

An Overview of Quality of Service Routing for Next-Generation High-Speed Networks: Problems and Solutions

Shigang Chen and Klara Nahrstedt, University of Illinois at Urbana-Champaign

Abstract

The upcoming gigabit-per-second high-speed networks are expected to support a wide range of communication-intensive real-time multimedia applications. The requirement for timely delivery of digitized audio-visual information raises new challenges for next-generation integrated services broadband networks. One of the key issues is QoS routing. It selects network routes with sufficient resources for the requested QoS parameters. The goal of routing solutions is twofold: (1) satisfying the QoS requirements for every admitted connection, and (2) achieving global efficiency in resource utilization. Many unicast/multicast QoS routing algorithms have been published recently, and they work with a variety of QoS requirements and resource constraints. Overall, they can be partitioned into three broad classes: (1) source routing, (2) distributed routing, and (3) hierarchical routing algorithms. In this article we give an overview of the QoS routing problem as well as the existing solutions. We present the strengths and weaknesses of different routing strategies, and outline the challenges. We also discuss the basic algorithms in each class, classify and compare them, and point out possible future directions in the QoS routing area.

The timely delivery of digitized audio-visual information over local or wide area networks is now becoming realistic, thanks to fruitful research in high-speed networks, image processing, and video/audio compression. On the other hand, the emerging distributed multimedia applications also raise new challenges for network research and development. For example, a video-on-demand application requires that its data throughput over the network be guaranteed at or above a certain rate.

In the current Internet, data packets of a session may follow different paths to the destination. The network resources (e.g., switch buffer and link bandwidth) are fairly shared by packets from different sessions. However, this architecture does not meet the requirements of future integrated services networks that will carry heterogeneous data traffic. First, it does not support resource reservation, which is vital for the

provision of guaranteed end-to-end performance (bounded delay, bounded delay jitter, and/or bounded loss ratio). Second, data packets may experience unpredictable delays and arrive at the destination out of order, which is undesirable for continuous real-time media. Hence, the next generation of high-speed wide area networks is likely to be *connection-oriented* for real-time traffic.¹ This article focuses on the routing problem of connection establishment. For unicast, the problem is to find a network *path* that meets the requirement of a connection between two end users. For multicast, the problem is to find a multicast *tree*, rooted at a sender, which covers all receivers with every internal path from the sender to a receiver satisfying the requirement.

The notion of quality of service (QoS) has been proposed

This work is supported by the ARPA grant under contract number F30602-97-2-0121 and the National Science Foundation Career grant under contract number NSF CCR 96-23867.

¹ At the transport layer, a connection (call) means (1) the logical association between end users, and (2) the correct, ordered delivery of data. At the network layer a connection means a network path, consisting of switches and links, which connects the end users. Data packets (or cells) of the same connection are sent along the path in first in first out (FIFO) order.

to capture the qualitatively or quantitatively defined performance contract between the service provider and the user applications. The QoS requirement of a connection is given as a set of constraints, which can be *link constraints*, *path constraints* [1], or *tree constraints*. A *link constraint* specifies a restriction on the use of links. A *bandwidth constraint* of a unicast connection requires, for instance, that the links composing the path must have certain amount of free bandwidth available. A *path constraint* specifies the end-to-end QoS requirement on a single path; a *tree constraint* specifies the QoS requirement for the entire multicast tree. The *delay constraint* of a multicast connection requires that the longest end-to-end delay from the sender to any receiver in the tree must not exceed an upper bound.

A *feasible path (tree)* is one that has sufficient residual (unused) resources to satisfy the QoS constraints of a connection. The basic function of QoS routing is to find such a feasible path (tree). In addition, most QoS routing algorithms consider the optimization of resource utilization measured by an abstract metric, *cost*. The cost of a link can be defined in dollars or as a function of the buffer or bandwidth utilization. The cost of a path (tree) is the total cost of all links on the path (tree). The optimization problem is to find the lowest-cost path (tree) among all feasible paths (trees).

The problem of QoS routing is difficult for a number of reasons. First, distributed applications such as Internet phone and distributed games have very diverse QoS constraints on delay, delay jitter, loss ratio, bandwidth, and so on. Multiple constraints often make the routing problem intractable. For example, finding a feasible path with two independent path constraints is NP-complete [2]. Second, any future integrated services network is likely to carry both QoS and best-effort traffic, which makes the issue of performance optimization complicated. It is hard to determine the best operating point for both types of traffic if their distributions are independent. Although the QoS traffic will not be affected due to resource reservation, the throughput of the best-effort traffic will suffer if the overall traffic distribution is misjudged. Third, the network state changes dynamically due to transient load fluctuation, connections in and out, and links up and down. The growing network size makes it increasingly difficult to gather up-to-date state information in a dynamic environment, particularly when wireless communication is involved. The performance of a QoS routing algorithm can be seriously degraded if the state information used is outdated.

Many QoS routing algorithms have been proposed recently with a variety of constraints considered. The purpose of this article is to provide a survey of recent development in this area. In the following, we present different routing problems, their challenges, the routing strategies, a classification and comparison of existing routing algorithms, and possible future directions. We refer to QoS routing simply as "routing" unless it is necessary to make clear distinction between QoS and best-effort routing.

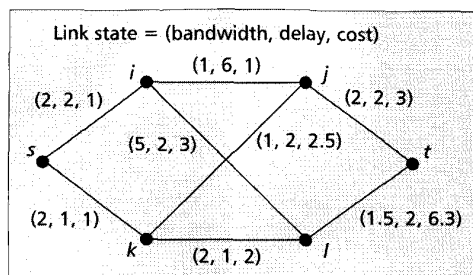


Figure 1. Network state.

The Weighted Graph Model

A network can be modeled as a graph (V, E) . Nodes (V) of the graph represent switches, routers, and hosts. Edges (E) represent communication links. The edges are undirected only if the communication links are always symmetric.

A symmetric link has the same properties (capacity, propagation delay, etc) and the same traffic volume in both directions. For most real networks the communication links are asymmetric, and hence every link is represented by two directed edges in the opposite directions. It should be noted that although the examples in this article use the undirected graphs for fewer edges, most routing algorithms under discussion were designed for asymmetric networks.

Every link has a state measured by the QoS metrics of concern. In Fig. 1 the link state is a triple consisting of residual bandwidth, delay and cost. Every node also has a state. The node state can be either measured independently or, as in this article, combined into the state of the adjacent links. In the latter case the residual bandwidth is the minimal of the link

bandwidth and the CPU bandwidth;² the delay of a link consists of the link propagation delay and the queuing delay at the node; the cost of a link is determined by the total resource consumption at the link and node.

Destination	i	j	k	l	t
Bandwidth	2	1.5	2	2	1.5
Next hop	i	k	k	i	i
Destination	i	j	k	l	t
Delay	2	3	1	2	4
Next hop	j	k	k	k	k
Destination	i	j	k	l	t
Cost	1	2	1	3	5
Next hop	i	i	k	k	i

Table 1. Global state in distance vectors at node s.

Maintenance of State Information

Routing consists of two basic tasks. The first task is to collect the state information and keep it up to date. The second task is to find a feasible path for a new connection based on the collected information. The performance of any routing algorithm directly depends on how well the first task is solved.

Local State — Each node is assumed to maintain its up-to-date local state, including the queuing and propagation delay, the residual bandwidth of the outgoing links, and the availability of other resources.

Global State — The combination of the local states of all nodes is called a *global state*. Every node is able to maintain the global state by either a link-state protocol [3–5] or a distance-vector protocol [6–8], which exchanges the local states among the nodes periodically. Link-state protocols broadcast the local state of every node to every other node so that each node knows the topology of the network and the state of every link (Fig. 1). Distance-vector protocols periodically exchange distance vectors among adjacent nodes. A distance vector has an entry for every possible destination, consisting of the property of the best path and the next node on the best path (Table 1).

The global state kept by a node is always an approximation of the current network state due to the nonnegligible delay of

