# Routing in Overlay Multicast Networks

Sherlia Y. Shi         Jonathan S. Turner

Department of Computer Science

Washington University in St. Louis

{*sherlia, jst*}*@cs.wustl.edu*

*Abstract*— **Multicast services can be provided either as a basic network service or as an application-layer service. Higher level multicast implementations often provide more sophisticated features, and since they don't require network support for multicast, they can provide multicast services, where no network layer support is available. *Overlay multicast networks* offer an intermediate option, potentially combining the flexibility and advanced features of application layer multicast with the greater efficiency of network layer multicast. Overlay multicast networks play an important role in the Internet. Indeed, since Internet Service Providers have been slow to enable IP multicast in their networks, Internet multicast is only widely available as an overlay service. Commercial overlay networks, which provide multicast services have also begun to appear. This paper, introduces several routing algorithms that are suitable for routing in overly multicast networks and evaluates their performance. The algorithms seek to optimize the end-to-end delay and the *interface bandwidth usage* at the routing sites within the overlay network. The interface bandwidth is typically a key resource for an overlay network provider, and needs to be carefully managed in order to maximize the number of users that can be served. Through simulations, we evaluate the performance of these algorithms under various traffic conditions and on various network topologies.**

*Keywords*— **overlay network, load-balancing routing, application-level multicast**

## I. Introduction

Multicast communication is an important part of many next generation networked applications, including video conferencing, video-on-demand, distributed interactive simulation (including large multiplayer games) and peer-to-peer file sharing. Multicast services allow one host to send information to a large number of receivers, without being constrained by its network interface bandwidth. This makes applications more scalable and leads to more efficient use of network resources. The limited network layer support for multicast in the Internet today, has made it necessary for applications requiring multicast services to obtain services at a higher level. In *application layer multicast*, hosts participating in an application session share responsibility for forwarding information to other hosts [1–5]. While highly flexible, this approach places a significant additional burden on hosts, and is not as

efficient as network-layer multicast. *Overlay multicast networks* provide multicast services through a set of distributed *Multicast Service Nodes* (MSN), which communicate with hosts and with each other using standard unicast mechanisms. Overlay networks effectively use the Internet as a lower level infrastructure, to provide higher level services to end users. The multicast backbone, Mbone [6], is the best-known multicast overlay network, but multicast services are also a part of commercial overlay network services, such as Akamai [7] and iBeam [8].

Because overlay multicast networks are built on top of a general Internet unicast infrastructure, rather than point-to-point links, the problem of managing their resource usage is somewhat different than in networks that do have their own links. One of the principal resources that an overlay network must manage is the access bandwidth to the Internet at the MSNs' interfaces. This interface bandwidth represents a major cost, and is typically the resource that constrains the number of simultaneous multicast sessions that an overlay network can support. Hence, the routing algorithms used by an overlay multicast network, should seek to optimize its use.

In addition to optimizing MSN interface bandwidth, a multicast routing algorithm should ensure that the routes selected for multicast sessions do not contain excessively long paths, as such paths can lead to excessively long packet delays. However, the objective of limiting delay in a multicast network can conflict with the objective of optimizing the interface bandwidth usage, so multicast routing algorithms must strike an appropriate balance between these objectives. References [9] introduced the overlay multicast routing problem and studied the performance of two algorithms, using simulation. In this paper, we will briefly review these two algorithms and introduce a new algorithmic strategy that takes a more direct approach to optimizing the MSN interface bandwidth. We describe several specific algorithms based on this strategy and examine the performance of two of them in detail. The algorithms are evaluated using simulation and a range of traffic conditions and network configurations.

The rest of the paper is organized as follows: in section II, we briefly present the two multicast routing algo-

rithms that were developed earlier. Our new strategy for overly multicast routing is presented in section III. In section IV, we compare the performance of these routing algorithms on different network topologies and under various traffic conditions. In section V, we discuss issues related to dynamic membership control and implementation issues, and in section VI we discuss some of the related works. Finally we conclude in section VII.

## II. Background

An overlay multicast network can be modeled as a complete graph since there exists a unicast path between each pair of MSNs. For each multicast session, we create a shared overlay multicast tree spanning all MSNs serving participants of a session, with each tree edge corresponding to a unicast path in the underlying physical network. The amount of available interface bandwidth at an MSN imposes a constraint on the degree of that node in the multicast tree. We let $d_{max}(v)$ denote this degree constraint at node $v$.

There are two natural formulations of the overlay multicast routing problem. The first seeks to minimize diameter while respecting the degree constraints.

*Definition 1:* **Minimum diameter, degree-limited spanning tree problem (MDDL)**

Given an undirected complete graph $G = (V, E)$, a degree bound $d_{max}(v) \in N$ for each vertex $v \in V$ and a cost $c(e) \in Z^+$ for each edge $e \in E$; find a spanning tree $T$ of $G$ of minimum diameter, subject to the constraint that $d_T(v) \leq d_{max}(v)$ for all $v \in T$.

The MDDL problem is NP-hard. Reference [9] introduced a heuristic for MDDL, referred to here as the *Compact Tree* (CT) algorithm. It is a greedy algorithm and builds a spanning tree incrementally. We let $\delta(v)$ denote the length of the longest path from vertex $v$ to any other node in the partial tree $T$, constructed so far. For each vertex $v$ that is not yet in the partial tree $T$ constructed so far, we maintain an edge $\lambda(v) = \{u, v\}$ to a vertex $u$ in the tree; $u$ is chosen to minimize $\delta(v) = c(\lambda(v)) + \delta(u)$. At each step, we select a vertex $v$ with the smallest value of $\delta(v)$ and add it and the edge $\lambda(v)$ to the tree. Then, for each vertex $v$, not yet in the tree, we update $\lambda(v)$.

