

Technical Report TREE-0109: Maximum Likelihood Identification of Network Topology from End-to-End Measurements

Rui Castro, Mark Coates and Robert Nowak*

May 3, 2002

Abstract

One of the predominant schools of thought in networking today is that monitoring and control of large scale networks is most practical at the edge. Identifying or estimating the routing topology is crucial to edge-based approaches. The focus of this work is a new Maximum Likelihood criterion for topology identification that makes use only of unicast measurements performed between host computers and requires no special support (e.g., ICMP responses) from internal elements. The measurement procedure used is based on a unicast probing strategy, and uses delay difference measures, not requiring clock synchronization. Theoretical analysis of the likelihood criterion leads to a characterization of the maximum likelihood tree, and simplifies significantly the inference task. Based on that characterization we develop a new, fast algorithm for topology identification. Simulation results as well as actual Internet measurements are presented, illustrating the potentials of our technique.

1 Introduction

The rapid growth of the Internet, combined with fast and unpredictable developments in applications and workloads, has rendered network modelling and control increasingly demanding tasks. One of the predominant schools of thought in networking today is that management and control of large scale

*Rui Castro and Robert Nowak are with the Department of Electrical and Computer Engineering, Rice University, Houston, Texas, USA. Mark Coates is with the Department of Electrical and Computer Engineering, McGill University, Canada. Corresponding Author: Rui Castro, E-mail: rcastro@rice.edu.

networks is only practical at the edge. Effective monitoring of the network internal state is crucial for edge-based control. One approach to gathering local information is to augment the network structure with special-purpose devices that collect and transmit local performance data.

In most cases this is highly impractical, requiring expensive, special purpose hardware and software at each router, raising privacy concerns and involving an enormous communication overhead because of the tremendous amount of information that must be transmitted to end-systems.

The alternative is to indirectly infer the network characteristics from edge-based measurements, with minimal cooperation from the network structure.

In this paper we consider the problem of topology inference. The knowledge of the topology and connectivity maps is a fundamental variable in the state of a network. Several groups have begun investigating methods for inferring of internal network behavior based solely on end-to-end measurements [1, 2, 3] . This methods require knowledge of the network topology, so, a first step of these approaches is the determination of the network topology, preferably with minimal cooperation of the network. Most existing tools for topology mapping, like `traceroute`, rely on the close cooperation of internal routers, and can only reveal areas where the network is functioning properly and prepared to reveal itself. These conditions are often not met and one expects them to become more uncommon as the Internet grows, due to its unregulated environment.

The focus of this paper is a new maximum likelihood approach to topology identification. We consider a single source, communicating with multiple receivers. The network topology can be represented as a graph, where each vertex represents a network device, and the edges represent the connection between two such devices (each connection may also include devices not represented in the graph). Each network device, also called element or node, corresponds to a physical point of the network with certain properties, namely, where traffic branches out, or where exists a bandwidth limitation (for example a switch). We assume that the routes from the sender to the receivers are fixed during the measurement period. In that case the topology is a tree-structured graph.

Related work: Methods for topology identification have been proposed for multicast network trees. These techniques where first developed for loss-based inference [4], relying on the fact that multicast receivers sharing a longer portion of that path from the source experience higher shared

losses. This work has been further generalized in [5, 6], defining the concept of metric-based topology identification, making the approach applicable to different types of measurements. Specific metrics have been proposed, based on loss rates, utilization (frequency of zero delay), mean delay and delay covariance. In [5] it was noted that loss and utilization based metrics provided the best results, but under different network load situations. In [7] an adaptive scheme was developed to take advantage of different measurement techniques. The loss-based multicast methods were adapted to unicast measurements in [8]. In [6] the topology identification problem was also formulated as a maximum likelihood estimation, but the study was restricted to loss based inference.

The previous methods depend on either loss or delay measurements. Under normal operating conditions, packet losses are rare, making the loss-based methods perform poorly. The delay-based inference techniques present a number of advantages, but also involve the challenging problem of accurate timing and synchronization of the clocks at the sender and receivers. In most cases, accurate timing and synchronization demands GPS access, restricting the application of these methods. Some advances are being made towards PC-based clock synchronization and accurate delay measurement [9], but the capabilities are far from widespread. In this paper we use a measurement scheme developed in [10] that relies only on delay differences, overcoming the synchronization issues.

Contribution: Most of the methods above rely on estimates of performance metrics from the end-to-end measurements. The measurements are noisy and that influences the quality of topology identification. We propose instead a Maximum Likelihood Estimation (MLE) approach that takes into account the uncertainty of the measurements.

The second contribution of this paper is a characterization of the Maximum Likelihood Tree (MLT). Finding the maximum likelihood tree involves the solution of a constrained continuous optimization problem, and a discrete optimization problem. We derive some properties of the maximum likelihood tree, and based on that characterization we formulate the maximum likelihood estimation again as two optimization problems, but now the continuous optimization is not constrained, and it can be solved very efficiently.

Determining the globally optimal tree through an exhaustive search of all the possible topologies is not practical. Even for a relatively small number of receivers, the number of potential topologies is

huge. This means that identification of the globally optimal tree by simple inspection is an extremely hard problem, and in most cases demands a prohibitive computational expense. The characterization of the maximum likelihood tree motivates a new, improved (but suboptimal) bottom-up algorithm based on the likelihood formulation. The method we propose can easily incorporate deterministic side information about the topology, making it suitable also for topology identification when partial topology information is available. As an example, such information might be generated by a portion of the network responding to `traceroute`.

The paper is organized as follows. In Section 2 we present the main modelling assumptions and framework. In Section 3 we formulate the topology identification as a maximum likelihood estimation problem. In Section 4 we introduce the measurement procedure. A characterization of the maximum likelihood tree is presented in Section 5. In Section 6, we describe a greedy inference algorithm motivated by the MLE approach, and a possible extension allowing incorporation of side information (obtained using `traceroute` for example). In Section 7 we present some simulation results, comparing the performance of the greedy inference algorithm described in Section 6 with that of previously described techniques. In section 8 we present the results of an Internet measurement experiment, demonstrating the applicability of our techniques in a realistic setting. Finally, in Section 9 we make some concluding remarks and indicate directions for future work.

2 Model and Framework

Consider a source transmitting information through a network to a set of receivers. Assume that the routes from the sender to the receivers are fixed. The problem we address is the identification of the network topology based on end-to-end measurements that measure the degree of correlation between receivers [4, 5, 6, 8, 10]. With this limited information, it is only possible to identify the so-called logical topology, defined by the branching points between paths to different receivers. This corresponds to a tree-structured topology with the sender at the root and the receivers at the leaves, as depicted in Figure 1.