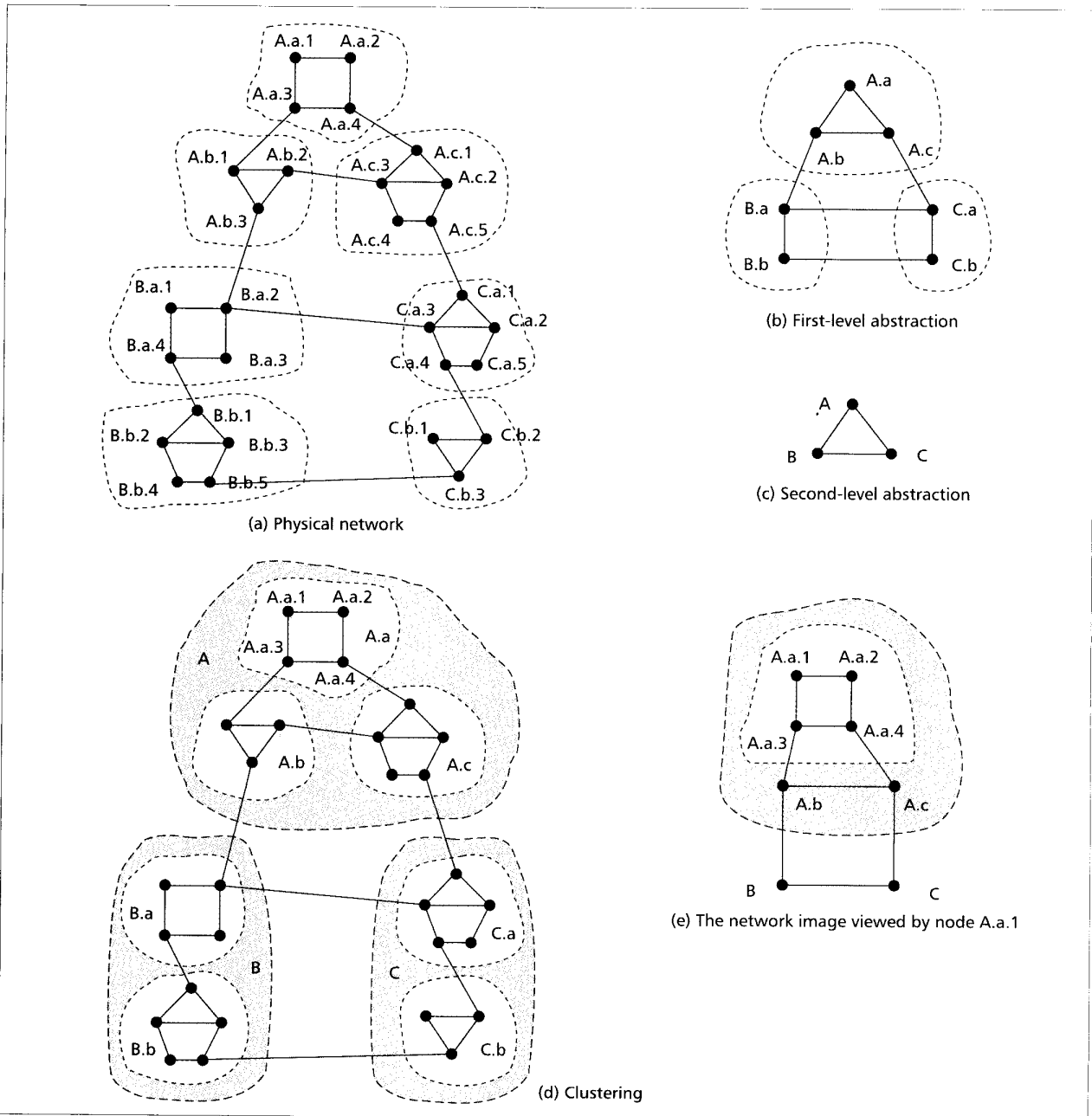
² The CPU bandwidth is defined as the maximum rate at which the node can pump data into the link.

propagating local states. As the network size grows, the imprecision increases.

Aggregated (Partial) Global State — A common approach to achieving scalability is to reduce the size of the global state by aggregating information according to the hierarchical structure of large networks. Figure 2 shows the hierarchical model used by [9, 10]. In Fig. 2a, nodes are clustered into the first-level groups. The nodes with at least one link crossing two groups are called *border nodes*. In Fig. 2b, each group is represented by a *logical node*. A physical node in the group is elected to act on behalf of the logical node and store the higher-level state information.

The links connecting logical nodes are *logical links*. The logical nodes are further clustered to form higher-level groups, which are abstracted by higher-level logical nodes (Fig. 2c). Figure 2d presents the overall clustering. On each hierarchy level, the nodes in a group are called the *children* of the logical node representing the group; the logical node is called the *parent*. An ancestor of a node is either its parent or an ancestor's parent. We have used the simplest topology aggregation, which abstracts a group by a single logical node. There are other types of aggregation using different simple topologies to replace a group. Their performance was studied in [11].

Each physical node maintains an *aggregated network image*.



■ Figure 2. The hierarchical network model.

The image maintained at node *A.a.1* is shown in Fig. 2e. It stores different portions of the network in different details. More specifically, the image is derived by starting from the highest hierarchy level and recursively replacing the ancestor of the node with the corresponding lower-level group. As the network topology is aggregated, the state information is aggregated as well. The state of each logical link is the combination of the states of many lower-level links. The link-state algorithm can be extended to collect the aggregated state information for every node [9]. As the state is aggregated, the imprecision is also aggregated.

Routing Problems

The routing problems can be divided into two major classes: *unicast routing* and *multicast routing*. The unicast routing problem is defined as follows: given a source node *s*, a destination node *t*, a set of QoS constraints *C*, and possibly an optimization goal, find the best feasible path from *s* to *t* which satisfies *C*. The multicast routing problem is defined as follows: given a source node *s*, a set *R* of destination nodes, a set of constraints *C* and possibly an optimization goal, find the best feasible tree covering *s* and all nodes in *R* which satisfies *C*. The two classes of routing problems are closely related. Multicast routing can be viewed as a generalization of unicast routing in many cases. These two problem classes can be further partitioned into subclasses as follows.

Unicast Routing

For some QoS metrics such as residual bandwidth and residual buffer space, the state of a path is determined by the state of the bottleneck link. For example, in Fig. 1 the bandwidth of path $s \rightarrow i \rightarrow j \rightarrow t$ is 1, which is the bandwidth of the bottleneck link (*i, j*). For these QoS metrics, two basic routing problems can be defined. One is called *link-optimization routing*. An example is bandwidth-optimization routing, which is to find a path that has the largest bandwidth on the bottleneck link. Such a path is called the *widest path* [12]. The other problem is called *link-constrained routing*. An example is bandwidth-constrained routing, which is to find a path whose bottleneck bandwidth is above a required value. The link-optimization routing problem can be solved by a slightly modified Dijkstra's algorithm [13] or Bellman-Ford algorithm [14]. The link-constrained routing problem can easily be reduced to the link-optimization problem.

For other QoS metrics such as delay, delay jitter, and cost, the state of a path is determined by the combined state over all links on the path. For example, in Fig. 1 the delay of path $s \rightarrow i \rightarrow j \rightarrow t$ is 10, which is the total delay of all links on the path. Two basic routing problems can be defined for this type of QoS metrics. One is called *path-optimization routing*. An example is least-cost routing, which is to find a path whose total cost is minimized. The other problem is called *path-constrained routing*. An example is the delay-constrained routing, which is to find a path whose delay is bounded by a required value. Both problems can be directly solved by a Dijkstra's (or Bellman-Ford) algorithm.

Many composite routing problems can be derived from the above four basic problems (Fig. 3). For example, the bandwidth-constrained least-delay routing problem belongs to the *link-constrained path-optimization routing* problem class. It is to find the least-delay path that has the required bandwidth. This problem can be solved by a shortest path algorithm on the graph where the links violating the bandwidth constraint have been removed. There are four other problem classes that are solvable in polynomial time by a modified shortest path algorithm. They are *link-constrained link-optimization routing*, *multi-link-constrained routing*, *link-constrained path-constrained*

routing, and *path-constrained link-optimization routing*. Figure 3 gives an example for each of them.

There are two NP-complete problem classes, path-constrained path-optimization routing (PCPO) and multi-path-constrained routing (MPC), which are of particular interest. An example of PCPO is delay-constrained least-cost routing, which is to find the least-cost path with bounded delay. An example of MPC is delay-delay-jitter-constrained routing; that is, to find a path with both bounded delay and bounded delay jitter. For the above problems to be NP-complete, we have two assumptions: the QoS metrics are independent, and are allowed to be real numbers or unbounded integer numbers. If all metrics except one take bounded integer values, the problems are solvable in polynomial time by running an extended Dijkstra's (or Bellman-Ford) algorithm [15].³ If all metrics are dependent on a common metric, the problems may also be solvable in polynomial time. For example, the worst-case delay and delay jitter are functions of bandwidth in networks using Weighted Fair Queuing (WFQ) scheduling. The delay-delay-jitter-constrained routing problem is solvable in polynomial time in such networks [16].

Multicast Routing

The hierarchy of multicast routing problems is defined similarly in Fig. 3. The difference is that an optimization or a constraint must be applied to the entire tree instead of a single path. For example, bandwidth-optimization routing asks for maximizing the bandwidth of the bottleneck link of the tree. Delay-constrained routing finds a tree in which the end-to-end delay from the sender to any destination is bounded by a given value.

There are several well-known multicast routing problems. The *Steiner tree problem* is to find the *least-cost tree*, the tree covering a group of destinations with the minimum total cost over all links. It is also called the *least-cost multicast routing* problem, belonging to the *tree-optimization routing* problem class (Fig. 3). The *constrained Steiner tree problem* is to find the least-cost tree with bounded delay. It is also called the *delay-constrained least-cost routing* problem, belonging to the *tree-constrained tree-optimization routing* problem class. Finding either a Steiner tree or a constrained Steiner tree is NP-complete [17]. The delay-delay-jitter-constrained multicast routing problem belongs to the *multi-tree-constrained routing* problem class. It is also NP-complete [18] under the assumptions that the metrics under constraints are independent and are allowed to take real numbers or unbounded integer numbers. However, this problem (and any other multi-tree-constrained routing problem) is solvable in polynomial time if all metrics except one take bounded integer values. If all metrics are dependent on a common metric, the problem may also be solvable in polynomial time. Figure 4 gives examples of constrained paths and constrained trees.