The second natural formulation of the overlay multicast routing problem seeks the "most balanced" tree, that satisfies an upper bound on the diameter. To explain what is meant by "most balanced", we first define the *residual degree* at node $v$ with respect to a tree $T$ as $res_T(v) = d_{max}(v) - d_T(v)$, where $d_T(v)$ is the degree of $v$ in $T$. To reduce the likelihood of blocking a future multicast session request, we should choose trees that maximize the smallest

residual degree. Since the sum of the degrees of all multicast trees is the same, this strategy works to "balance" the residual degrees of different vertices. Any tree that maximizes the smallest residual degree is a "most balanced" tree.

*Definition 2:* **Bounded diameter, residual-balanced spanning tree problem (LDRB)**

Given an undirected complete graph $G = (V, E)$, a degree bound $d_{max}(v)$ for each $v \in V$, a cost $c(e) \in Z^+$ for each $e \in E$ and a bound $B \in Z^+$; find a spanning tree $T$ of $G$ with diameter $\leq B$ that maximizes $\min_{v \in V} res_T(v)$, subject to the constraint that $d_T(v) \leq d_{max}(v)$, for all $v \in V$.

Like the MDDL problem, the LDRB problem is NP-hard. Reference [9] introduced a heuristic for LDRB, referred to here as the *Balanced Compact Tree* (BCT) algorithm. The algorithm can be viewed as a generalization of the CT algorithm. Like the CT algorithm, it builds the tree incrementally. However, at each step it first finds the $M$ vertices that have the smallest values of $\delta(v)$ and from this set, it selects a vertex $v$ with $\lambda(v) = \{u, v\}$, which maximizes the smaller of $res_T(u)$ and $res_T(v)$, where $T$ is the current partial tree. The parameter $M$ may be varied to trade-off the goals of residual degree balancing and diameter minimization. Specifically, when $M = 1$, it is equivalent to the CT algorithm and when $M$ is equal to the number of vertices in the multicast session, BCT concentrates on balancing the residual degrees. Simulation studies have shown that fairly small values of $M$ are effective in achieving good balance, without violating the diameter bound.

We evaluate overlay multicast routing algorithms using a simulation in which new multicast sessions start and end at random times. The primary performance metric is the fraction of sessions that are rejected because no multicast route can be found, either due to the failure to satisfy the diameter bound or due to the exhaustion of interface bandwidth of at least one MSN. We also obtain a lower bound on the rejection probability using a simulation in which a multicast session is rejected only if the MSN interface bandwidth required by the session exceeds the total unused interface bandwidth at all MSNs (including those not involved in the session). Results reported in [9] showed that by distributing the load more evenly across servers, the BCT algorithm rejects substantially fewer multicast sessions than the CT algorithm on the same network configuration. At the mean time, there remained a significant gap between the achieved performance and the potential suggested by the lower bound. While the bound was not expected to be tight, the size of the gap suggested that there was room for improvement. In the next section, we intro-

duce a new strategy for overlay multicast routing that leads to an algorithm with uniformly better performance.

## III. BALANCED DEGREE ALLOCATION

The incremental nature of the BCT algorithm prevents it from achieving the best-possible residual degree balance. *Balanced Degree Allocation* (BDA) is a strategy for constructing multicast trees that approaches the problem in a fundamentally different way. It starts by determining the ideal degree of each node in the multicast session, with respect to the objective of maximizing the smallest residual degree.

To state the strategy precisely, we need to stretch our definition of the residual degree of a vertex. Let $k$ denote the multicast session fanout. First, we define a *degree allocation* $d_A$ to be a function from the vertices of a multicast session to the positive integers that satisfies two properties: (1) $\sum_v d_A(v) = 2(k-1)$ and (2) there are at least two vertices $u$ and $v$ with $d_A(u) = d_A(v) = 1$. A partial degree allocation is a similar function in which the first property is replaced with $\sum_v d_A(v) \leq 2(k-1)$. Now, define the residual capacity of a vertex with respect to a partial degree allocation $d_A$ as $res_A(v) = d_{max}(v) - d_A(v)$.

We can compute a degree allocation that maximizes the smallest residual degree as follows.
• For each vertex $v$ in the multicast session, initialize $d_A(v)$ to 1.
• While $\sum_v d_A(v) < 2(k-1)$ select a vertex $v$ that maximizes $res_A(v)$ and increment $d_A(v)$.
This procedure actually does more than maximize the smallest residual degree. It produces the most balanced possible set of residual degrees, as illustrated in Figure 1. It shows how the degree allocation increases with the residual degree of a node.
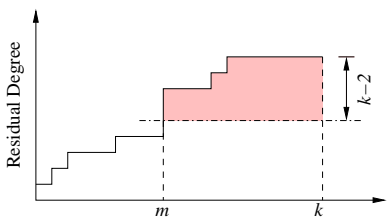


Fig. 1. Balanced Degree Allocation

Given a degree allocation for a tree, we would like to construct a tree in which the vertices have the assigned degrees and which satisfies the limit on the diameter. There is a general procedure for generating a tree with a given degree allocation, which is described below.