Let $\mathcal{T} = (V, L)$ denote a logical tree with nodes V and links L . Denote the source node (the root

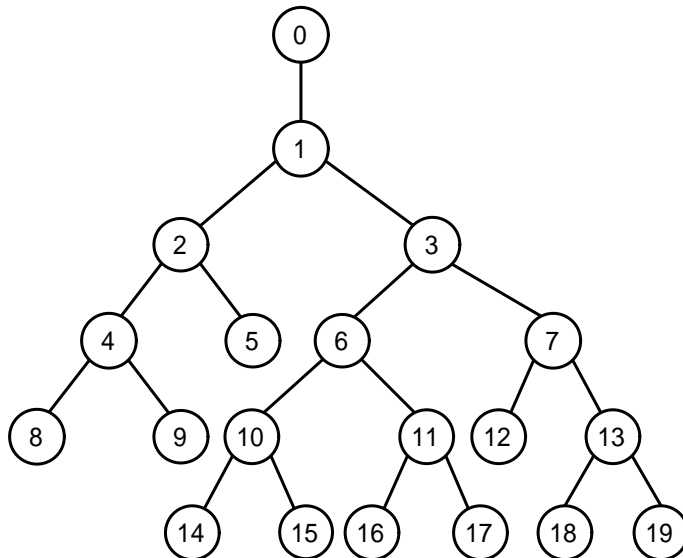


Figure 1: A binary logical tree topology.

of the tree) by 0. Identify the leaf nodes with the set of receivers R . Every node has at least two descendants, apart from the root node (which has one, denoted by 1), and the leaf nodes (which have none). If all internal nodes have exactly two descendants then the tree is called binary. For each internal node let $f(i)$, $i \in V \setminus \{0\}$, denote the parent of i (for example, in Figure 1, $f(10) = 6$). We can identify each link with its end node, i.e., $(f(i), i) \sim i$, $(f(i), i) \in L$. Denote by \mathcal{P}_i the path from the source node to i , e.g., $\mathcal{P}_{12} = \{0, 1, 3, 7, 12\}$. Let $a(i, j)$, $i, j \in V$, denote the nearest ancestor of the pair of nodes (i, j) , e.g., $a(15, 17) = 6$. We define also $d(i)$, $i \in V$, as the set of children nodes of i . For a given node k we denote by $\mathcal{T}(k)$ the subtree rooted at k , and by $R(k)$ the set of receivers of $\mathcal{T}(k)$.

For each node in the tree we can associate a metric γ_k , $k \in V$. The metric value is related to the extent of the (unique) path from the root to node k . The considered metrics must have the following

Monotonicity Property: Let $i, j \in V$ be any two nodes and let $\mathcal{P}_i, \mathcal{P}_j$ denote the paths from the source to i and j respectively. If \mathcal{P}_i is a proper subpath of \mathcal{P}_j then $\gamma_i \leq \gamma_j$.

For each pair of receivers $i, j \in R$ we can associate a metric value, $\gamma_{a(i,j)}$. For simplicity we consider

the notation $\gamma_{ij} \equiv \gamma_{a(i,j)}$. The value of γ_{ij} is a measure of the "length" of the shared portion of the paths to i and j .

Knowledge of the metric values for each pair of receivers and above monotonicity property are sufficient for identification of the underlying topology. For example, referring to Figure 1, the metric $\gamma_{18,19}$ will be greater than $\gamma_{i,19}$ for all $i \in R \setminus \{18,19\}$, revealing that receivers 18 and 19 have a common parent in the logical tree. The property can be exploited in this manner to devise simple and effective bottom-up merging algorithm that identifies the complete, logical topology [4, 5, 6, 8].

Metrics possessing the Monotonicity Property can be estimated from a number of different end-to-end measurements including counts of losses, counts of zero delay events (utilization), and delay correlations [4, 5, 6, 8]. These estimated metrics, denoted $\{x_{ij}\}$, $i, j \in R$, can be interpreted as statistics derived from repeated measurements. Randomness in network conditions leads to variability in the measurements and hence variability in the estimated metrics. Most of the previous work in this area does not incorporate this variability (which can be assessed directly from the measurements) into the identification process. We claim that this variability can have a major impact on the performance of topology identification algorithms.

Several methods have been proposed for topology identification in both unicast and multicast settings [4, 5, 6, 8], but all have a very similar structure. The DBT algorithm proposed in [5] is a representative example. The algorithm is a recursive selection and merging/aggregation process that generates a binary tree from the bottom-up (receivers to sender). In this paper, we describe a new algorithm that specifically addresses the uncertainty in estimated metrics, providing substantial performance improvement in certain cases.

3 Likelihood formulation

To address the issue of metric variability and uncertainty we pose topology identification as a maximum likelihood estimation problem. We select the MLE approach for its well known asymptotic optimality properties.

The estimated metrics $\mathbf{x} \equiv \{x_{ij} : i, j \in R\}$ can be interpreted as observations of the true metric

values $\boldsymbol{\gamma} \equiv \{\gamma_k : k \in V \setminus \{0, R\}\}$ contaminated by some randomness or noise. We model this contamination probabilistically. The estimated metrics are randomly distributed according to a density (whose precise form depends on the contamination model) that is parameterized by the underlying topology \mathcal{T} and the set of true metric values, written as $p(\boldsymbol{x}|\boldsymbol{\gamma}(\mathcal{T}), \mathcal{T})$. The metric values are constrained by the tree \mathcal{T} , but in order to simplify the mathematical expressions we will usually not indicate this dependence. The estimated metrics \boldsymbol{x} are observed and hence fixed, and when $p(\boldsymbol{x}|\boldsymbol{\gamma}, \mathcal{T})$ is viewed as a function of \mathcal{T} and $\boldsymbol{\gamma}$ it is called the likelihood of \mathcal{T} and $\boldsymbol{\gamma}$. The maximum likelihood tree is given by

$$\mathcal{T}^* = \arg \max_{\mathcal{T} \in \mathcal{F}} \max_{\boldsymbol{\gamma}(\mathcal{T}) \in \mathcal{G}} p(\boldsymbol{x}|\boldsymbol{\gamma}(\mathcal{T}), \mathcal{T}), \quad (1)$$

where \mathcal{F} denotes the *forest* of all possible tree topologies connecting the sender to the receivers and \mathcal{G} denotes the set of all metrics satisfying the Monotonicity Property, that is

$$\mathcal{G} = \left\{ \boldsymbol{\gamma} \in \mathbb{R}^{\#W} : \gamma_{f(k)} \leq \gamma_k, \forall k \in W \right\}, \quad W = V \setminus \{0, 1, R\}. \quad (2)$$

The set $W \in V$ corresponds to all the internal links, that is, all the links except the root link and the receiver links. In rigor we should write $W(\mathcal{T})$, but for notational simplicity we do not usually show this dependence explicitly.

Occasionally we refer to the pair $(\mathcal{T}, \boldsymbol{\gamma}(\mathcal{T}))$ as a tree, including both the topology and metric values. Together with the measurement values \boldsymbol{x} , $(\mathcal{T}, \boldsymbol{\gamma}(\mathcal{T}))$ completely defines the likelihood value.

4 Sandwich Measurements

To illustrate our approach we will focus on one type of metric. In earlier work we proposed a topology identification method based on delay differences [10], that method relied on a measurement scheme called sandwich probing. We briefly describe the measurement technique, details can be found in [10]. Each sandwich probe consists of three packets, and gives information about the shared path among two receivers. Figure 2 shows the details of the probing scheme.