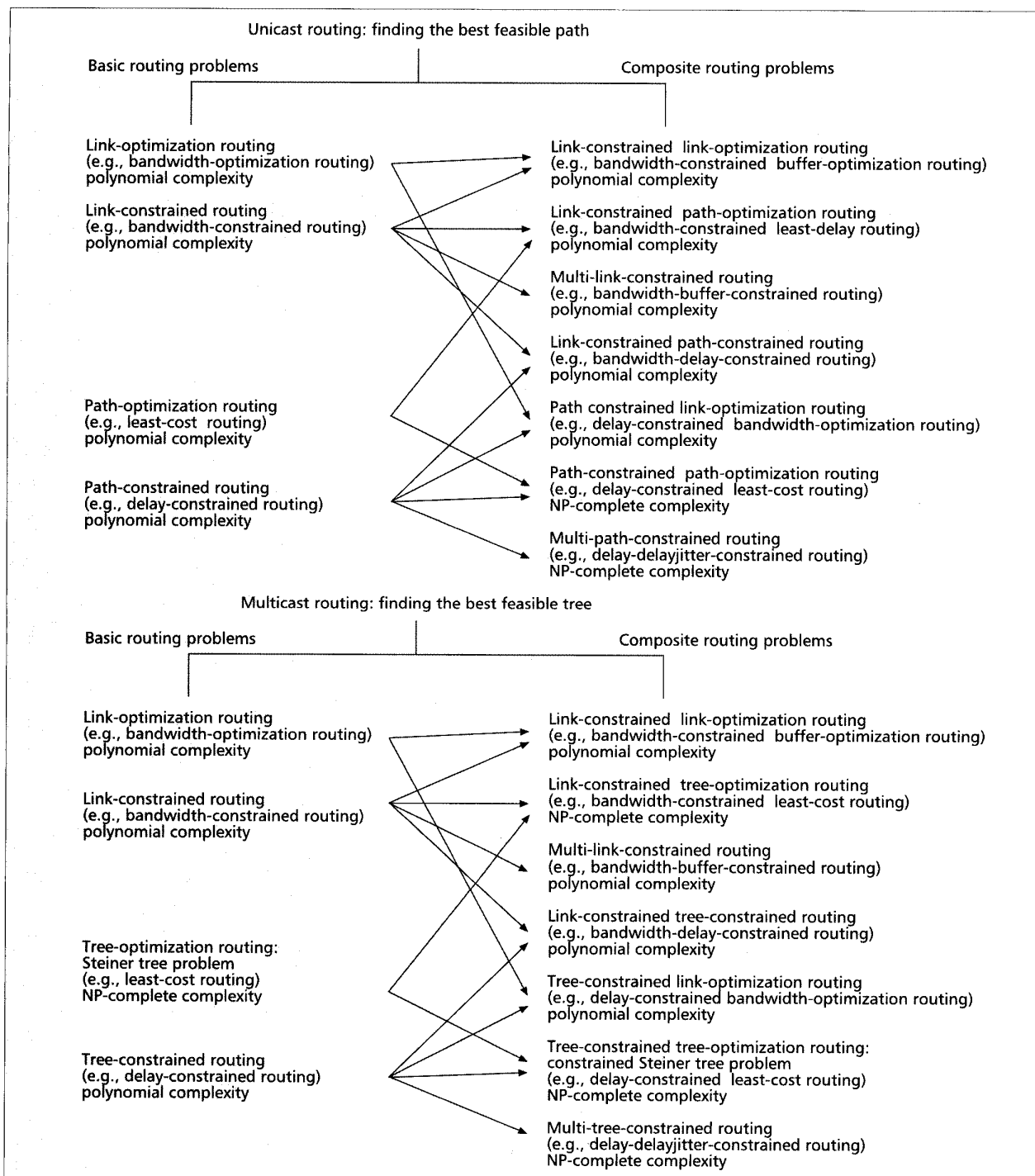
QoS Routing and Other Network Components

QoS Routing vs. Best-Effort Routing — QoS routing is different from the traditional best-effort routing. The former is normally connection-oriented with resource reservation to provide the guaranteed service. The latter can be either connection-oriented or connectionless with dynamic performance subject to the current availability of shared resources. Meeting the QoS requirement of each individual connection and reducing the call-blocking rate are important in QoS routing, while fairness, overall throughput, and average response time are the essential issues in traditional routing.

³ If all metrics (e.g., delay) except one take unbounded integer values but the maximum constraints (e.g., delay bound requirement) are bounded, the problems are also solvable in polynomial time.

QoS Routing and Resource Reservation — QoS routing and resource reservation [19, 20] are two important, closely related network components. In order to provide the guaranteed services, the required resources (CPU time, buffer, bandwidth, etc.) must be reserved when a QoS connection is established. Hence, the data transmission of the connection will not be affected by the traffic dynamics of other connections shar-

ing the common links. Before the reservation can be done, a path with the best chance to satisfy the resource requirement must be selected. That is the job of routing. While routing is decoupled from resource reservation in most existing schemes, some recent proposals combine routing and resource reservation in a single multipath message pass from the source to the destination [21].



■ Figure 3. Routing problems.

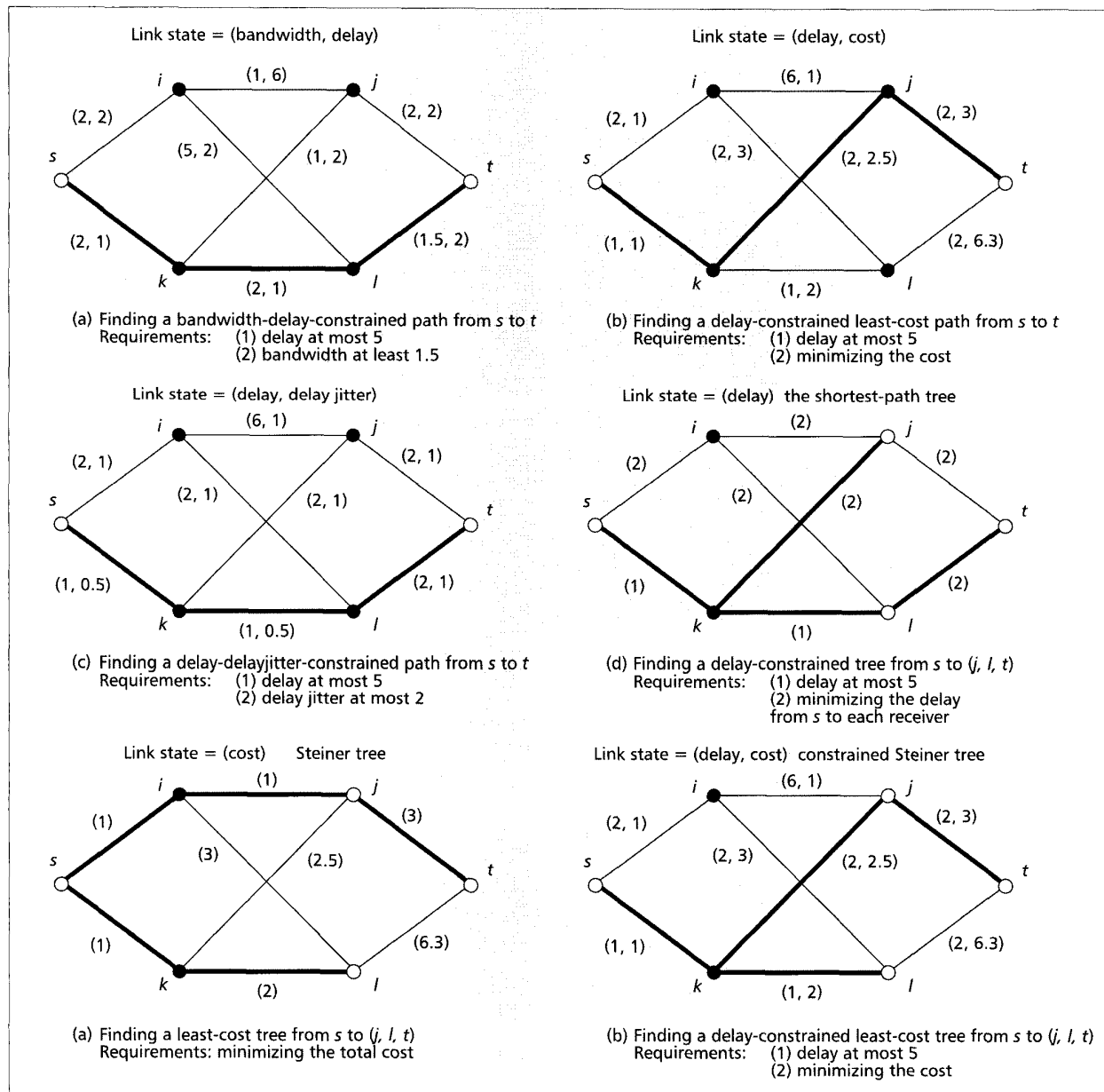
QoS Routing and Admission Control — The task of *admission control* is to determine whether a connection request should be accepted or rejected. Once a request is accepted, the required resources must be guaranteed. The admission control is often considered a by-product of QoS routing and resource reservation. If resource reservation is successfully done along the route(s) selected by the routing algorithm, the connection request is accepted; otherwise, the request is rejected.

QoS Routing and QoS Negotiation — A QoS routing algorithm may fail to find a feasible path for a new connection, either because a feasible path does not exist, or because the searching space of a heuristic approach does not cover any existing feasible path. When this happens, the system can either reject the

connection or negotiate with the application for a looser QoS constraint. QoS routing can assist the negotiation by finding the best available path and returning the QoS bounds supported. If the negotiation is successful according to the provided bounds, the best available path can be used immediately.

Routing Strategies

Routing involves two basic tasks: collecting the state information and keeping it up to date, and searching the state information for a feasible path. In order to find an optimal path which satisfies the constraints, the state information about the intermediate links between the source and destination(s) must be known. The search for feasible paths greatly depends on



■ Figure 4. Constrained paths and constrained trees.

how the state information is collected and where the information is stored.