The procedure builds a tree by selecting *eligible pairs* of vertices, and adding the edge joining them to a set of edges $F$ that, when complete, will define the tree. The spare degree of a vertex $v$, is $d_A(v)$ minus the number of edges in $F$ that are incident to vertex $v$. At any point in the algorithm, the edges in $F$ define a set of connected components. A pair of vertices $\{u, v\}$ is an eligible pair if the following conditions are satisfied.
• $u$, $v$ are not in the same component;
• both $u$, $v$ have spare degree $\geq 1$;
• either, there are only two components remaining, or the summation of the spare degree of the vertices in the components containing $u$ and $v$ is greater than 2.

This process is guaranteed to produce a tree with the given degree allocation, and all trees with the given degree allocation are possible outcomes of the process. We can get different specific tree construction algorithms by providing different rules for selecting the vertices $u$ and $v$.

One simple and natural rule is to select the closest pair $u$ and $v$. Call this the *Closest Pair* (CP) algorithm. Since the CP algorithm does nothing to directly address the objective of diameter minimization, it may not produce a tree that meets the diameter bound. An alternative selection rule is to select the pair $\{u, v\}$ that results in the smallest diameter component in the collection of components constructed as the algorithm progresses. This algorithm is referred to as the *Compact Component* (CC) algorithm.

One can also use a selection rule that builds a single tree incrementally. In this rule, we consider all eligible pairs of vertices $u$ and $v$, that either $u$ is in the tree built so far, or $v$ is in the tree, but not both. Among all such pairs, we pick the one that results in the smallest diameter tree. This procedure is repeated for every choice of initial vertex, and the smallest diameter tree kept. This algorithm is the same as the CT algorithm described before, but with the original input degree constraint $d_{max}$ replaced with the much tighter degree allocation $d_A$.

Any of the selection rules described, will produce a tree with the desired degree allocation. However, the resulting tree may not satisfy the bound on diameter. We can reduce the diameter of the trees generated by using a less balanced degree allocation. To reduce the diameter, we can increase the degree allocation of "central vertices" while decreasing the degree allocation of "peripheral vertices." Unfortunately, we have found that natural strategies for adjusting degree allocations lead to only marginal improvement in the diameters of the resulting trees. It appears difficult to find good degree allocations, independently of the tree building process. We have found that a more productive approach is to use a "loose degree allocation" and allow the tree-building process to construct a suitable tree satisfying the degree limits imposed by the loose allocation. Loose degree allocations are derived from the most balanced allocation by allowing small increases in the degrees
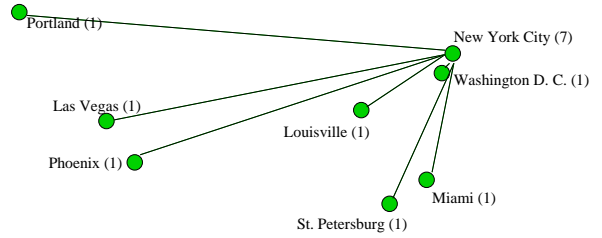
of vertices.

By combining the tree building procedure, with a specific rule for selecting eligible pairs and a procedure for "loosening" a degree allocation, we can define iterative algorithms for the overlay multicast routing problem. We start with a balanced degree allocation and build a tree using that allocation. If the resulting tree satisfies the diameter bound, we stop. Otherwise we loosen the degree allocation and build a new tree. If this tree satisfies the diameter bound, we stop. Otherwise we continue to loosen the bound, stopping when we find a tree with small enough diameter, or until a decision is made to terminate the process and give up.
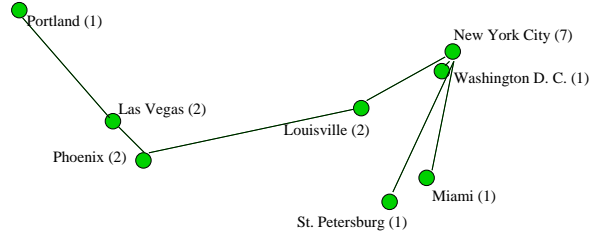
The degree loosening procedure increases by 1, the degree allocation of up to $b$ vertices, where $b$ is a parameter. The first application of the procedure adds 1 to the degree allocation of the $b$ most central vertices. If incrementing the degree allocation for a vertex $v$ would cause its degree allocation to exceed the degree bound for $v$, then $v$'s degree allocation is left unchanged. A vertex $u$ is more central than a vertex $v$ if its radius $max_w c(\{u, w\}) < max_w c(\{v, w\})$. The second round of the procedure adds 1 to the next $b$ most central vertices (again, so long as this would not cause their degree bounds to be exceeded). Subsequent applications affect the degree bounds of successive groups of $b$ vertices, and the process wraps around to the most central vertices, after all vertices have been considered. The process can be stopped after some specified number of applications of the degree loosening procedure, or after all degree allocations have been increased to the smaller of their degree bounds and $k - 1$ where $k$ is the number of vertices.

If we select eligible pairs for which the connecting edge has minimum cost, the resulting algorithm is called *Iterative Closest Pair* (ICP) algorithm. If we select eligible pairs so as to minimize the diameter of the resulting component, the algorithm is called *Iterative Compact Component* (ICC) algorithm. In the *Iterative Compact Tree* (ICT) algorithm, we select eligible pairs, with one vertex of each pair in a single tree being constructed, and the other selected to minimize the tree diameter. This procedure is repeated for all possible initial vertices.