The metrics used are the mean delay differences. Since we only need local delay differences (at

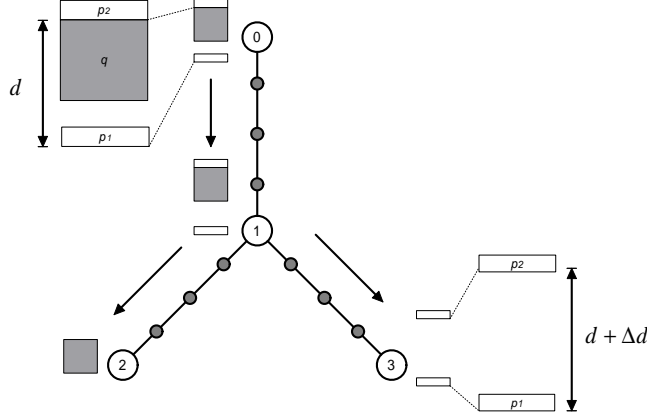


Figure 2: An example of sandwich probe measurement. The large packet is destined for node 2, the small packets for node 3. The black circles on the links represent physical queues where no branching occurs. In the absence of cross-traffic, the initial spacing between the small probes d is increased along the shared path from nodes 0 to 1 because the second small probe p_2 queues behind the large packet. The measurement $x_{2,3}$ for this receiver pair is equal to $d + \Delta d$. A larger initial spacing d reduces the chance of p_2 catching p_1 because of a bottleneck or cross-traffic on the path from node 1 to 3.

the sender and the receiver), there is no need for clock synchronization between the source host and the receivers. In the case where there is no cross traffic the measurement Δd is directly related to the bandwidth of the shared queues. We assume the cross-traffic is stationary, and that the initial spacing of the two small packets d is large enough so that neither the large packet nor the second small packet queue behind the first small packet at any time. We send the each probe far apart in time, so that we can assume that the outcomes of different measurements are independent.

The delay differences provide (noisy) measurements of a metric related to the number of shared queues in the paths to two receivers. Let x_{ij} be the sample mean of repeated delay difference measurements for pair $i, j \in R$. Under reasonable assumptions the measurements are statistically independent and have finite variance, hence, according to the Central Limit Theorem, the distribution of each empirical mean tends to a Gaussian. This motivates the following (approximate) model:

$$x_{ij} \sim \mathcal{N}(\gamma_{ij}, \sigma_{ij}^2), \quad (3)$$

where σ_{ij}^2 is sample variance of the measurements divided by the number of measurements ($\sigma^2 \equiv \{\sigma_{ij}^2\}$), x_{ij} is the sample mean of the measurements, and $\mathcal{N}(\gamma, \sigma^2)$ denotes the Gaussian density with

mean γ and variance σ^2 . Notice that we are not assuming that the delay differences are normally distributed, but only their empirical means. Under the above assumptions, as the number of measurements increases, the model accuracy increases. Thereon we will refer to \mathbf{x} and σ^2 as estimated metrics and estimate variances, respectively.

The likelihood function in this case is a product of Gaussian densities, one factor for each pair of receivers. We assume that the each sandwich probe measurements is statistically independent of the others the measurements. The log likelihood becomes then

$$\log p(\mathbf{x}|\boldsymbol{\gamma}(\mathcal{T}), \mathcal{T}) = - \sum_{i \in V} \sum_{j \in V \setminus \{i\}} \frac{(x_{ij} - \gamma_{ij}(\mathcal{T}))^2}{2\sigma_{ij}^2} + C, \quad (4)$$

where $C \in \mathbb{R}$ is a constant. Notice that $\log p(\mathbf{x}|\boldsymbol{\gamma}(\mathcal{T}), \mathcal{T})$ is a concave functional of $\boldsymbol{\gamma}(\mathcal{T})$. Notice also that there is no assumption of spatial independence of delay in the network.

We can also consider a different parameterization for the metric, motivated by the additivity of the delay differences. For each link k we can associate a metric value θ_k . Let $\boldsymbol{\theta} \equiv \{\theta_k\}$. The relationship between parameters γ_i and $\boldsymbol{\theta}$ is

$$\gamma_i = \sum_{k \in \mathcal{P}(i)} \theta_k.$$

The parameters $\boldsymbol{\theta}$ are called link-levels parameters. Note that the constraint set written in terms of parameters $\boldsymbol{\theta}$ is

$$\mathcal{G} = \left\{ \boldsymbol{\theta} \in \mathbb{R}^{\#W} : \theta_k \geq 0, \forall k \in W \right\}.$$

5 Characterization of the Maximum Likelihood Tree

The maximizations involved in (1) are quite formidable. We are not aware of any method for computing the global maximum except by a brute force examination of each tree in the forest. Consider a network with N receivers. A very loose lower bound on the size of the forest \mathcal{F} is $N!/2$. For example, if $N = 10$ then there are more than 1.8×10^6 trees in the forest. This explosion of the search space precludes the brute force approach in all but very small (logical) networks. Moreover, the inner maximization is non-trivial because it involves a constrained optimization over \mathcal{G} . The following theorem

establishes a key property of the MLE solution that leads to an important simplification.

Theorem 1. *Let $\tilde{\mathcal{T}}$ be a tree such that*

$$\arg \max_{\gamma \in \mathbb{R}^{\#W}} \log p(\mathbf{x}|\gamma, \tilde{\mathcal{T}}) \neq \arg \max_{\gamma \in \mathcal{G}} \log p(\mathbf{x}|\gamma, \tilde{\mathcal{T}}). \quad (5)$$

Then another tree (\mathcal{T}, γ) , $\gamma \in \mathcal{G}$, can be explicitly constructed from $\tilde{\mathcal{T}}$ such that

$$\log p(\mathbf{x}|\gamma, \mathcal{T}) > \max_{\gamma \in \mathcal{G}} \log p(\mathbf{x}|\gamma, \tilde{\mathcal{T}}). \quad (6)$$

In particular, if \mathcal{T}^ is the solution to (1), i.e., the MLT, we have*

$$\arg \max_{\gamma \in \mathbb{R}^{\#W}} \log p(\mathbf{x}|\gamma, \mathcal{T}^*) = \arg \max_{\gamma \in \mathcal{G}} \log p(\mathbf{x}|\gamma, \mathcal{T}^*). \quad (7)$$

Remark 1: Consider an arbitrary tree $\tilde{\mathcal{T}}$. Suppose that the vector γ maximizing $\log p(\mathbf{x}|\gamma, \tilde{\mathcal{T}})$ changes depending on consideration of constrained ($\gamma \in \mathcal{G}$) or unconstrained ($\gamma \in \mathbb{R}^{\#W}$) maximization, that is expression (5) holds. The theorem says that in that case we can construct another tree (\mathcal{T}, γ) from $\tilde{\mathcal{T}}$ such that it yields a higher likelihood than any tree $(\tilde{\mathcal{T}}, \gamma)$, $\gamma \in \mathcal{G}$. Consequently the tree $\tilde{\mathcal{T}}$ cannot be the maximum likelihood tree (7).

Remark 2: The second part of the theorem (7) shows that it is unnecessary to perform the constrained optimization. For each tree, we can compute the unconstrained optimization, which simply involves calculating a weighted sum of metrics, and check if the resulting maximizer lies in the set \mathcal{G} . The computation of the parameters γ can be done in a simple bottom-up procedure.

