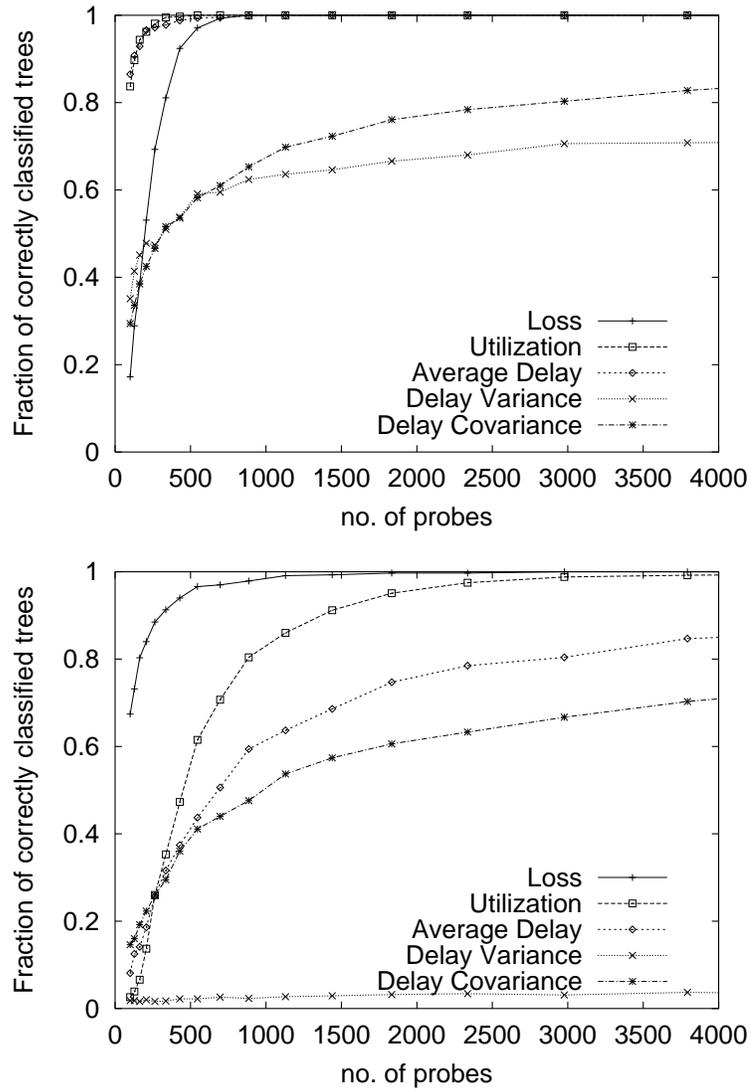
Figure 3: MODEL SIMULATION. Fraction of correctly classified topologies for different classifiers as function of the number of probes: (top) light load scenario; (bottom) heavy load scenario.

## 4.1 Misclassification Probability for the Different Classifiers

The foregoing analysis can be instantiated with the different estimators to obtain the estimated misclassification probability of the topology classifiers. In the following we compute the asymptotic behavior of the different classifiers by substituting the proper expressions in (11). In general, the calculation of the infimum in (11) is quite difficult since $\sigma_D^2(j,k,l)$ is a complex function of both the topology and the distribution $\omega$. Here, we capture the dominant modes of misclassification in asymptotic regime of small loss and delay. The results in [2] and Theorem 3 in [6] suggest that in this regime, the curve $\log P^f$ vs. $n$ is asymptotically linear with negative slope as follows:

1. for the loss based classifier

$$\frac{n}{2}\inf_{i\in V\setminus R}\mathsf{P}[x_i=0] \tag{12}$$

   The same expression applies to the two-state versions of the marking classifier.