An example execution of the ICT algorithm appears in Figure 2. For simplicity, we used geographical distance as routing cost and a diameter bound of 8000 km. Initially, the BDA output dictates the creation of a star topology with New York City as the center with degree of 7; this exceeds the diameter bound. In the second round, the degree allocation of the three nodes with small radius, Las Vegas, Phoenix, Louisville, is loosened by 1; this results in a smaller diameter tree satisfying the diameter bound.



(a) Longest Path: Portland → New York → Phoenix; Length = 9416.67 km



(b) Longest Path: Portland → Las Vegas → Phoenix → Louisville → New York → St. Petersburg; Length = 7823.64 km

Fig. 2. An Example of the ICT Algorithm with Degree Adjustment

We observe that the actual degree allocation is still close to the balanced degree allocation.

## IV. EVALUATION

This section reports simulation results for the overlay multicast routing algorithms described above. We report results for three network topologies and a range of multicast session sizes. The principal performance metric is the multicast session rejection rate. We also evaluate the multicast tree diameter and computation times of the algorithms. Comparisons with the CT and BCT algorithms are also included.

### A. Simulation Setup

We have selected three overlay network configurations for evaluation purposes. The first (called the metro configuration) has an MSN at each of the 50 largest metropolitan areas in the United States [10]. The "traffic density" at each node is proportional to the population of the metropolitan area it serves. We use a Poisson session arrival process and the session holding times follow a Pareto distribution. Session fanouts follow a truncated binomial distribution with a minimum of 2 and maximum of 50, and means varied in different result sets. All multicast sessions are assumed to have the same bandwidth. Different MSNs were assigned different interface bandwidths, depending

on their traffic density and their location. MSNs in more central locations are assigned higher interface bandwidths than those in less central locations, since it is more efficient for multicast sessions to branch out from these locations than from the more peripheral locations. The assignment of interface bandwidth at MSNs is critical to the performance of the routing algorithms. We have dimensioned the network to best carry a projected traffic load given a specific routing algorithm. Details of the MSN dimensioning process can be found in [9].
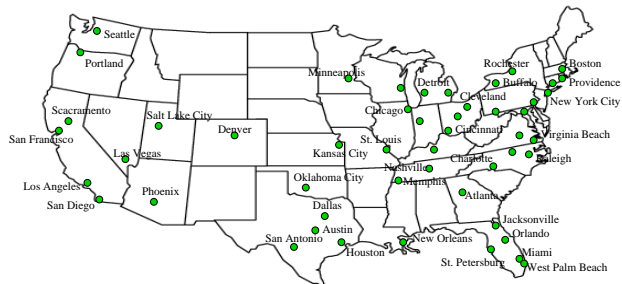
The metro configuration was chosen to be representative of a realistic overlay multicast network. However, like any realistic networks, it is somewhat idiosyncratic, since it reflects the locations of population centers and the differing amounts of traffic they produce. The other two configurations were chosen to be more neutral. The first of these consists of 100 randomly distributed nodes on a disk and the second consists of 100 randomly distributed nodes on the surface of a sphere. In both cases, all nodes are assumed to have equal traffic densities. In the disk, as in the metro configuration, the MSN interface bandwidths must be dimensioned, but in this case it is just a node's location that determines its interface bandwidth. In the sphere configuration, all nodes are assigned the same interface bandwidth, since there is no node that is more central than any other.

The three network configurations are illustrated in Figure 3. In all configurations, the geographical distance between two nodes is taken as the cost of including an edge between those nodes in the multicast session tree.
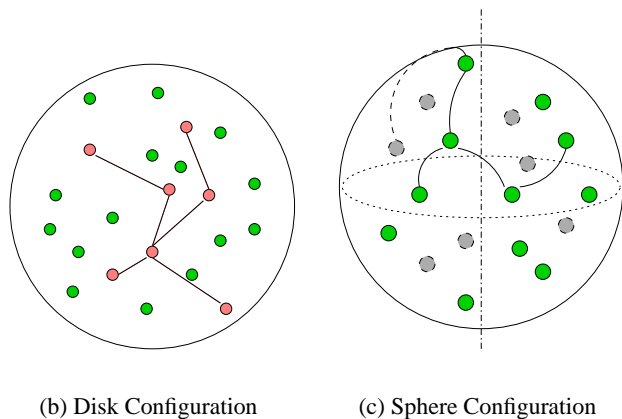
### B. Comparison of Tree Building Techniques

In the previous section, we suggested three basic tree building techniques: selecting the closest pair (CP), selecting the pair that minimizes the component diameter (CC), and selecting the pair that minimizes the single tree diameter (CT). The iterative versions of these algorithms, namely ICP, ICC and ICT, seek to satisfy the diameter bound by loosening the degree allocation produced by BDA. In this section, we examine their performance sensitivities to different diameter bounds and to the number of rounds allowed for degree adjustment. The simulation uses the metro configuration as the network topology and a session fanout of 10.

Figure 4 shows the session rejection rates versus the ratio of the diameter bound to the maximum inter-city delay (6000 km). In this simulation, we allow each algorithm to loosen the degree allocation as many rounds as possible, until it reaches the smaller of nodes' degree bounds or $k - 1$, where $k$ is the session size. The horizontal line labeled as BDA, shows the rejection rate using the balanced



(a) 50 Largest U.S. Metropolitan Areas



(b) Disk Configuration         (c) Sphere Configuration

Fig. 3. Overlay Network Configurations

degree allocation strategy, but ignoring the diameter of the resulting tree.