The set of trees that satisfy that property in the theorem is defined as

$$\mathcal{F}' = \left\{ \mathcal{T} \in \mathcal{F} : \arg \max_{\gamma(\mathcal{T}) \in \mathbb{R}^{\#W(\mathcal{T})}} p(\mathbf{x}|\mathcal{T}, \gamma(\mathcal{T})) \in \mathcal{G} \right\}.$$

Note that the maximum likelihood tree belongs to this set, i.e., $\mathcal{T}^* \in \mathcal{F}'$.

Remark 3: From the proof technique (in Appendix A) we also note that we need only to consider binary trees, because for any non-binary tree we can construct a corresponding binary tree yielding

the same likelihood value, obtained from the former adding links with metric value zero. Therefore, without loss of generality, we can consider only binary trees in the forest \mathcal{F} .

It's worth noticing that this result can be applied to any search technique on the forest, searching only among trees in \mathcal{F}' . This can simplify the computations and perhaps improve performance of the Monte Carlo methods proposed in [10], by reducing the size of the search space.

6 Likelihood-based Binary Tree (LBT) Algorithm

Determining the globally optimal tree through exhaustive searching over the possible topology trees is not practical, even for a small number of receivers. Recall that even with 10 receivers there are over 1.8×10^6 of possible trees. The theorem above motivates a new, improved (but suboptimal) bottom-up merging algorithm based on our likelihood formulation of the problem. The new approach, called LBT algorithm, is partially motivated by the following observation. Consider the problem of identification of an unknown network topology with four receivers (receivers 1, 2, 3 and 4), and suppose that it was known that two receivers have a common parent node, for example receivers 1 and 2. Can we rewrite the initial topology identification problem as a similar problem for three receivers, where one of the receivers corresponds to the aggregation of receivers 1 and 2? The answer to this question is yes. In fact the more general result below is true.

Consider an arbitrary tree with receiver set R and the sets $A, B \subseteq R$, with $A \cap B = \emptyset$. Define

$$\begin{aligned} N(A, B) &= \sum_{i \in A} \sum_{j \in B} \frac{x_{ij}}{\sigma_{ij}^2} \\ D(A, B) &= \sum_{i \in A} \sum_{j \in B} \frac{1}{\sigma_{ij}^2}. \end{aligned}$$

Proposition 1. *Let $\mathcal{F}(R)$ denote the forest of trees with receiver set R . Let \mathbf{x}, σ^2 be the estimated metrics corresponding to a tree with receiver set R . Let $R_0 \subset R$ be a subset of receivers. Define a subset \mathcal{F}' of the forest \mathcal{F} , corresponding to the trees such that there exists an internal node k satisfying*

$R(k) = R_0$. In other words, k is the root of a subtree with receiver set R_0 . Formally

$$\mathcal{F}' = \{\mathcal{T} \in \mathcal{F}(R) : \exists_{k \in V(\mathcal{T})} \text{ such that } R(k) = R_0\}.$$

Aggregate all the receivers in R_0 into a single node r , denoted by aggregated receivers.

Then

$$\arg \max_{\mathcal{T} \in \mathcal{F}'} \max_{\boldsymbol{\gamma} \in \mathbb{R}^{\#W(\mathcal{T})}} p(\mathbf{x}, \boldsymbol{\sigma}^2 | \boldsymbol{\gamma}, \mathcal{T}) = \arg \max_{\mathcal{T} \in \mathcal{F}(R')} \max_{\boldsymbol{\gamma} \in \mathbb{R}^{\#W(\mathcal{T})}} p(\mathbf{x}', \boldsymbol{\sigma}'^2 | \boldsymbol{\gamma}, \mathcal{T}), \quad (8)$$

where $R' = (R \setminus R_0) \cup \{r\}$, and \mathbf{x}' , $\boldsymbol{\sigma}'^2$ are given by

$$x'_{ir} = \frac{N(\{i\}, R_0)}{D(\{i\}, R_0)}, \quad x'_{ri} = \frac{N(R_0, \{i\})}{D(R_0, \{i\})}, \quad i \in R' \setminus \{r\} \quad (9)$$

$$\sigma'^2_{ir} = \frac{1}{D(\{i\}, R_0)}, \quad \sigma'^2_{ri} = \frac{1}{D(R_0, \{i\})}, \quad i \in R' \setminus \{r\}. \quad (10)$$

$$x'_{ij} = x_{ij}, \quad \sigma'^2_{ij} = \sigma^2_{ij}, \quad i, j \in R' \setminus \{r\}.$$

This proposition tells us how we can merge measurements in order to preserve the likelihood structure. It states that solving the problem (1) (without constraints on the estimated metrics) conditional on the fact that some subtree of receivers is fixed, is equivalent to solving a smaller problem, where the fixed subtree with receiver set R_0 is replaced by a single node (the aggregated receiver r), and the measurements are aggregated according to (9) and (10). The proof of Proposition 1 is tedious but not difficult, and follows from equation (12).

Proposition 1 suggests a similar approach to the DBT algorithm proposed in [7] with some key modifications, to account for the variability of the measurements.

We start with an unknown topology, and a set aggregated receivers R' (initially set $R' = R$, the set of physical receivers). We then find the pair $\{i, j\} \subset R'$ that maximizes

$$\left(\frac{x_{ij}}{\sigma^2_{ij}} + \frac{x_{ji}}{\sigma^2_{ji}} \right) / \left(\frac{1}{\sigma^2_{ij}} + \frac{1}{\sigma^2_{ji}} \right). \quad (11)$$

We infer that the receivers i, j have a common parent, denoting it by k . This way of choosing a pair enforces the monotonicity property (we choose the pair with the largest estimated metric). Now, using

LBT Algorithm

1. **Input:** Set of receivers R , measurements $\{x_{ij}\}$ and $\{\sigma_{ij}^2\}$, $i, j \in R$.
2. **Initialization:** $R' := R$ and $V = R$.
3. Find the pair (i, j) such that

$$(i, j) = \arg \max_{\substack{i, j \in R' \\ i \neq j}} \left(\frac{x_{ij}}{\sigma_{ij}^2} + \frac{x_{ji}}{\sigma_{ji}^2} \right) / \left(\frac{1}{\sigma_{ij}^2} + \frac{1}{\sigma_{ji}^2} \right).$$

4. Denote by k the new aggregated receiver. Set $V := V \cup \{k\}$, $R' := R' \cup \{k\} \setminus \{i, j\}$. Set also $f(i) := k$; $f(j) := k$. Compute the aggregated measurements

$$\begin{aligned} x_{kl} &= \frac{x_{il} \cdot \sigma_{jl}^2 + x_{jl} \cdot \sigma_{il}^2}{\sigma_{il}^2 + \sigma_{jl}^2} & x_{lk} &= \frac{x_{li} \cdot \sigma_{lj}^2 + x_{lj} \cdot \sigma_{li}^2}{\sigma_{li}^2 + \sigma_{lj}^2} \\ \sigma_{kl}^2 &= \frac{\sigma_{il}^2 \cdot \sigma_{jl}^2}{\sigma_{il}^2 + \sigma_{jl}^2} & \sigma_{kl}^2 &= \frac{\sigma_{il}^2 \cdot \sigma_{jl}^2}{\sigma_{il}^2 + \sigma_{jl}^2}, \quad \text{where } l \in R' \setminus \{k\}. \end{aligned}$$

5. If $|R'| > 1$ go back to 3.
6. Set $V := V \cup \{0\}$ and $f(k) = 0$.
7. **Output:** The node set V and the parent function $f : V \setminus \{0\} \rightarrow V$.