There are three routing strategies: *source routing*, *distributed routing* and *hierarchical routing*. They are classified according to how the state information is maintained and how the search of feasible paths is carried out. In source routing, each node maintains the complete global state, including the network topology and state information of every link. Based on the global state, a feasible path is locally computed at the source node. A control message is then sent out along the selected path to inform the intermediate nodes of their precedent and successive nodes. A link-state protocol is used to update the global state at every node. In distributed routing, the path is computed by a distributed computation. Control messages are exchanged among the nodes, and the state information kept at each node is collectively used for the path search. Most distributed routing algorithms need a distance-vector protocol (or a link-state protocol) to maintain a global state in the form of distance vectors (Table 1) at every node. Based on the distance vectors, the routing is done on a hop-by-hop basis. In hierarchical routing, nodes are clustered into groups, which are further clustered into higher-level groups recursively, creating a multilevel hierarchy. Each physical node maintains an aggregated global state. This state contains detailed state information about the nodes in the same group and aggregate state information about the other groups. Source routing is used to find a feasible path on which some nodes are *logical nodes* representing groups. A control message is then sent along this path to establish the connection. When the border node of a group represented by a logical node receives the message, it uses the source routing to expend the path through the group.

Strengths and Weaknesses of Source Routing

Source routing achieves its simplicity by transforming a distributed problem into a centralized one. By maintaining complete global state, the source node calculates the entire path locally. It avoids dealing with distributed computing problems such as distributed state snapshot, deadlock detection, and distributed termination. It guarantees loop-free routes. Many source algorithms are conceptually simple and easy to implement, evaluate, debug, and upgrade. In addition, it is much easier to design centralized heuristics for some NP-complete routing problems than to design distributed ones.

Source routing has several problems. First, the global state maintained at every node has to be updated frequently enough to cope with the dynamics of network parameters such as bandwidth and delay. This makes the communication overhead excessively high for large-scale networks. Second, the link-state algorithm can only provide *approximate* global state due to the overhead concern and nonnegligible propagation delay of state messages. As a consequence, QoS routing may fail to find an existing feasible path due to the imprecision in the global state used [22]. Third, the computation overhead at the source is excessively high. This is especially true in the case of multicast routing or when multiple constraints are involved. In summary, source routing has a scalability problem. It is impractical for any single node to have access to detailed state information about all nodes and all links in a large network [10].

Strengths and Weaknesses of Distributed Routing

In distributed routing, the path computation is distributed among the intermediate nodes between the source and the destination. Hence, the routing response time can be made shorter, and the algorithm is more scalable. Searching multiple paths in parallel for a feasible one is made possible, which increases the chance of success. Most existing distributed routing algorithms [12, 23, 24] require each node to maintain global network state (distance

vectors), based on which the routing decision is made on a hop-by-hop basis. Some flooding-based algorithms do not require any global state to be maintained. The routing decision and optimization is done based entirely on the local states [25, 26].

Distributed routing algorithms which depend on the global state share more or less the same problems of source routing algorithms. The distributed algorithms which do not need any global state tend to send more messages. It is also very difficult to design efficient distributed heuristics for the NP-complete routing problems, especially in the case of multicast routing, because there is no detailed topology and link-state information available. In addition, when the global states at different nodes are inconsistent, loops may occur. A loop can easily be detected when the routing message is received by a node for the second time. However, loops generally make the routing fail because the distance vectors do not provide sufficient information for an alternative path.

Strengths and Weaknesses of Hierarchical Routing

Hierarchical routing has long been used to cope with the scalability problem of source routing in large internetworks [11, 27]. The private network-network interface (PNNI) [9] standard for routing in ATM networks is also hierarchical. Hierarchical routing scales well because each node only maintains a partial global state where groups of nodes are aggregated into logical nodes. The size of such an aggregated state is logarithmic in the size of the complete global state. Well-studied source routing algorithms are directly used at each hierarchical level to find feasible paths based on the aggregated state. Hence, hierarchical routing retains many advantages of source routing. It also has some advantages of distributed routing because the routing computation is shared by many nodes.

However, because the network state is aggregate additional imprecision is introduced, which has a significant negative impact on QoS routing [10]. Recall that a logical node in an aggregate network image may represent a large subnet with complex internal structure, and a logical link may be the abstraction of multiple physical links. Consider the aggregate network image in Fig. 2e. It is hard to estimate the end-to-end delay from *A.a.1* to a node in the group represented by *C*, because the internal structure of *C* is hidden. More specifically, although the actual delay between physical nodes in *A.c* and physical nodes in *C* may vary, there is a single delay from *A.c* to *C* in the aggregate state. Such an abstraction inevitably results in imprecision. The same thing happens to all other logical links, (*A.a.3, A.b*), (*A.a.4, A.c*), (*A.b, A.c*), (*A.b, B*) and (*B, C*).

The problem becomes more complicated when multiple QoS constraints are involved. Figure 5 shows an example. Two QoS metrics, *bandwidth* and *delay*, are considered. The pair of numbers beside a link is the residual bandwidth and the delay of the link, respectively. Four nodes, *D.1*, *D.2*, *D.3*, and *D.4*, form a group *D*. Suppose after aggregation the internal bandwidth and delay of *D* are merged into those of links (*D, F*) and (*D, G*). Consider the problem of determining the bandwidth and delay of link (*D, F*). A naive way is to find the path with the largest bandwidth from *D.1* to *D.2*, which is $P_1 = D.1 \rightarrow D.3 \rightarrow D.4 \rightarrow D.2$ with bandwidth 3. The bandwidth of (*D, F*) is the minimum of 3 and the bandwidth of (*D.2, F*), and the result is 3. Similarly, find the path with the smallest delay, which is $P_2 = D.1 \rightarrow D.2$ with delay 1. The delay of (*D, F*) is the summation of 1 and the delay of (*D.2, F*), and the result is 2. Such an optimistic approach is, however, incorrect, because P_1 and P_2 are not the same path. In general, there exist many different paths between two border nodes of a group. Some paths have better bandwidth availability, and others have smaller delay. There may not exist a path with the best properties in both terms. How to aggregate such information is still an open problem.

Unicast Routing Algorithms

We describe the unicast source, distributed, and hierarchical routing algorithms in this section. We discuss the problems and solutions, present the existing algorithms, compare them, and discuss their pros and cons. See Table 2 for a summarizing comparison. Algorithms are referred to by the authors' names and a reference to their article.

Source Routing Algorithms

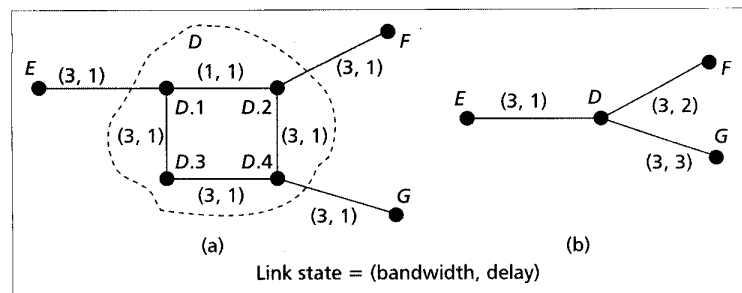
The Wang–Crowcroft Algorithm [12] — This algorithm finds a *bandwidth-delay-constrained path* by Dijkstra's shortest-path algorithm. First, all links with bandwidths less than the requirement are eliminated so that any paths in the resulting graph will satisfy the bandwidth constraint. Then, the shortest path in terms of delay is found. The path is feasible if and only if it satisfies the delay constraint.

The Ma–Steenkiste Algorithm [16] — Ma and Steenkiste showed that when a class of WFQ-like scheduling algorithms [28–31] are used, the end-to-end delay, delay-jitter, and buffer space bounds are not independent. They are functions of the reserved bandwidth, the selected path, and the traffic characteristics. Therefore, the problem of finding a path satisfying bandwidth, delay, delay-jitter, and buffer space constraints, which is NP-complete in general [2, 12], can be simplified. It can be solved by a modified version of the Bellman–Ford algorithm in polynomial time by taking those functional relationships into consideration. A much deeper study of QoS routing in rate-based scheduling networks was done recently by Orda [32].

The Guerin–Orda Algorithm [10]⁴ — Guerin and Orda studied the bandwidth-constrained and delay-constrained routing problems with imprecise network states. The model of imprecision is based on the probability distribution functions. Every node maintains, for each link l , the probability $p_l(w)$ of link l having a residual bandwidth of w units, $w \in [0, c_l]$, where c_l is the capacity of the link. The goal of bandwidth-constrained routing is to find the path that has the highest probability to accommodate a new connection with a bandwidth requirement of x units. This problem can be solved by a standard shortest path algorithm with each link l weighted by $(-\log p_l(x))$.

The goal of delay-constrained routing is to find a path that has the highest probability to satisfy a given end-to-end delay bound. Suppose every node maintains, for each link l , the probability $p_l(d)$ of link l having a delay of d units, where d ranges from zero to the maximum possible value. It is NP-hard to find the path that has the highest probability of satisfying a given delay constraint [10], but various special cases (e.g., symmetric networks and tight constraints) can be solved in polynomial time. Heuristic algorithms were proposed for the NP-hard problem. The idea is to transform a global constraint into local constraints. More specifically, it splits the end-to-end delay constraint among the intermediate links in such a way that every link in the path has an equal probability of satisfying its local constraint. The heuristics then try to find the path with the best multiplicative probability over all links.