As the large cities in this map are along the coastal areas, the majority of the sessions will span across the continent. Therefore, it is difficult to find a multicast tree for these sessions when the diameter bound is tight, resulting in very high rejection rates for all algorithms. However, as the diameter bound is relaxed, the rejection rate improves for all the algorithms, with the iterative algorithms all achieving essentially the same performance for diameter bounds of more than 1.8 times the maximum inter-city distance. At intermediate diameter bounds, the ICT and BCT algorithm performs better than the ICC and ICP algorithms. This suggests that building from a single tree, as both ICT and BCT do, is more effective in minimizing the tree diameter. The BCT algorithm does not allocate its node degree before building the tree; rather, it seeks to maximize the residual bandwidth as the algorithm progress. For large diameter bounds, BCT is not able to reduce its rejection rate as much as the algorithms using balanced degree allocation.

Figure 5 shows the rejection ratio versus the maximum number of degree adjustment rounds allowed in ICP, ICC and ICT. The diameter bound is fixed at 8000 km for this
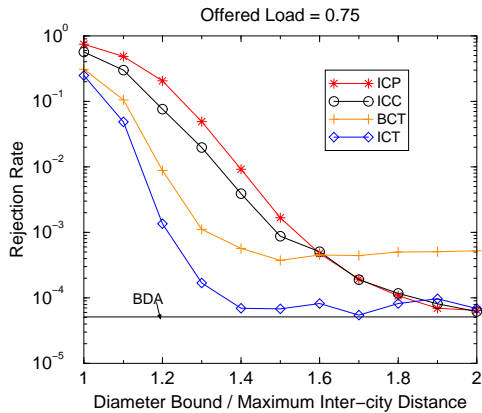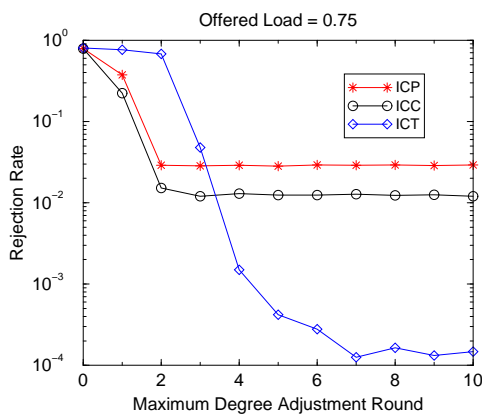
Fig. 4. Sensitivity to Diameter Bound



Fig. 5. Sensitivity to Degree Adjustment Round

simulation. In each degree adjustment round, the number of vertices being adjusted is 3. Generally, the ICP and ICC algorithms only benefit from the very first few rounds of degree adjustment, which allows nearby nodes to be joined together to form collections of small forests. The additional rounds of degree adjustment have no effect on them and the rejection rates remain relatively high. Contrarily, the ICT algorithm benefits greatly from the additional rounds of degree adjustment. It is able to utilize the increased degree allocation at the centrally located nodes and form smaller diameter trees.

We conclude in this subsection that the ICT algorithm, when combined with the degree loosening procedure, is more effective at producing small diameter trees than ICP and ICC. In the rest of this paper, we will focus our evaluation mainly on the ICT algorithm. However, ICT's greater effectiveness comes with a cost of added complexity, as it iterates through each possible starting vertex in order to find the best tree. We will analyze the computational cost of the ICT algorithm later in the section.

## C. Performance Results – Rejection Rate

Figure 6 shows the session rejection rates versus offered load for a subset of the overlay multicast routing algorithms presented earlier. The charts also include a lower bound on the rejection fraction that was obtained by running a simulation in which a session is rejected only if the sum of the degree bounds at all nodes in the network is less than $2(k-1)$ where $k$ is the number of nodes in the session being set up. The lower bound curves are labeled LB. The results, labeled BDA, are obtained using the balanced degree allocation strategy, and ignoring the diameter of the resulting tree. We conjecture that this also represents a lower bound on the best possible rejection fraction that can be obtained by any on-line routing algorithm. It is certainly a lower bound for algorithms based on the balanced degree allocation strategy.

For these charts, the diameter bound for the metro configuration is 8000 km which is approximately 1.5 times the maximum distance between nodes. For the disk and sphere topology, the bound is two times the disk diameter and three times the sphere diameter; each is about twice the maximum distance between any two nodes. The total interface bandwidth for all MSNs is 10,000 times the bandwidth consumed by a single edge of a multicast session tree. So, in the sphere, each MSN can support a total of 100 multicast session edges and for the metro configuration, the average number is 200.

Overall, the results show that the algorithms that seek to balance the residual degree usually perform much better than the CT algorithm, which merely seeks to minimize the diameter subject to a constraint on the maximum degree bound of a node $(d_{max}(v))$. The performance of CT is particularly poor in the fanout 20 case, since it tends to create nodes with large fanout leading to highly unbalanced residual degree distributions.