Figure 3: LBT algorithm

the results from the proposition we reduce the size of the receiver set, i.e., $R' \leftarrow R' \cup \{k\} \setminus \{i, j\}$. The aggregation of the measurements follows from (9) and (10). The new (aggregated) metric estimate value relating node k to any other node $l \in R'$ is given by

$$x_{kl} = \left(\frac{x_{il}}{\sigma_{il}^2} + \frac{x_{jl}}{\sigma_{jl}^2} \right) / \left(\frac{1}{\sigma_{il}^2} + \frac{1}{\sigma_{jl}^2} \right),$$

and similarly for x_{lk} . The variances for the new aggregated metrics must be updated in a similar fashion

$$\sigma_{kl}^2 = (\sigma_{il}^2 \sigma_{jl}^2) / (\sigma_{il}^2 + \sigma_{jl}^2),$$

and similarly for σ_{lk}^2 . This procedure is iterated until there is only one aggregated receiver left (i.e. $|R'| = 1$). This procedure is outlined in Figure 3 in algorithmic format. The algorithm outputs the nodes set V and the parent function $f : V \rightarrow V$, that uniquely defines a tree.

Notice that the number of aggregated receivers is decreased by one at each step of the algorithm, guaranteeing the algorithm to stop. This algorithm ensures that the estimated tree satisfies the monotonicity property. It can be shown, using a similar reasoning as in [7], that this procedure yields a consistent estimator of the true topology tree, in the sense that the probability of correct identification of the true tree approaches one as the measurement variances tend to zero.

Aside from two simple, yet crucial modifications to the aggregation and decision step, that account for the metric estimates variance (i.e. confidence), the algorithm operates in the same manner as the DBT algorithm. The performance improvements provided by these modifications are examined in the next section.

6.1 Incorporating Side-Information

Due to its structure, the LBT algorithm is particularly well suited to incorporate deterministic side information. Such information can be provided for example by `traceroute`. Using `traceroute` one can get an estimate of the routing topology (each node in estimated topology corresponds to a router). Missing links are due to hosts that either did not send ICMP (Internet Control Message Protocol) time exceeded messages or sent them with a TTL (Time-to-live) too small to reach the `traceroute` source. Using the sandwich probing procedure, we can refine the above logical topology estimate by probing only parts of the network where we have incomplete information. Hence, we need only to send probes to a subset of pairs of receivers, instead of all possible pairs. That can reduce significantly the traffic overhead induced by probing, and make the overall process faster.

7 Simulation Results

We conducted some simple simulations in order to compare the performance of the MLE, the LBT and the DBT approaches. The DBT was designed for multicast topology identification, but can also be applied to unicast topology identification. In the DBT algorithm the "metrics" used are simply the average of each pair of empirical means (i.e. $\hat{\gamma}_{ij} = (x_{ij} + x_{ji})/2$), and empirical variances are not employed.

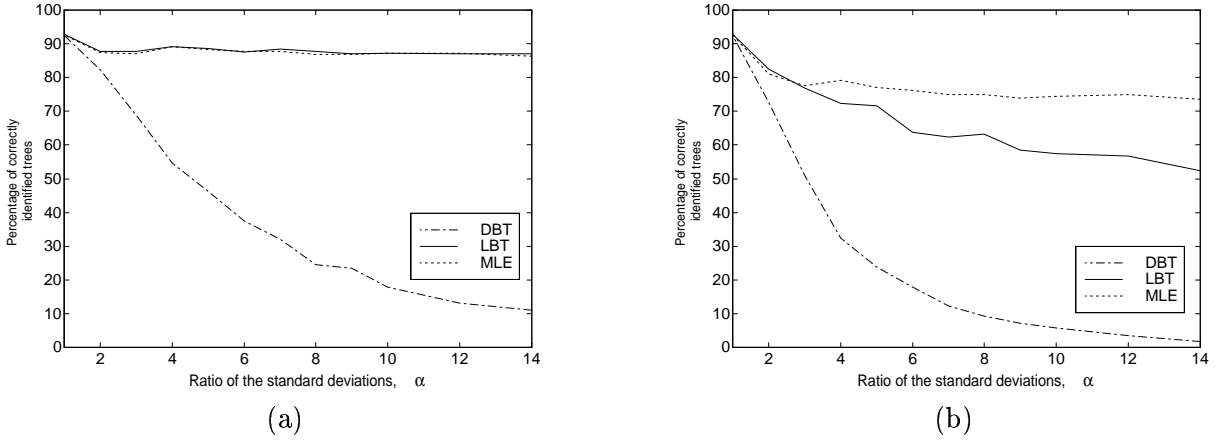


Figure 4: (a) Performance results considering different measurement variance for one receiver (b) Performance results considering different measurement variance for two receiver

To contrast the characteristics of the DBT and LBT algorithms we performed some simple simulation experiments. These simulations are not an exhaustive performance comparison of the algorithms, but are intended to indicate the benefits and drawbacks of the various methods.

We consider a randomly chosen six receiver topology, such that the link link-level parameters θ are the same for all links ($\theta_k = 1$, for all $k \in W$). For each pair of receivers $(i, j) \in R$ we generate 100 independent measurements with variance $100\sigma_{ij}^2$. In the first experiment we let $\sigma_{ij}^2 = \sigma^2/100$ for all $i \neq 1$, and $\sigma_{1j}^2 = \alpha^2\sigma^2/100$, where $\sigma^2 = 25$ and $\alpha > 1$. In other words, the measurements taken at any receiver have all the same variance, except for receiver 1 (since the trees are chosen randomly, we can consider a particular receiver, without loss of generality). This illustrates the case where the link corresponding to one of those receivers has problems, due to congestion or any other factor. For each value of α we evaluate the three methods on 1000 randomly generated trees.

In Figure 4(a) we plot the performance of the three methods as a function of the standard deviation ratio α . As we can observe, the MLE and LBT performance are very similar, and remain almost constant for the range of α considered. On the other hand, the performance of the DBT algorithm degrades significantly as α increases. The reason for this is that the metrics $\hat{\gamma}_{1j}$ in the DBT are strongly affected by the highly variable empirical means x_{1j} , for $j \in R$. In the case of the LBT, those empirical means do not affect the performance because the measurements are weighted according to their variance. Hence, for high values of α , the measurements x_{1j} have little impact on the estimated

tree.