The Guerin–Orda algorithm works with imprecise information and is suitable to be used in hierarchical routing. One of the heuristic algorithms was extended by the authors to make routing based on the aggregate network state of the hierarchical model. A further study of QoS routing with imprecise state based on the probability model was done by Lorenz and Orda [33].



■ Figure 5. a) The internal state; b) incorrect aggregation on link (D, F).

The Chen–Nahrstedt Algorithm [15] — Chen and Nahrstedt proposed a heuristic algorithm for the NP-complete multi-path-constrained routing problem. We already know that if all metrics except one take bounded integer values, the multi-path-constrained routing problem is solvable in polynomial time. Consider delay-cost-constrained routing. The idea is to map the cost (or delay) of every link from an unbounded real number to a bounded integer. This reduces the original NP-complete problem to a simpler problem solvable in polynomial time. Let C be the cost requirement and x a small integer. The algorithm first maps the cost of every link to an integer bounded by $x + 1$. Real numbers in $[0, C]$ are mapped into integers in $[0, x]$, real numbers in $(C, \infty]$ are mapped to $x + 1$, and the cost bound C is mapped to x . See Fig. 6 for an example. The new problem with the link cost bounded by $x + 1$ can be solved in polynomial time by an extended Dijkstra's algorithm (EDSP) or an extended Bellman–Ford algorithm (EBF) [15]. It was proved that a feasible path of the new problem must also be a feasible path of the original problem. The performance of the algorithm is tunable by choosing the value of x : a larger x results in a higher probability of finding a feasible path and a higher overhead.

The Awerbuch et al. Algorithm [42] — Awerbuch et al. proposed a *throughput-competitive routing* algorithm for bandwidth-constrained connections. The algorithm tries to maximize the amortized (average) throughput of the network over time. It combines the functions of admission control and routing. Every link is associated with a cost function that is exponential to the bandwidth utilization. A new connection is admitted into the network only if there exists a path whose *accumulated cost over the duration of the connection* does not exceed the *profit* measured by the bandwidth-duration product of the connection. It was proved that such a path satisfies the bandwidth constraint. Let T be the maximum connection duration and v the number of nodes in the network. The algorithm achieves a throughput that is within $O(\log vT)$ factor of the highest possible throughput achieved by the best off-line algorithm that is assumed to know all of the connection requests in advance. The competitive routing for connections with unknown duration was studied in [43]. A survey for the competitive routing algorithms was done by Plotkin [44].

Summary — All the above algorithms require global state to be maintained at every node. Most algorithms transform the routing problem to a shortest-path problem and then solve it by Dijkstra's or the Bellman–Ford algorithm. We summarize the distinctive properties of some algorithms as follows: the Ma–Steenkiste algorithm provides a routing solution to rate-based networks; the Guerin–Orda algorithm works with imprecise information, and hence is suitable to be used in

⁴ The Guerin–Orda algorithm was designed to be used in hierarchical routing, although we present it as an independent source routing algorithm in this article.

hierarchical routing; the performance of the Chen–Nahrstedt algorithm is tunable by trading overhead for success probability; the Awerbuch *et al.* algorithm takes the connection duration into account, which allows more precise cost-profit comparison. All the above algorithms are executed at the connection arrival time on a per-connection basis, which may cause the overall computational overhead to be excessively high. Path precomputation and *caching* were studied to make a trade-off between processing overhead and routing performance [45–48].

Distributed Routing Algorithms

The Wang–Crowcroft Algorithm [12] — Wang and Crowcroft proposed a *hop-by-hop* distributed routing scheme. Every node pre-computes a forwarding entry for every possible destination. The forwarding entry, which is updated periodically, stores the *next hop* on the routing path to the destination. After the forwarding entries at every node are computed, the actual routing simply follows the entries.

Given two end nodes, the path with the maximum bottleneck bandwidth is called the *widest path*. If there are several such paths, the one with the smallest delay is called the *shortest-widest path*. A link-state protocol is used to maintain complete global state at every node. Based on the global state, the forwarding

entry for the *shortest-widest* path to each destination is computed by a modified Bellman–Ford (or Dijkstra’s) algorithm [12]. A routing path is the combination of the forwarding entries indexed by the same destination at all intermediate nodes. The path is loop-free if the state information at all nodes is consistent. However, in a dynamic network the path may have a loop due to the contradicting state information at different nodes.

The Salama et al. Algorithm [23] — Salama *et al.* proposed a distributed heuristic algorithm for the NP-complete delay-constrained least-cost routing problem. A *cost vector* and a *delay vector* are maintained at every node by a distance-vector protocol. The cost (delay) vector contains for every destination the next node on the least-cost (least-delay) path. A control message is sent from the source toward the destination to construct a delay-constrained path. Any node *i* at the end of the partially constructed path can select one of only two alternative outgoing links. One link (*i, j*) is on the least-cost path directed by the cost vector, and the other (*i, k*) is on the least-delay path directed by the delay vector. Link (*i, j*) has the priority to be chosen, as long as adding the least-delay path from *j* to the destination does not violate the delay constraint.

Loops may occur as the control message chooses the least-

Algorithm	Solving routing problem	Routing strategy	Time complexity	Communication complexity	
				Maintaining state	Routing
Wang–Crowcroft [12]	Bandwidth-delay-constrained	Source	$O(v \log v + e)$	Global	Zero
Ma–Steenkiste [16]	Bandwidth-constrained	Source	$O(v \log v + e)$	Global	Zero
	Multi-constrained ¹	Source	$O(kve)$ ¹	Global	Zero
Guerin–Orda [10]	Bandwidth-constrained	Source	$O(v \log v + e)$	Imprecise global	Zero
	Delay-constrained	Source	Polynomial ²	Imprecise global	Zero
Chen–Nahrstedt [15]	Bandwidth-cost-constrained	Source	$O(xve)$ ³	Global	Zero
Wang–Crowcroft [12]	Bandwidth-optimization	Distributed	$O(ve)$	Global	$O(v)$
Salama <i>et al.</i> [23]	Delay-constrained least-cost	Distributed	$O(v^3)$	Global	$O(v^3)$ ⁴
Sun–Landgendorfer [24]	Delay-constrained least-cost	Distributed	$O(v)$	Global	$O(v)$
Cidon <i>et al.</i> [21]	Generic ⁵	Distributed	$O(e)$	Global	$O(e)$ ⁶
Shin–Chou [26]	Delay-constrained	Distributed	$O(e)$	Local	$O(e)$
Chen–Nahrstedt [25]	Generic ⁵	Distributed	$O(e)$	Local	$O(e)$
PNNI [9]	Generic ⁵	Hierarchical	Polynomial ⁷	Aggregated	$O(v)$

v is the number of nodes, and *e* is the number of edges.

After a source routing algorithm selects a path, a control message needs to be sent along the path to establish the connection, which has a worst-case communication overhead of $O(v)$.

¹ Ma and Steenkiste studied routing with constraints on delay, delay jitter, and buffer space in rate-based scheduling networks. *k* in the time complexity is the number of all possible residual bandwidth a link may have.

² Heuristics with different assumptions have different polynomial time complexities.

³ *x* is a constant in the algorithm. A larger *x* results in a higher probability of finding a feasible path and a higher overhead.

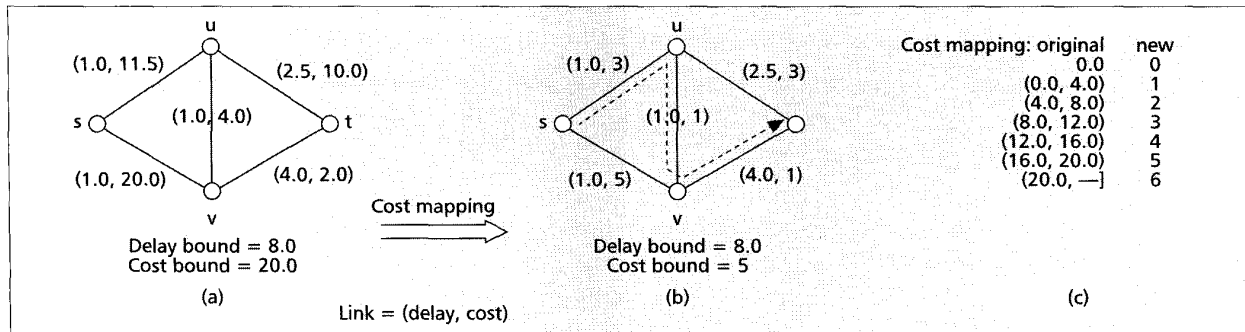
⁴ It was shown that the average overhead is substantially less than the worst-case overhead.

⁵ A routing framework was proposed from which algorithms on different QoS constraints can be derived.

⁶ Variants of the algorithm may have higher worst-case overhead.

⁷ The time complexity of a hierarchical routing algorithm depends on which source routing algorithm is used to route the connection through every group.

■ Table 2. Unicast routing algorithms.