2. for the utilization based classifier

$$\frac{n}{2}\inf_{i\in V\setminus R}\mathsf{P}[x_i=1] \tag{13}$$

3. for the average based classifier

$$\frac{n}{2}\inf_{\substack{i\in V\setminus R \\ k\in V\,|\,i,k\neq a(i,k)}}\frac{\left(\sum_{i\preceq j\prec a(i,k)}\mathsf{E}(x_i)\right)^2}{\sum_{i\preceq j\prec a(i,k)}\mathsf{E}(x_i^2)} \tag{14}$$

4. for the variance based classifier

$$\frac{n}{2}\inf_{\substack{i\in V\setminus R \\ k\in V\,|\,i,k\neq a(i,k)}}\frac{\left(\sum_{i\preceq j\prec a(i,k)}\mathsf{Var}(x_i)\right)^2}{\sum_{i\preceq j\prec a(i,k)}\mathsf{E}(x_i^4)} \tag{15}$$

For the covariance-based approach we used experience from experiments to identify the events that dominate misclassification. As $n$ increases, it was found that the most likely way to misclassify a tree is by incorrectly identifying the link with the smallest link variance. When it comes to grouping the children of the terminating node $k$, of such a link, this results is mistakenly grouping one of the children with the sibling of $k$. This suggests the following approximation for the covariance approach

$$P^f\approx e^{-(n/2)\frac{\mathsf{Var}^2(x_j)}{\sigma_D^2(h(i),h^*(i),j)}}, \tag{16}$$

where $j=\arg\min_{i\in V\setminus R}\mathsf{Var}(x_i)$.

15

# 5 Simulation Evaluation and Algorithm Comparison

In this section we compare the performance of the different classification algorithms through two types of simulation. In *model simulations* delay and loss are chosen to follow our statistical model, allowing us to test algorithm performance in the setting on which our analysis is based. *Network simulations*, using the `ns` [8] simulator, test the algorithms in a more realistic setting, where delay and loss are due to queueing delay and buffer overflows at nodes as multicast probes compete with background TCP/UDP traffic. We did not conduct simulation for the case of marking with ECN, since we did not have a model for endpoint congestion control using the marks in the multicast case.

## 5.1 Model Simulation

In the model simulations, at each link a probe is either lost, or encounters no delay, or suffers an exponentially distributed delay. We conducted 1000 simulations over randomly generated 15 node binary trees. In Figure 3 we plot the fraction of correctly classified topologies as a function of the number of probes for the different classifiers. We considered two regimes: a light load regime with low loss (1%) and utilization (randomly chosen between 10% and 40%), and a heavy load regime with higher loss (randomly chosen between 1% and 20%) and utilization (randomly chosen in between 30% and 80%). In both cases, mean delays were randomly chosen between 0.2ms and 2ms. We adopted a delay granularity of 1ms. (Although this seems large compared with mean delays, the delay predictions are quite accurate; see [6]).

The loss based classifier is found to be most accurate in general. The exception is with small numbers of probes at small loss rates. Then rare losses provide insufficient data points, and accuracy is greater for the utilization and average delay classifiers. The utilization based classifier has the best accuracy among the delay distribution based classifiers. This was expected because the delay average and variance approaches use estimates of the entire delay distribution while the utilization approach uses estimates of only the first bin; estimation of the weights of lower bins being found to be more accurate. Similarly, delay average is more accurate than delay variance since it attaches less weight to higher delay bins.

From the plots, the utilization based approach appears to work very well with low links utilizations, while its performance degrades with higher utilizations, which is in contrast with (13). This can be explained by observing that the theorem holds for the limit behavior as utilization goes to 0; in our experiments, we found that (13) captures well the misclassification behavior for utilization up to $20\%$. On the other hand, as link utilizations increase, the number of events used by the algorithm, namely those of minimum end-to-end delay, decreases rapidly. This results in increased estimator variance.