In Figure 4(b) we plot the results of an experiment similar to the previous one, but the measurements taken at two different receivers have higher variance. Precisely $\sigma_{ij}^2 = \sigma^2/100$ for $i \neq \{1, 2\}$ and $\sigma_{1j}^2 = \sigma_{2j}^2 = \alpha^2\sigma^2$ (as before we can consider a fixed pair of receivers, without loss of generality). In this case we observe that the performance of the three methods degrades as α increases. The DBT algorithm exhibits the same trends as before, degrading considerably as α increases, but we note also the LBT performance degrades considerably as compared to the MLE, although it is still much better than that of the DBT. The reason for this is that the estimate of γ_{12} has a high variance, since both σ_{12}^2 and σ_{21}^2 are large, for large α . There is a higher probability of mispairing receivers 1 and 2, since there is a higher chance of choosing that particular pair in the greedy decision step. Unlike the LBT, the MLE algorithm is able to account the higher variance of those measurements, with respect to all the other measurements, consequently improving its performance.

The second experiment shows that, although the LBT algorithm provides a significant improvement over the DBT, its greedy nature (based on local decisions), makes it perform poorly (relative to the MLE) in certain scenarios. On the other hand, the MLE algorithm estimates a tree by solving a global optimization problem. The drawback is the heavy complexity that a search over all possible trees has. In [10] a search strategy was developed based on Monte-Carlo methods. This technique is computationally more demanding, but it is not greedy. Given enough computation time this algorithm will find the actual MLT. Currently the authors have some ongoing work on deterministic search strategies, exploiting the characterization in Theorem 1.

8 Internet Experiments

We have implemented a software tool called `nettom` that performs sandwich probing measurements and estimates the topology of a tree-structured network. The software has been implemented as a set of C programs using the Unix socket library and the programs have been ported to Solaris, FreeBSD, and (some) Linux platforms. The topology estimation (data collection and inference) is performed at a source host. The program at the source sends UDP sandwich probes to a set of remote clients,

which are required to run a low overhead receiver task during the measurement period. The receiver task primarily time-stamps a received UDP small packet (with the local time) and returns the time-stamped packet to the sender via a dedicated TCP connection. The sender software maintains a log of the returned packets, and at the completion of the measurement period, calculates the delay differences, the associated metrics and their variability.

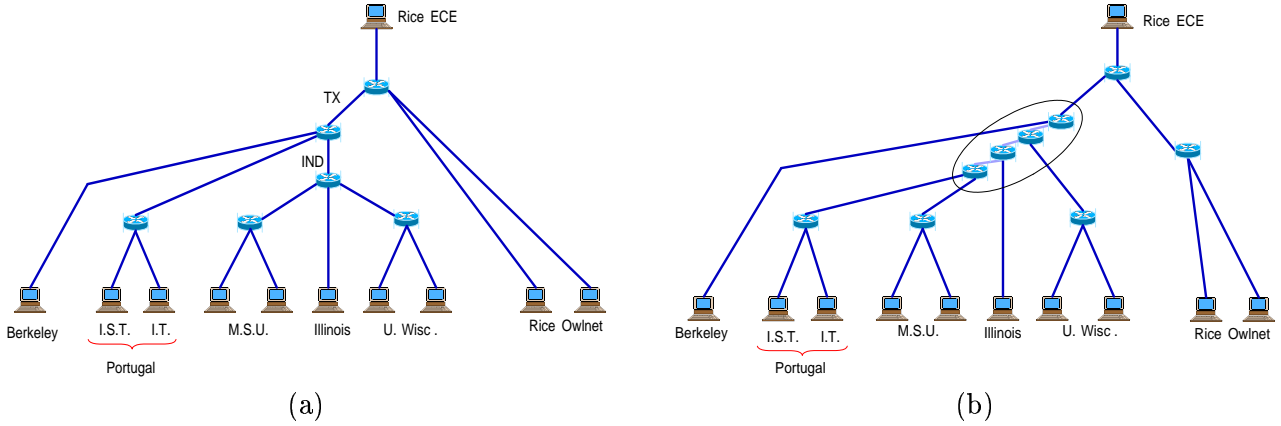


Figure 5: (a) The topology of the network used for Internet experiments, obtained using traceroute. (b) Estimated topology using the LBT algorithm. The signaled links have link-parameter values θ one order of magnitude smaller than all the other links. Those links can be collapsed, that is, the three routers inside the circle can be identified as one unique router.

We conducted Internet experiments using several hosts in the United States and abroad. The topology inferred from `traceroute` is depicted in Figure 5(a). The source for the experiments was located at Rice University. There are ten receiver clients, two located on different networks at Rice, two at separate hosts in Portugal, and six located at four other US universities.

The experiment was conducted for a period of eight minutes, during which a sandwich probe was sent to a randomly chosen receiver-pair once every 50 ms. Without any loss, the maximum number of probes available is 8600. This corresponds to less than 200 probes per pair, hence the traffic overhead on any link is very low.

We applied the LBT algorithm to the measurements collected and the result is depicted in Figure 5(b). Since the procedure is suited only for binary trees, it adds some extra links with small metric value ($\theta \approx 0$ for those links, indicating possibly high bandwidth links). Typically those links have link-level parameters one order of magnitude smaller than all the other ones, hence, if we col-

lapse those links, identifying they end-nodes with each other, we get a non-binary tree that provides a more truthful estimate of the network. Using this pruning procedure we get a very good estimate of the true network topology, although it fails to detect the backbone connection between Texas and Indianapolis. We expect that the latter connection is very high speed and the queuing effects on the constituent links are too minor to influence measurements. The estimated topology also places an extra element shared between the Rice computers. Although that element is not a router, hence it is not shown in the topology estimate using `traceroute`, it corresponds to a real physical device. To the best of our knowledge the detected element is a bandwidth limitation device.

We also applied the DBT algorithm to the measurements, but the performance was very poor. The reason for that behavior is that the measurements taken at the I.S.T. host had a very high variance and, as we observed in Section 7 that increases the probability of misclassification of the DBT.

9 Conclusions and Future Work

In this paper we considered the problem of network topology identification using end-to-end measurements. We formulated the problem in a probabilistic setting, as a maximum likelihood estimation. We particularize this approach considering a measurement technique introduced in [10]. This probing scheme is delay-based, but only requires local delay difference measurements at each receiver host, hence, no clock synchronization is needed among the different hosts. We showed that within this setting the maximum likelihood tree has some interesting properties, simplifying considerably the inference task. Those properties also motivate a new, likelihood-based, bottom-up procedure, called LBT. It has low complexity and, although sub-optimal, performs well under a variety of scenarios. One important difference of the LBT, compared to most of the previous methods in the literature, is that it accommodates differences in the variability of the measurements, making it more robust than the previous techniques in certain adverse scenarios.

We illustrated how the developed method can be used in conjunction with deterministic side information, obtained for example using `traceroute`. The use of this kind of structure can reduce

the probing time and overhead, making our techniques suitable to be used in large scale topology identification problems.

Some Internet experiments were conducted to evaluate the performance of the methods developed in a realistic setting. The results obtained were promising and demonstrate the potential of our techniques.

So far we only considered binary tree topologies. In [5] a pruning technique is suggested. The main idea is that at the first stage one estimates the binary topology that best describes the measured data, and then at the second stage links that are not significant are eliminated. Although this is asymptotically optimal there is a danger associated with forcing a binary tree structure when the true topology is not a binary tree. One way to overcome this problem is to consider a penalized likelihood approach, penalizing complex topologies. This approach was pursued in [10].