■ Figure 6. Chen–Nahrstedt heuristic ($C = 20.0$ and $x = 5$): a) The original problem is to find a path from s to t such that the delay is bounded by 8.0 and the cost is bounded by 20.0 ; b) the costs of links are mapped to integers in $[1..6]$. For link (s, u) , the cost 11.5 is mapped to 3 . The new problem is find a path from s to t such that the delay is bounded by 8.0 and the cost is bounded by 5 . A feasible path is $s \rightarrow u \rightarrow v \rightarrow t$, which as expected is also a feasible path for the original problem. c) The cost-mapping table.

cost and least-delay paths alternatively. A loop is detected if the control message visits a node twice; whenever this happens, the routing process is rolled back until reaching a node from which the least-cost path was followed. The routing process resumes from there by changing the next hop along the least-delay path. It was proven that such a mechanism removes all loops, provided that the delay and cost vectors at all nodes are up to date (or at least consistent), a condition that does not hold sometimes in a dynamic network.

The Sun–Langendorfer Algorithm [24] — Sun and Langendorfer improved the worst-case performance of the Salama *et al.* algorithm by avoiding instead of detecting and removing loops. A control message is sent to construct the routing path. The message travels along the least-delay path until reaching a node from which the delay of the least-cost path satisfies the delay constraint. From that node on, the message travels along the least-cost path all the way to the destination. The difference between the Sun–Langendorfer and Salama *et al.* algorithms is illustrated in Fig. 7. It was proven that the algorithm constructs loop-free paths, provided that the state information at all nodes is updated (or consistent). In a dynamic network, different nodes may have inconsistent information. The least-cost (least-delay) path computed based on such inconsistent information may contain a loop, which makes the control message unable to reach the destination.

*The Cidon *et al.* Algorithm [21]* — The distributed multi-path routing algorithms proposed by Cidon *et al.* combine the process of routing and resource reservation. Every node maintains the topology of the network and the cost of every link. When a node wishes to establish a connection with certain QoS constraints, it finds a subgraph of the network which contains links that lead to the destination at a “reasonable” cost. Such a subgraph is called a *diroute*. A link is *eligible* if it has the required resources. Reservation messages are flooded along the eligible links in the diroute toward the destination and reserve resources along different paths in parallel. When the destination receives a reservation message, a routing path is established. The algorithm releases resources from segments of the diroute as soon as it learns that these segments are inferior to another segment. Variants of the above algorithm were proposed to make a trade-off between routing time and path optimality. Reserving resources on multiple paths makes the routing faster and more resilient to the dynamic change of network state. However, it also increases the level of resource contention.

The Shin–Chou Algorithm [26] — Shin and Chou proposed a distributed routing scheme for establishing delay-constrained

connections. No global state need be maintained at any node. The algorithm floods routing messages from the source toward the destination. Each message accumulates the total delay of the path it has traversed so far. When a routing message is received by an intermediate node, the message is forwarded only when one of the following two conditions is satisfied:

- It is the first such message received by the node.
- It carries a better accumulated delay than the previously received messages.

If either condition is true, the message will be forwarded along the outgoing links whose delay plus the message’s accumulated delay does not exceed the end-to-end delay requirement. Once a message reaches the destination, it finds a delay-constrained path, which is the one it has traversed. It was shown that when certain scheduling policies [49] are used and the routing messages are set to the appropriate priority, there will be at most one message sent along every link. Another flooding-based routing algorithm was proposed by Hou [50]. It routes virtual circuits with delay requirements in ATM networks.

The Chen–Nahrstedt Algorithm

Selective Probing [25] — Chen and Nahrstedt proposed a distributed routing framework based on *selective probing*. After a connection request arrives, probes are flooded selectively along those paths which satisfy the QoS and optimization requirements. Every node only maintains its local state, based on which the routing and optimization decisions are made collectively in the process of probing. As in the Shin–Chou algorithm, each probe arriving at the destination detects a feasible path.

Algorithms were derived from the framework to route connections with a variety of QoS constraints on bandwidth, delay, delay jitter, cost, and their combinations. Several techniques were developed to overcome the high communication overhead of the Shin–Chou algorithm. First, probes are only allowed to be forwarded to a *subset* of outgoing links selected based on topological distance to the destination. Second, *iterative probing* is used to further reduce the overhead. At the first iteration, probes are sent only along the shortest paths. If the first iteration fails, probes are allowed to be sent along paths with increasing lengths in the following iterations. Simulation shows that with two iterations the Chen–Nahrstedt algorithm achieves substantial overhead reduction.

Ticket-Based Probing [51] — If every node maintains a global state, which is allowed to be imprecise, the *ticket-based probing* is used to improve the performance of selective probing. A certain number of tickets is issued at the source according to the contention level of network resources. Each probe must contain at least one ticket in order to be valid. Hence, the

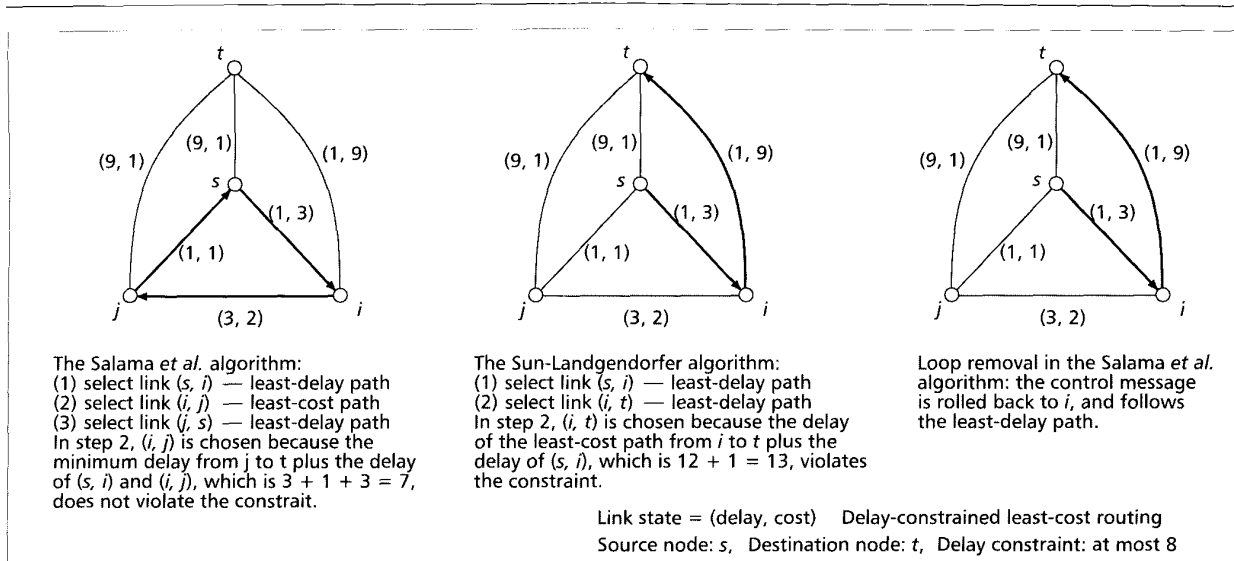


Figure 7. The Salama *et al.* algorithm vs. the Sun-Landendorfer algorithm.

maximum number of probes is bound by the total number of tickets, which limits the maximum number of paths to be searched. The algorithm utilizes the imprecise state at intermediate nodes to guide the limited tickets (the probes carrying them) along the best possible paths to the destination. In such a way, the probability of finding a feasible path is maximized with limited probing overhead.

Summary — The distinctive properties of the above algorithms are summarized as follows:

- The Salama *et al.* and Sun-Landendorfer algorithms provide efficient distributed solutions to the NP-complete delay-constrained least-cost routing problem.
- The Cidon *et al.* Shi-Chou, and Chen-Nahrstedt algorithms are multipath routing algorithms.⁵
- The Cidon *et al.* algorithm combines routing with resource reservation.
- The Shi-Chou and Chen-Nahrstedt's selective probing algorithms require only the local state to be maintained at each node.
- Chen-Nahrstedt's iterative probing substantially reduces routing overhead at the cost of longer routing time.

Hierarchical Routing Algorithms

PNNI — PNNI is a hierarchical link-state routing protocol [9]. Its hierarchical model was discussed earlier. We use an example to illustrate the routing process. The network in Fig. 8a has a two-level hierarchy with three groups. The aggregated topology maintained at $A.1$, $B.1$ and $C.1$ are shown in Fig. 8b, c, and d, respectively. Suppose every link has an available bandwidth of one. Consider a connection request arriving at $A.1$ with destination $C.2$. Let the bandwidth requirement be one. The routing process is described as follows. Based on the aggregated state, the source node $A.1$ finds a path ($A.1 \rightarrow A.2$) within its group and a logical path ($A \rightarrow B \rightarrow C$) on the higher hierarchy level. The logical path, together with the destination $C.2$, is sent to the next group B on the path. When the border node $B.1$ receives the information, it selects a path ($B.1 \rightarrow B.2 \rightarrow B.3$) within its group and then passes the logical path and the destination to group C . Finally, the border node $C.1$ of the destination group completes the routing by selecting $C.1 \rightarrow C.2$. It may happen that a link on the selected path does not have sufficient resources. Figure 8e gives an example, where link $B.3 \rightarrow B.2$ does not have enough band-

width for the connection due to traffic dynamics. In this case, the routing process is cranked back to $B.1$ and resumes with an alternative path $B.1 \rightarrow B.2$.

Multicast Routing Algorithms

Most existing work on multicast routing focuses on the following problems:

- Bandwidth-constrained multicast routing
- Delay-constrained multicast routing
- Delay-constrained least-cost multicast routing (the constrained Steiner tree problem)
- Delay-delay-jitter-constrained multicast routing

We describe the algorithms in this section. A summarizing comparison can be found in Table 4.