Looking down each column, we see that the lower bound increases with the fanout. This simply reflects the fact that each session consumes a larger fraction of the total interface bandwidth. For the sphere, the rejection fraction also increases with fanout for the BDA curve. This makes sense intuitively, since as the fanout grows, one expects it to be more difficult to find balanced trees with small enough diameter. In the disk and metro configurations, it is less clear why the BDA curve changes with fanout as it does. Part of the explanation for the observed behavior is that the network dimensioning process is based on an assumed traffic load, and in particular, an assumed multicast session fanout of 10. When the simulated traffic has the same fanout distribution as the one used to dimension the network, we get smaller rejection rates. However,
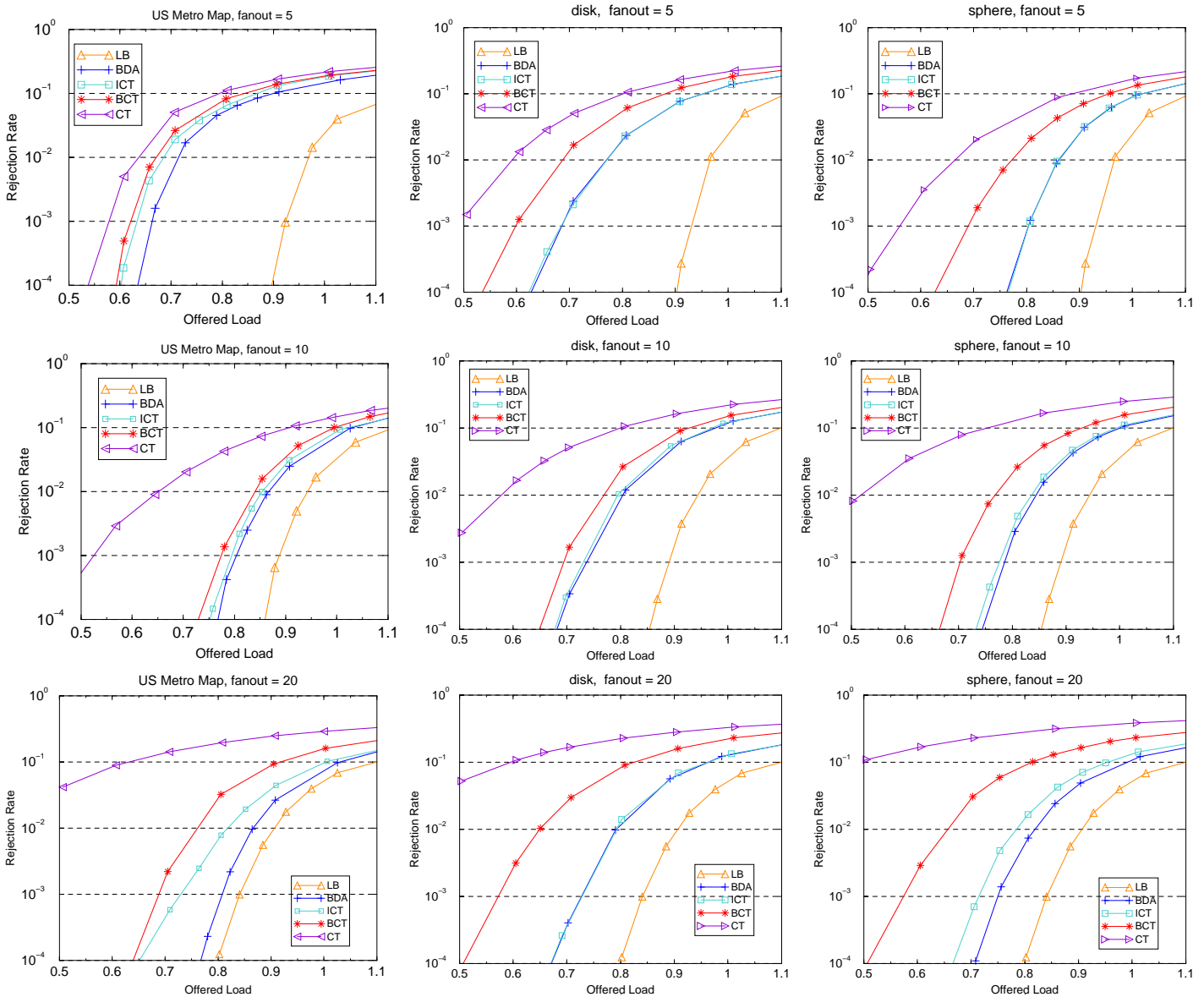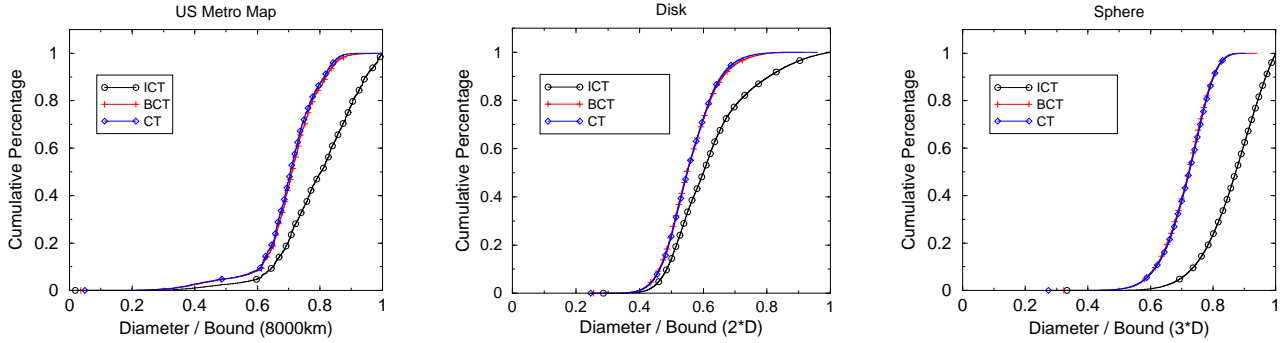
Fig. 6. Rejection Fraction Comparison

there is a somewhat surprising deterioration of the rejection rate for small fanout, particularly in the metro network case. The apparent explanation for this is that with small fanouts, we often get sessions involving MSNs near the east and west coasts, but none in the center of the country. Such sessions are unable to exploit the ample unused bandwidth designed into the more central MSNs, based on a larger average fanout. A similar effect is observed with the disk, but it is more extreme in the metro configuration because of the greater population densities on the coasts, and also the smaller ratio of the diameter bound to the maximum inter-MSN distance (1.5 vs. 2).