The solution of the maximum likelihood estimation (1) is still an open problem. In [10] a Markov Chain Monte Carlo method was developed to search the space of possible tree topologies in an efficient way. Currently the authors are investigating deterministic forest search methods yielding fast algorithms to find the maximum likelihood tree.

Another direction for future work is considering a multiple receiver setting. In that case the topology is not a tree. This raises a number of problems that might be overcome by slight changes in the measurement procedure.

A Proof of theorem 1

The proof is done by construction. Given a tree $\tilde{\mathcal{T}}$ satisfying (5) we construct a tree \mathcal{T} with larger likelihood.

Let $\tilde{\mathcal{T}}$ be a tree satisfying (5). Notice that since (4) is a concave functional, the RHS of (5) is given by the solution of the equation

$$\nabla_{\boldsymbol{\gamma}} \log p(\mathbf{x}|\boldsymbol{\gamma}, \tilde{\mathcal{T}}) = \mathbf{0}, \tag{12}$$

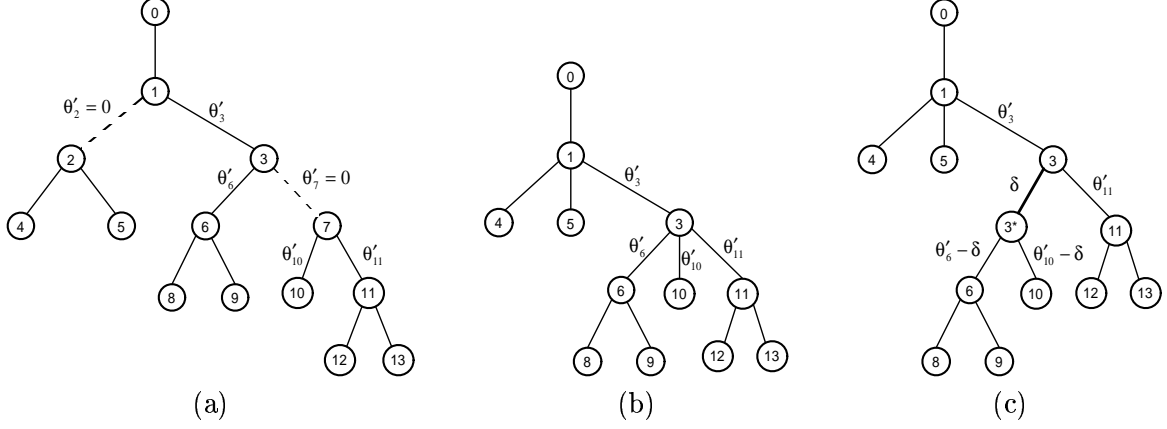


Figure 6: (a) Original tree $(\tilde{\mathcal{T}}, \tilde{\gamma})$ (b) Collapsed tree (\mathcal{T}', γ') (c) Constructed tree (\mathcal{T}, γ) .

where ∇ denotes the gradient operator. Let $\tilde{\gamma}$ be

$$\tilde{\gamma} = \arg \max_{\gamma \in \mathcal{G}} \log p(\mathbf{x} | \gamma, \tilde{\mathcal{T}}). \quad (13)$$

Again using the concavity of the log-likelihood we have that

$$\nabla_{\gamma} \log p(\mathbf{x} | \gamma, \tilde{\mathcal{T}}) \Big|_{\gamma = \tilde{\gamma}} \neq \mathbf{0}.$$

In particular, there exists a link $l \in V(\tilde{\mathcal{T}}) \setminus R(\tilde{\mathcal{T}})$ such that $\tilde{\gamma}_l = \tilde{\gamma}_{f(l)}$ and

$$\frac{\partial}{\partial \gamma_l} \log p(\mathbf{x} | \gamma, \tilde{\mathcal{T}}) \Big|_{\gamma = \tilde{\gamma}} \neq 0.$$

Since $\tilde{\gamma}_l = \tilde{\gamma}_{f(l)}$ we have that $\tilde{\theta}_l = 0$.

Computing the gradient of (4) we then get

$$\begin{aligned} \frac{\partial}{\partial \gamma_l} \log p(\mathbf{x} | \gamma, \tilde{\mathcal{T}}) \Big|_{\gamma = \tilde{\gamma}} &= \sum_{i \in d_{\tilde{\mathcal{T}}}(l)} \sum_{j \in d_{\tilde{\mathcal{T}}}(l) \setminus \{i\}} N(R_{\tilde{\mathcal{T}}}(i), R_{\tilde{\mathcal{T}}}(j)) \\ &- \tilde{\gamma}_l \left[\sum_{i \in d_{\tilde{\mathcal{T}}}(l)} \sum_{j \in d_{\tilde{\mathcal{T}}}(l) \setminus \{i\}} D(R_{\tilde{\mathcal{T}}}(i), R_{\tilde{\mathcal{T}}}(j)) \right] \neq 0. \end{aligned} \quad (14)$$

Consider now the tree obtained collapsing all the links such that $\tilde{\theta}_l = 0$ and keeping the value

of the remaining link-level parameters unchanged. Denote it by (\mathcal{T}', γ') (see Figure 6(a)(b)). Note that the parameters γ' satisfy the constraints (2), and that (\mathcal{T}', γ') yields the same likelihood value as $(\tilde{\mathcal{T}}, \tilde{\gamma})$, that is

$$\log p(\mathbf{x}|\tilde{\gamma}, \tilde{\mathcal{T}}) = \log p(\mathbf{x}|\gamma', \mathcal{T}'). \quad (15)$$

The assumption on the maximality of the log likelihood of $(\tilde{\mathcal{T}}, \tilde{\gamma})$ (from (13)) yields

$$\gamma'(\mathcal{T}') = \arg \max_{\gamma \in \mathcal{G}(\mathcal{T}')} \log p(\mathbf{x}|\gamma, \mathcal{T}').$$

From (15) (\mathcal{T}', γ') yields the maximum log likelihood for a tree \mathcal{T}' , and γ' is in the interior of the constraint set (note that the constraint set \mathcal{G} is a closed set). Hence we have [11]

$$\gamma'(\mathcal{T}') = \arg \max_{\gamma \in \mathbb{R}^{\#W(\mathcal{T}')}} \log p(\mathbf{x}|\gamma, \mathcal{T}'). \quad (16)$$

We now return to expression (14). Since $\tilde{\theta}_l = 0$, the link l was collapsed and does not appear in the tree (\mathcal{T}', γ') . There exists an ancestor k of node l (in $\tilde{\mathcal{T}}$) such that $\tilde{\gamma}_k = \tilde{\gamma}_l$ and the corresponding link in $\tilde{\mathcal{T}}$ was not collapsed in \mathcal{T}' (notice that k is not necessarily the parent of l). Hence there exists a pair $(i, j) \in d_{\mathcal{T}'}(k)$ of children nodes of k in \mathcal{T}' such that

$$\frac{N(R_{\mathcal{T}'}(i), R_{\mathcal{T}'}(j)) + N(R_{\mathcal{T}'}(j), R_{\mathcal{T}'}(i))}{D(R_{\mathcal{T}'}(i), R_{\mathcal{T}'}(j)) + D(R_{\mathcal{T}'}(j), R_{\mathcal{T}'}(i))} \neq \gamma'_k. \quad (17)$$