Source Routing Algorithms

MOSPF [34] — MOSPF is a multicast extension of the unicast link-state protocol OSPF [3]. It was based on Deering's work [52]. In addition to a global state, the protocol maintains at every node the membership information of every multicast group in the routing domain. Group membership changes in a subnetwork are detected by a local router, which broadcasts the information to all other nodes. Given full knowledge of network state and group membership, any node can compute the shortest-path multicast tree from a source to a group of destinations, using Dijkstra's algorithm. Such a protocol can be easily used for delay-constrained multicast routing.

Steiner Tree Problem — A Steiner tree is the least-cost tree that spans a given subset of nodes. Strictly speaking, finding a Steiner tree is not a QoS routing problem. However, the heuristics for constructing a Steiner tree have a direct impact on how to construct a constrained Steiner tree. In the following, we briefly discuss two algorithms. A nice survey on the Steiner problem can be found in [53].

The Kou *et al.* Algorithm [35] — In this algorithm, a network is abstracted to a complete graph, where the nodes represent the source and destinations, and the edges represent the shortest paths between these nodes. Prim's algorithm [54] is

⁵ Search multiple paths for a feasible one.

used to construct a minimum spanning tree in the complete graph. Then the Steiner tree of the original network is obtained by expending the edges of the minimum spanning tree into the shortest paths they represent. Any loops caused by the expansion are removed.

The Takahashi-Matsuyama Algorithm [36] — This algorithm finds a Steiner tree by an incremental approach called nearest destination first (NDF). Initially, the nearest destination (in terms of cost) to the source is found and the least-cost path between them selected. Then at each iteration the *nearest unconnected destination* to the partially constructed tree is found and added into the tree. This process is repeated until all destinations are included in the tree.

The Constrained Steiner Tree Problem — The problem of finding a delay-bound least-cost multicast tree, called a *constrained Steiner tree*, is NP-complete [37]. Heuristic source routing algorithms were proposed for this problem [37-40]. A performance evaluation of these algorithms was done by Salama *et al.* through extensive simulation [17].

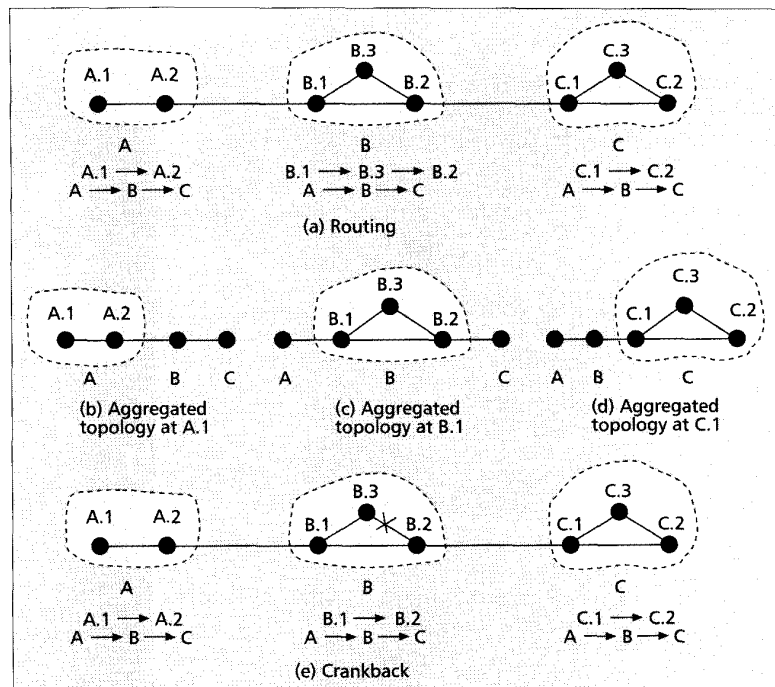
The Kompella *et al.* Algorithm [37] — A source routing heuristic was proposed by Kompella *et al.* to construct a constrained Steiner tree. The first step is to create a complete graph, where the nodes represent the source and destinations, and the edges represent the delay-constrained least-cost paths between these nodes. The link delays are assumed to be integers, and the delay constraint always bounded, so such a complete graph can be constructed in polynomial time. The second step is to construct a delay-constrained spanning tree of the complete graph. Starting with the source node, the tree is incrementally expanded by adding an edge each time until every destination node is included. The selected edge is the one that:

- Connects a node in the tree and a node outside of the tree
- Does not violate the delay constraint
- Minimizes a selection function

Two selection functions are considered. One is simply the cost of the edge, and the other tries to make a trade-off between minimizing the cost and minimizing the delay. The third step is to expand the edges of the constrained spanning tree into the delay-constrained least-cost paths they represent. Any loops caused by this expansion are removed.

The Sun-Langendoerfer Algorithm [38] — Sun and Langendoerfer proposed an algorithm which constructs an approximated constrained Steiner tree by Dijkstra's algorithm. It first computes the shortest-path tree in terms of cost. Namely, the cost of every path in the tree from the source to a destination is minimized. Then the tree is modified to satisfy the delay constraint. If the end-to-end delay to any destination in the tree violates the delay constraint, the minimum-delay path is used to replace the minimum-cost path. The advantage of the algorithm is its low time complexity, $O(v \log v)$, which is the same as Dijkstra's algorithm.

The Widyono Algorithm [39] — Widyono proposed several heuristic algorithms for the constrained Steiner tree problem. The one with the best performance is called the *constrained*



■ Figure 8. An example of PNNI routing.

adaptive ordering heuristic. At each step, a constrained Bellman-Ford algorithm is used to find a delay-constrained least-cost path from the source to a destination that is not yet in the tree. The found path as well as the destination are then inserted into the tree. The cost of links in the tree is set to zero. The above process repeats until the tree covers all destinations.

The Zhu *et al.* Algorithm [40] — Zhu *et al.* proposed a source routing heuristic to construct the constrained Steiner tree. The algorithm allows variable delay bounds on destinations. A shortest-path tree in terms of delay is first constructed by Dijkstra's algorithm. If the delay constraint cannot be satisfied for any destination, it must be renegotiated; otherwise, the algorithm proceeds by iteratively refining the tree for lower cost. The basic idea is to replace a path in the tree by another path with lower cost unless such a replacement cannot be found. Figure 9 gives a replacement example. Heuristics were proposed for finding such a replacement. The algorithm always finds a delay-constrained tree (probably not least-cost), if one exists, because it starts with a shortest-path tree.

The Rouskas-Baldine Algorithm [18] — Rouskas and Baldine proposed a heuristic for constructing a delay-delay-jitter-constrained multicast tree. The tree must have bounded delay along the paths from the source to the destinations and bounded variation among the delays along these paths. The shortest-path tree T_0 in terms of delay is first constructed by Dijkstra's algorithm. If the tree does not meet the delay jitter constraint, the algorithm finds the largest-delay path in T_0 from the source to a destination, and starts from that path to incrementally construct a feasible tree. At each iteration, a "good" path from a node in the tree to a destination out of the tree is found and added into the tree. The path must be completely disjoint from the tree, and must not cause the tree to violate the constraints. The above process repeats until all destinations are included in the tree. The authors showed that the heuristic demonstrates good average-case behavior in terms of the maximum interdestination delay variation.

Summary — All the above algorithms require global state to be maintained at every node. Most heuristic algorithms for the NP-complete multicast routing problems construct a constrained tree incrementally by adding one destination into the tree each time based on certain selection criteria. The Kou *et al.* and Kompella *et al.* algorithms reduce the original problem to a spanning tree problem by constructing a logical complete graph among the source node and the destination nodes. The Zhu *et al.* algorithm iteratively refines the multicast tree by replacing paths in the tree for lower cost. The Rouskas-Baldine algorithm constructs a multicast tree with both bounded delay and bounded delay jitter, which is very useful in interactive audio-visual communication such as teleconferencing. Among the four algorithms for the constrained Steiner tree problem, Salama's simulation [17] showed that

- The Zhu *et al.* algorithm achieves the best average performance in terms of minimizing the cost of the tree.
- The Sun-Langendoerfer algorithm has the shortest execution time.

Distributed Routing Algorithms

The Kompella et al. Algorithm [41] — Kompella *et al.* proposed a distributed heuristic algorithm for constructing the constrained Steiner tree [41]. The algorithm requires every node to maintain a distance vector, storing the minimum delay to every other node. Starting with the source node, the algorithm constructs the multicast tree iteratively by adding a link into the tree each time. Each iteration of the algorithm consists of three phases of message passing. In the first phase, the source node broadcasts a Find message down the partially constructed tree. When a node receives the message, it finds the adjacent link that leads to a destination out of the tree, does not violate the delay constraint, and minimizes a selec-

tion function. In the second phase the selected links are sent to the source node, where the best link l that minimizes the selection function is chosen. In the third phase an ADD message is sent to add l to the tree. This procedure continues until every destination is included in the tree. The above algorithm requires intensive multipass message exchange. The worst-case message complexity is $O(v^3)$.

The Chen-Nahrstedt Algorithm [25] — Chen and Nahrstedt extended their distributed unicast routing algorithms [25] (discussed earlier) for multicast routing. Probes (routing messages) are flooded from the source toward the destinations of a multicast group. Probes proceed only along the paths that lead to at least one destination and have sufficient resources to guarantee the end-to-end QoS. As probes traverse toward a group of destinations, a multicast tree is built in a distributed manner. Every node maintains only its local state. The worst-case message complexity is $O(e)$ for constructing the entire tree. This approach only works for a multicast group whose membership is fixed and a priori known. The dynamic membership problem is handled by receiver-initiated probing. When a new destination joins in a multicast group, it sends probes toward the multicast tree. Probes proceed only along the paths that do not violate QoS and optimization requirements. Once a probe reaches any node in the multicast tree, a feasible extension of the tree is found. The worst-case message complexity of the above receiver-initiated probing is $O(e)$ for a single receiver.