The two best performing algorithm, (loss and utilization based) have the smallest com-
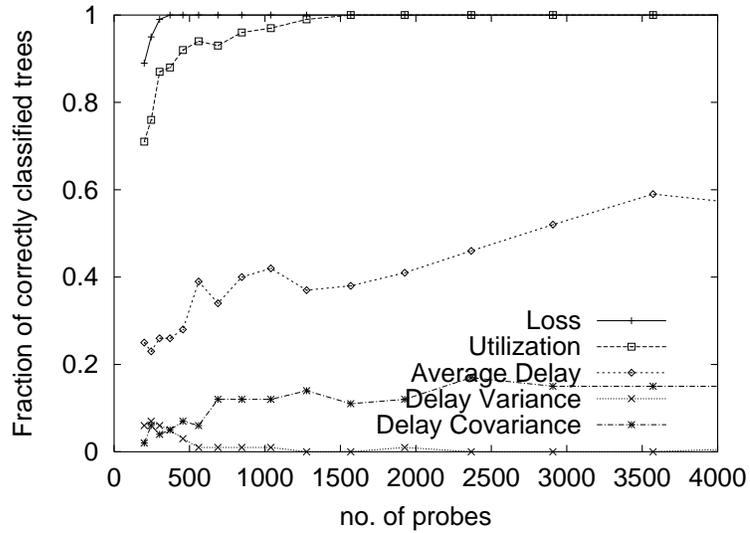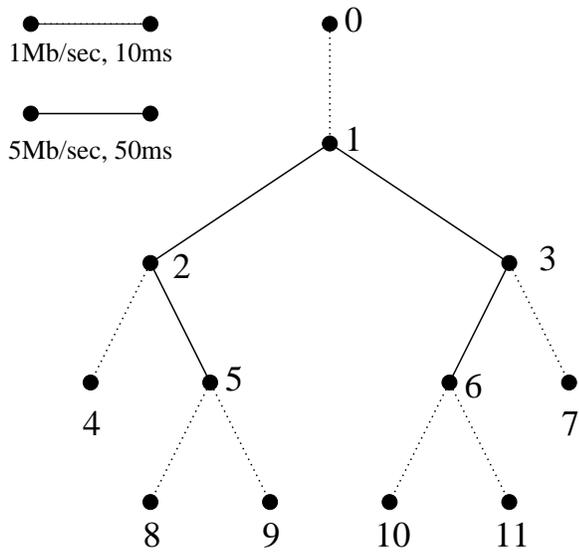
Figure 4: ns SIMULATION: (top) simulation topology; (bottom) fraction of correctly classified topologies for different classifiers as function of the number of probes

putational complexity. All algorithms require $O(\#R^3)$ node pairs computations. Each of these is $O(n)$ for the loss, utilization and covariance based estimators. These are considerably more complex for the delay average and variance classifiers since the whole delay distribution must be calculated, by recursively solving quadratic equations, the number of which is inversely proportional to the bin size $q$.

## 5.2  TCP/UDP Network Simulation

The `ns` simulations used the topology shown in Figure 4(top). To capture the heterogeneity between edges and core of a WAN, interior links have higher capacity (5Mb/sec) and propagation delay (50ms) than at the edge (1Mb/sec and 10ms). Each link is modeled as a FIFO queue with a 4-packet capacity.

   The root node 0 generates probes as a 20Kbit/s stream comprising 40 byte UDP packets according to a Poisson process with a mean interarrival time of 16ms; this represents $2\%$ of the smallest link capacity. The background traffic comprises a mix of infinite data source TCP connections (FTP) and exponential on-off sources using UDP. Averaged over the different simulations, the link loss ranges between $1\%$ and $11\%$ and link utilization ranges between $20\%$ and $60\%$. The average delay ranges between 1 and 2ms for the slower links and between 0.2 and 0.5ms for the faster links. The delay distributions were computed using a bin size of 1ms.

   In Figure 4(bottom) we plot the fraction of correctly identified topologies over 100 simulations. The relative accuracy among the different classifiers is in agreement with the results from the model simulation with the loss based algorithm having the best performance with no misclassification for more than 500 probes. The rather poor performance of the delay based algorithms, with the exception of the utilization classifiers, is largely due to the presence of spatial correlation. In our simulations, a multicast probe is more likely to experience a similar level of congestion on consecutive links or on sibling links than is dictated by the independence assumption. This has negative impact on the accuracy of the delay estimates which accounts for the observed performance.