It is easy to show the (17) implies that

$$\frac{N(R_{\mathcal{T}'}(l), R_{\mathcal{T}'}(m)) + N(R_{\mathcal{T}'}(m), R_{\mathcal{T}'}(l))}{D(R_{\mathcal{T}'}(l), R_{\mathcal{T}'}(m)) + D(R_{\mathcal{T}'}(m), R_{\mathcal{T}'}(l))} > \gamma'_k. \quad (18)$$

Using the chosen pair (18), we construct a new tree \mathcal{T} (refer to Figure 6(c)) adding an extra node k^* and the link (k, k^*) . Node k^* has children l and m . Loosely speaking we are pulling the pair of nodes l and m down, adding a new node k^* . The parameter values for this new tree, denoted by γ , are adjusted such that $\gamma_{k^*} = \gamma'_k + \delta$, $\gamma_l = \gamma'_l - \delta$, and $\gamma_m = \gamma'_m - \delta$, $\delta > 0$. All the other metric values remain the same. Note that $\delta > 0$, but small enough so that the tree (\mathcal{T}, γ) still satisfies the

constraints (2). In terms of the link-level parameters $\boldsymbol{\theta}$ we have $\theta_{k^*} = \delta$, $\theta_l = \theta'_l - \delta$ and $\theta_m = \theta'_m - \delta$.

The log likelihood of $(\mathcal{T}, \boldsymbol{\gamma})$ is identical of the one from $(\tilde{\mathcal{T}}, \tilde{\boldsymbol{\gamma}})$, except for the term involving γ_{k^*} , thus

$$\begin{aligned}
\log p(\mathbf{x}|\boldsymbol{\gamma}, \mathcal{T}) - \log p(\mathbf{x}|\boldsymbol{\gamma}', \mathcal{T}') &= \sum_{i \in R_{\mathcal{T}'(l)}} \sum_{j \in R_{\mathcal{T}'(m)} \setminus i} \left[\frac{(x_{ij} - \gamma'_k)^2}{2\sigma_{ij}^2/n_{ij}} - \frac{(x_{ij} - \gamma_{k^*})^2}{2\sigma_{ij}^2/n_{ij}} \right] \\
&= \sum_{i \in R_{\mathcal{T}'(l)}} \sum_{j \in R_{\mathcal{T}'(m)} \setminus i} \left[\frac{(x_{ij} - \gamma'_k)^2}{2\sigma_{ij}^2/n_{ij}} - \frac{(x_{ij} - \gamma'_k - \delta)^2}{2\sigma_{ij}^2/n_{ij}} \right] \\
&= \sum_{i \in R_{\mathcal{T}'(l)}} \sum_{j \in R_{\mathcal{T}'(m)} \setminus i} \left[\frac{x_{ij} - \gamma'_k}{\sigma_{ij}^2/n_{ij}} \delta + o_{ij}(\delta) \right], \text{ as } \delta \rightarrow 0 \\
&= \sum_{i \in R_{\mathcal{T}'(l)}} \sum_{j \in R_{\mathcal{T}'(m)} \setminus i} [N(R_{\mathcal{T}'(l)}(i), R_{\mathcal{T}'(m)}(j)) - \gamma'_k D(R_{\mathcal{T}'(l)}(i), R_{\mathcal{T}'(m)}(j))] \delta + o(\delta), \text{ as } \delta \rightarrow 0 \\
&> 0, \quad \text{for small enough } \delta > 0.
\end{aligned} \tag{19}$$

In the above, $o(\delta)$, as $\delta \rightarrow 0$, is the standard notation for a function such that $\lim_{\delta \rightarrow 0} o(\delta)/\delta = 0$. Notice that the conclusion in step (19) is due to (18). Hence, for a suitable value of δ , we have $\log p(\mathbf{x}|\boldsymbol{\gamma}, \mathcal{T}) > \log p(\mathbf{x}|\boldsymbol{\gamma}', \mathcal{T}')$, and hence

$$\log p(\mathbf{x}|\boldsymbol{\gamma}, \mathcal{T}) > \log p(\mathbf{x}|\boldsymbol{\gamma}^*, \tilde{\mathcal{T}}).$$

For the second part of the theorem note that, if $\tilde{\mathcal{T}}$ is the MLT, then, together with $\tilde{\boldsymbol{\gamma}}$ as defined above, we have that $(\tilde{\mathcal{T}}, \tilde{\boldsymbol{\gamma}})$ yields the maximum likelihood and hence (5) doesn't hold, yielding (7). This concludes the proof.

References

- [1] R. Cáceres, N. Duffield, J. Horowitz, and D. Towsley, "Multicast-based inference of network-internal loss characteristics," *IEEE Trans. Info. Theory*, vol. 45, no. 7, pp. 2462–2480, November 1999.
- [2] F. Lo Presti, N.G. Duffield, J. Horowitz, and D. Towsley, "Multicast-based inference of network-internal delay distributions," Tech. Rep. CMPSCI 99-55, University of Massachusetts, 1999.

- [3] M. Coates and R. Nowak, “Sequential Monte Carlo inference of internal delays in nonstationary data networks,” to appear in *IEEE Trans. Signal Processing, Special Issue on Monte Carlo Methods for Statistical Signal Processing*, 2002.
- [4] S. Ratnasamy and S. McCanne, “Inference of multicast routing trees and bottleneck bandwidths using end-to-end measurements,” in *Proceedings of IEEE INFOCOM 1999*, New York, NY, March 1999.
- [5] N.G. Duffield, J. Horowitz, F. Lo Presti, and D. Towsley, “Multicast topology inference from end-to-end measurements,” in *ITC Seminar on IP Traffic, Measurement and Modelling*, Monterey, CA, Sep. 2000.
- [6] N.G. Duffield, J. Horowitz, F. Lo Presti, and D. Towsley, “Multicast topology inference from measured end-to-end loss,” *IEEE Trans. Info. Theory*, vol. 48, no. 1, pp. 26–45, January 2002.
- [7] N.G. Duffield, J. Horowitz, and F. Lo Presti, “Adaptive multicast topology inference,” in *Proceedings of IEEE INFOCOM 2001*, Anchorage, Alaska, April 2001.
- [8] A. Bestavros, J. Byers, and K. Harfoush, “Inference and labeling of metric-induced network topologies,” Tech. Rep. BUCS-2001-010, Computer Science Department, Boston University, June 2001.
- [9] Pásztor A. and Veitch D., “PC based precision timing without GPS,” in *Proc. ACM Sigmetrics 2002*, Marina del Rey, California, July 2002.
- [10] R. Castro, M.J. Coates, M. Gadhiok, R. King, R. Nowak, E. Rombokas, and Y. Tsang, “Maximum likelihood network topology identification from edge-based unicast measurements,” Tech. Rep. TREE0107, Department of Electrical and Computer Engineering, Rice University, Oct. 2001.
- [11] M. R. Hestenes, *Optimization Theory: The Finite Dimensional Case*, John Wiley & Sons, Inc., New York, 1975.