The curves for the ICT algorithm are generally quite close to the BDA curves, leaving little apparent room for improvement. There are small but noticeable gaps in a few cases. For the metro configuration, the gaps for fanout 5 and 20 are most probably due to the difference between the average fanout of the simulated traffic and the fanout used for dimensioning. This explanation does not account for the gap in the sphere configuration for fanout 20, since in the sphere, all nodes have the same interface bandwidth. The most plausible explanation seems to be that with large fanout, it just becomes intrinsically more difficult to find trees that satisfy the diameter bound.

In three out of the nine cases shown, the BCT algorithm performs nearly as well as the ICT algorithm. Its performance relative to ICT is worst for the sphere and best for the metro configuration.

| | metro | | disk | | sphere | |
|------|---------|-------|---------|-------|---------|-------|
| | mean | std. | mean | std. | mean | std. |
| CT | 5581.22 | 14.4% | 222.38 | 13.9% | 429.16 | 10.0% |
| BCT | 5627.77 | 14.5% | 222.67 | 14.5% | 428.17 | 10.1% |
| ICT | 6347.36 | 15.7% | 248.17 | 19.9% | 516.18 | 10.4% |

Fig. 7. End-to-end Delay performance

## D. Performance Results – Tree Diameter

Next, we investigate the performance of these algorithms in terms of the tree diameter they created. Figure 7 shows the cumulative distribution of the tree diameter scaled to the diameter bound used in the algorithms. The fanout used here is 10 per session. Also, we show in a table the mean and variance of the tree diameter using unscaled value.

We observe that the diameter performance of the BCT algorithm is as good as that of the CT algorithm; the difference is almost indiscernible. The ICT algorithm does, however, show a performance degradation on the tree diameter. Especially in the sphere, ICT generates trees with diameter much larger than those generated by the other two algorithms. This is due to that traffic load tends to be more evenly distributed in the sphere and BDA is very effective at equalizing the available capacity at each MSN, that it tends to produce degree allocations with large numbers of vertices of degree 2, resulting in long paths.

The different performance perspective of ICT vs. BCT suggests that it is possible to use BCT for application sessions that require better end-to-end delay performance, and use ICT algorithm for others. Although we have not investigated the system utilization when both algorithms are deployed, we conjecture that it should probably do better than the BCT algorithm.

## E. Complexity of the ICT Algorithm

Although the ICT algorithms gives superior performance on the overall system utilization and can satisfy the diameter bound in most cases, it potentially has the disadvantage of higher computational complexity, which in reality may be hard to implement. Especially with a more stringent diameter bound, as the algorithm continues to loosen the nodes' degree allocation in a round by round basis in order to search for an appropriate tree, the complexity goes higher as well as the rejection ratio.
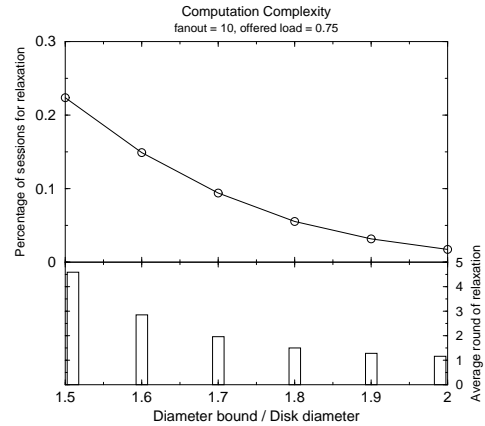


Fig. 8. Computation Complexity of the ICT Algorithm

In Figure 8, we plotted the computation requirement for the ICT algorithm on a disk topology with session fanout equal to 10. The top half of the figure shows the percentage of sessions that require degree adjustment; and the bottom half shows the average number of rounds iterated for these sessions. We observe that the actual number of additional rounds is reasonably low. For instance, if the diameter bound is 1.7 times the disk diameter, 10% of the sessions require an average of 2 rounds of degree adjustment. However, if the diameter bound is too stringent, the extra rounds of degree adjustment do not help much in reducing the session rejection rate (not shown). Therefore, it

is important to pick a suitable diameter bound for a topology so that the algorithm can operate more efficiently and more effectively. From our experience, a bound of twice the maximum distance between any pair of nodes seem to be a good choice in all three network configurations.

## V. DISCUSSIONS

In our current simulation model, we have not considered the dynamic session memberships. The main difference when incorporating dynamic member join and leave requests, is that when a node no longer has directly attached end users but only serves as a transit node, there can be different choices as whether to remove the server and if remove, when to remove it. The removal operation should be such that it does not cause any service disruptions. This could be difficult if the node has a large fanout. Conceivably, the node can only be removed when all its upstream and downstream nodes are connected in an optimal way. Or, we can decide to remove a node only if it has a fanout of two. In any case, the time it takes to stabilize the new connections is crucial to the performance, yet this time duration also depends on a variety of parameters and is hard to simulate in a meaningful way.