   We also observed temporal correlation among successive probes that encountered the same congestion events. However, it can be shown that the presence of short-term correlation does not affect estimator consistency, although the convergence rate may be slowed.

## 6  Conclusions

In this paper we have presented a general framework for the inference of the multicast tree topologies from end-to-end measurements. In contrast with tools such as `mtrace` [7], cooperation of intervening network nodes is not required.

   We specified an algorithm which reconstructs the topology of multicast tree in presence of any packet performance measure that: (i) monotonically increases as the packet traverses

down the tree; and (ii) can be estimated on the basis of end-to-end measurements at the receivers. Building on previous results in [1, 5, 6], we were able to specify several instances of this algorithm based on the performance measures of packet loss, link utilization, delay average and delay variance. The scheme can also be used to infer topology based on marks set in packets by Explicit Congestion Notification, subject to some assumptions on how this would be extended to multicast. This last case is potentially useful for congestion control schemes that use marks rather than loss or delay to indicate congestion to endpoints. In this case, loss and delay may be too infrequent to infer topology reliably.

We investigated the statistical properties of the algorithms, and showed that, under mild assumptions, they are consistent and computed their convergence rate. We evaluated our classifiers though simulation. We found out that the two algorithms with the lowest computational complexity, namely, the loss based and the utilization based algorithm, also have the best performance, with the loss based algorithm being in general the most accurate except when the number of probes and the loss rate are both small. Moreover, both algorithms seemed to be robust and exhibit good convergence in real traffic simulations, in spite of violation of the independence assumption of our model.

Finally, the algorithms described in this paper are each based on a different performance metric. We have recently extended this work by formulating algorithms which take advantage of the available measurements by integrating the different performance metrics we have here separately considered; see [4].

# References

[1] R. Caceres, N.G. Duffield, J.Horowitz and D. Towsley, "Multicast-Based Inference of Network Internal Loss Characteristics", IEEE Trans. on Information Theory, November 1999.

[2] R. Caceres, N.G. Duffield, J.Horowitz F. Lo Presti and D. Towsley, "Statistical Inference of Multicast Network Topology", in Proc. 1999 IEEE Conf. on Decision and Control, Phoenix, AZ, Dec. 1999.

[3] Cooperative Association for Internet Data Analysis, "Internet Measurement Efforts," `http://www.caida.org/Tools/taxonomy.html#InternetMeasurement`

[4] N.G. Duffield, J. Horowitz, F. Lo Presti, "Adaptive Multicast Topology Inference", submitted for publication, July 2000.

[5] N.G. Duffield and F. Lo Presti, "Multicast Inference of Packet Delay Variance at Interior Network Links", Proceedings IEEE Infocom 2000, Tel Aviv, March 2000.

[6] F. Lo Presti, N.G. Duffield, J.Horowitz and D. Towsley, "Multicast-Based Inference of Network-Internal Delay Distributions", submitted for publication, September 1999.

[7] `mtrace` – Print multicast path from a source to a receiver. See `ftp://ftp.parc.xerox.com/pub/net- research/ipmulti`

[8] ns – Network Simulator. See `http://www-mash.cs.berkeley.edu/ns/ns.html`

[9] S. Paul, et al. "Reliable multicast transport protocol (RMTP)", *IEEE JSAC*, Vol. 15, No. 3, pp. 407–421, April 1997.

[10] V. Paxson, J. Mahdavi, A. Adams, M. Mathis, "An Architecture for Large-Scale Internet Measurement," *IEEE Communications*, Vol. 36, No. 8, pp. 48-54, August 1998.

[11] K. Ramakrishnan and S. Floyd, "A Proposal to add Explicit Congestion Notification (ECN) to IP," RFC 2481, January 1999.

[12] S. Ratnasamy & S. McCanne, "Inference of Multicast Routing Tree Topologies and Bottleneck Bandwidths using End-to-end Measurements", Proc. IEEE Infocom'99, New York, NY (1999)