The issue of implementing the algorithms in a distributed fashion also needs to be addressed. If implemented in a fully distributed fashion, the proposed algorithms require synchronized update at the end of each node addition to the tree, which is potentially inefficient and unscalable. Alternatively, if it is implemented in a centralized way, i.e. let one of the MSNs compute the routes and inform others to connect as a tree, we can eliminate many message exchanges and other hazards for state synchronization necessary for a distributed computation. We should point out that the centralized version does not create a single point of failure, as each session is likely to select a different delegate for tree computation. More over, as overlay network servers potentially have imprecise information from end-to-end measurements, a distributed implementation only increases the risk of computing an erroneous tree, such as a loop or partition in the tree. Therefore, we think that a delegation-based scheme is more appropriate for overlay network routing.

## VI. RELATED WORK

The Multicast Backbone (Mbone) [6], is the best known and widely adopted multicast overlay network. The Mbone is implemented as tunnels at the network layer and implements distance vector multicast routing protocol [11]. Other standard routing protocols include: core-based tree (CBT) [12], protocol independent multicast (PIM) [13] and most recently, source specific multicast (SSM) [14]. All of these routing protocols build shortest path trees from data sources or from the core node of a session, to minimize network cost.

There are a number of application-level multicast services appearing in recent literatures, mostly due to the dwindling usage on the Mbone and the slow deployment of network multicast services. The flexibility of the application-level multicast services allow the routing policy to be changed based on the target application requirement. For example, Scattercast [15] uses delay as the routing cost and builds shortest path tree from data sources; Overcast [4] explicitly measures available bandwidth on an end-to-end path and builds a multicast tree by maximizing the available bandwidth from the source to the receivers; and Endsystem multicast [1] uses a combination of delay and available bandwidth, and prioritizes available bandwidth over delay when selecting a routing path.

In this paper, we have defined *interface bandwidth* as our primary routing metric. The path selection policies seek to optimize the usage of interface bandwidth of MSNs while satisfying the end-to-end delay performance of individual session. As a result, our routing algorithms differ fundamentally from all of the aboves.

## VII. CONCLUSIONS

In this paper, we presented several multicast routing algorithms that are specifically designed for overlay networks, where the routing metric and cost model differ drastically from conventional networks. We address these differences mainly as two parameters: load balancing and end-to-end delay bound. Our evaluation showed that it is possible to achieve a large gain on the system utilization without sacrificing much on the end-to-end delay performance and the algorithms are robust under traffic variations and topology variations.

As the overlay network concept is applicable in many areas, there are more varieties in the actual routing algorithms, as they have to consider many different parameters, such as fault tolerance and route redundancy. We hope our work will provide some ground work for future algorithm design and evaluations.

## REFERENCES

[1] Y. hua Chu, S. G. Rao, S. Seshan, and H. Zhang. Enabling Conferencing Applications on the Internet Using an Overlay Multicast Architecture. In *Proc. ACM SIGCOMM 2001*, San Diago, CA, August 2001.

[2] P. Francis. Yallcast: Extending the Internet Multicast Architecture. http://www.yallcast.com, September 1999.

[3] D. Pendarakis, S. Shi, D. Verma, and M. Waldvogel. ALMI: An **A**pplication **L**evel **M**ulticast **I**nfrastructure. In *3rd USENIX Symposium on Internet Technologies and Systems (USITS'01)*, San Francisco, CA, March 2001.

[4] J. Jannotti, D. K. Gifford, K. L. Johnson, M. F. Kaashoek, and J. W. O. Jr. Overcast: Reliable Multicasting with an Overlay Network. In *Proc. OSDI'01*, 2000.

[5] S. Zhuang, B. Zhao, A. D. Joseph, R. H. Katz, and J. Kubiatowicz. Bayeux: An architecture for wide-area, fault-tolerant data dissemination. In *Proc. NOSSDAV'01*, June 2001.

[6] H. Erikson. MBONE: The Multicast Backbone. *Communication of the ACM*, pages 54–60, August 1994.

[7] Akamai Technologies, Inc. http://www.akamai.com.

[8] iBeam Broadcasting Corp. http://www.ibeam.com.

[9] S. Shi, J. Turner, and M. Waldvogel. Dimension Server Access Bandwidth and Multicast Routing in Overlay Networks. In *11th International Workshop on Network and Operating System Support for Digital Audio and Video (NOSSDAV'01)*, June 2001.

[10] U.S. Census Bureau. http://www.census.gov/population/www/estimates/metropop.html.

[11] F. Baker. Distance Vector Multicast Routing Protocol - DVMRP. RFC 1812, June 1995.

[12] A. Ballardie. Core Based Trees (CBT version 2) Multicast Routing - Protocol Specification. RFC 2189, September 1997.

[13] D. Estrin, V. Jacobson, D. Farinacci, L. Wei, S. Deering, M. Handley, D. Thaler, C. Liu, S. P., and A. Helmy. Protocol Independent Multicast-Sparse Mode (PIM-SM): Motivation and Architecture. Internet Engineering Task Force, August 1998.

[14] H. Holbrook and B. Cain. Source Specific Multicast. IETF draft, draft-holbrook-ssm-00.txt, March 2000.

[15] Y. Chawathe, S. McCanne, and E. Brewer. An Architecture for Internet Content Distribution as an Infrastructure Service. Unpublished, available at http://www.cs.berkeley.edu/ yatin/papers.