The Carlberg-Crowcroft Algorithm [55] — The *spanning-joins* approach was proposed by Carlberg and Crowcroft for the construction of multicast trees across different domains [55]. A new group member broadcasts a *join-request* message. When

Algorithm	Solving routing problem	Routing strategy	Time complexity	Communication complexity	
				Maintaining state	Routing
MOSPF [34]	Least-delay	Source	$O(v \log v)$	Global	Zero
Kou <i>et al.</i> [35]	Least-delay	Source	$O(gv^2)$	Global	Zero
Takahashi-Matsuyama [36]	Least-delay	Source	$O(gv^2)$	Global	Zero
Kompella <i>et al.</i> [37]	Delay-constrained least-cost	Source	$O(v^3 \Delta)^1$	Global	Zero
Sun-Landendorfer [38]	Delay-constrained least-cost	Source	$O(v \log v + e)$	Global	Zero
Widyono [39]	Delay-constrained least-cost	Source	Exponential ²	Global	Zero
Zhu <i>et al.</i> [40]	Delay-constrained least-cost	Source	$O(kv^3 \log v)^3$	Global	Zero
Rouskas-Baldine [18]	Delay-constrained least-cost	Source	$O(klgv^4)^3$	Global	Zero
Kompella <i>et al.</i> [41]	Delay-constrained least-cost	Distributed	$O(v^3)$	Global	$O(v^3)$
Chen-Nahrstedt [25]	Generic	Distributed	$O(ge)$	local	$O(ge)$

v is the number of nodes, e is the number of edges, and g is the number of destinations. After a source routing algorithm constructs a multicast tree, a control message needs to be sent down the tree to establish the connection, which has a worst-case communication overhead of $O(e)$.

¹ Δ is the delay requirement. The time complexity is polynomial if Δ is a bounded integer.

² The Widyono algorithm uses the constrained Bellman-Ford (CBF) algorithm. Widyono pointed out that there are cases where the running time of CBF grows exponentially. However, simulation shows that its average performance is comparable to other algorithms that construct constrained Steiner trees [17].

³ k and l are constants in the algorithm. A larger k (or l) results in a higher probability of finding a feasible tree and a higher overhead.

■ Table 4. Multicast routing algorithms.

an on-tree node receives the message, it sends a unicast reply message back to the new member. The path of the reply message is determined by the existing unicast routing algorithm. The message may collect the QoS properties and resource availability of the path as it traverses. The new member may receive multiple reply messages that correspond to multiple candidate paths connecting to the multicast tree. It selects the best path according to the QoS information carried by the received reply messages. *Reverse path multicasting*, a *time-to-live field*, and *directed spanning joins* are used to reduce the message overhead [55]. Excellent recent work was done by Faloutsos *et al.* [56], which improves the performance of spanning-joins by the help of a *Manager router*. Based on the topology information, the Manager router selects a subset of the on-tree nodes to send the reply messages without receipt of the join-request message.

The probes in the Chen-Nahrstedt algorithm are only allowed to be sent along the paths that have sufficient resources to support the required QoS [25]. Hence, a feasible path is determined when a probe reaches the multicast tree. The join-request messages in the Carlberg-Crowcroft algorithm are sent along paths that may or may not have enough resources; hence, multiple candidate paths are used. Reply messages are sent along them to collect the QoS information, based on which the new member determines whether a feasible path exists.

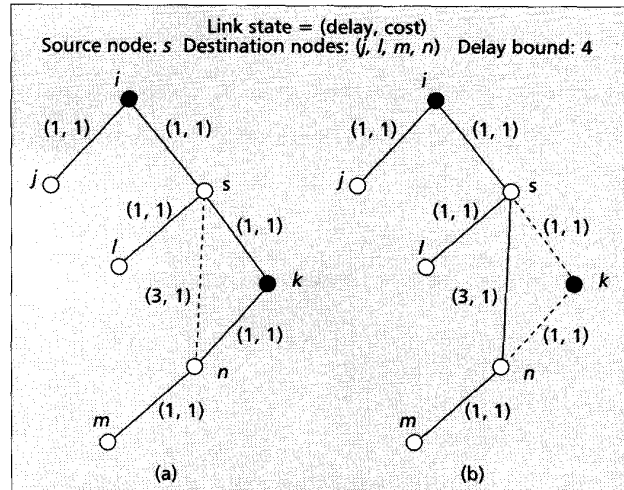
Summary — The Kompella *et al.* algorithm provides a distributed solution to the NP-complete constrained Steiner tree problem. Its communication overhead is high, and every node needs to maintain global state. The Chen-Nahrstedt algorithm requires every node to maintain only the local state. It is suitable to construct the shortest-path tree, but not the constrained Steiner tree, because the probes search for the shortest paths individually without cooperation to reduce the overall cost.

Future Directions

Efficient Routing Algorithms — Most source heuristic algorithms for the NP-complete routing problems (Fig. 3) are not scalable due to prohibitively high time complexity, especially in the case of multicast routing. New efficient algorithms are required to achieve a good balance between the computation time and the connection-success ratio so that the time complexity can be reduced to the shortest-path computation range while the success ratio is still acceptable [17].

Routing with Imprecise State Information — Most existing routing algorithms assume the availability of precise state information. However, state information is inherently imprecise in a distributed network environment. The imprecision directly affects routing performance. Therefore, the design of routing algorithms for large networks should take the information imprecision into consideration [10, 33, 51].

Distributed and Hierarchical Routing — Source routing based on the complete global state is generally not scalable because of the following reasons. The communication overhead to maintain the global state is proportional to the size of the network and the frequency of broadcasting local states. The storage overhead to store the state is proportional to the size of the network. The computation overhead of calculating the feasible paths is polynomial to the size of the network and proportional to the arrival frequency of connection requests. The precision of the global state at a node is in inverse proportion to the diameter of the network and the frequency of broadcasting



■ Figure 9. Solid links are in the multicast tree; dashed links are not in the tree. Path $s \rightarrow k \rightarrow n$ in a) is replaced by path $s \rightarrow n$ in b); the cost is reduced by 1.

local states. As a network grows large, the communication, storage, and computation overhead grows accordingly. Reducing the updating frequency does not solve the problem because the precision of the global state will decrease.

Distributed and hierarchical algorithms offer solutions for the scalability problem. In particular, the distributed algorithm based on selective probing [25] uses only local states, and no shortest-path computation is conducted at a single node. The ticket-based probing algorithm [51] works with imprecise state information, which allows relatively infrequent state updates. Hierarchical routing provides a clean solution to the scalability problem. It maintains an aggregate global state whose size is logarithmic to the network size if the (logical) nodes are clustered into groups of roughly uniform size. However, state aggregation leads to further imprecision, especially when multiple QoS metrics are involved. The design and evaluation of hierarchical routing algorithms should take this into account.

Multipath Routing — When the traffic load is light, network resources are readily available. QoS routing is of less importance in searching feasible paths but of more importance in balancing the traffic. Balanced traffic distribution helps to increase the call admission ratio of future connections and improve the response time of the best-effort traffic. However, when the network load is heavy and dynamic, efficient algorithms for finding feasible paths are critical. *Multipath routing* can be used to increase the probability of accepting a connection under resource contention. There are two interpretations of multipath routing.

One interpretation is to search multiple paths for a feasible one. PNNI [9] uses *crankback* to search multiple paths *sequentially*. When the selected path does not meet the requirement, the routing process is cranked back and resumes with an alternative path. This approach works well with network dynamics; the disadvantage is longer routing time. *Parallel* multipath routing was proposed to overcome this problem [21, 26]. Routing messages are sent along multiple paths in parallel and reserve resources along the way. If more than one message arrives at the destination, the best path is selected and resources reserved on the other paths are released. An alternative approach is to reserve resources only on a primary path. The messages sent along the other (secondary) paths only check resource availability. If the reservation on the pri-

mary path fails, a secondary path is picked for resource reservation.

The other interpretation of multipath routing is to select a set of paths instead of a single one for a connection. When there does not exist a feasible path with sufficient resources, the algorithm tries to find multiple paths whose combined resources satisfy the requirement [57–59]. Transmitting contiguous data (audio and video) along multiple paths raises the problem of synchronization. In addition, it demands more buffer space at the receiving end to absorb the delay jitter between different paths.

Routing QoS and Best-Effort Traffic — QoS and best-effort traffic coexist in most real-world networks [61, 62]. A primary task of routing is to maximize the *resource efficiency*, which is measured by two goals. One goal is to maximize the number of QoS flows that are admitted into the network, which is equivalent to minimizing the *call-blocking ratio*. The other goal is to optimize the throughput and responsiveness of best-effort traffic. The two goals may contradict each other because:

- The first goal considers only QoS traffic.
- The second goal considers only best-effort traffic.
- However, the two types of traffic may have very different distributions.

Generally speaking, QoS traffic will not be affected by best-effort traffic due to resource reservation. However, the throughput of the best-effort traffic will suffer if the overall traffic is misjudged. For example, links with light QoS traffic may have heavy best-effort traffic. By many QoS routing algorithms, these links are often considered as good candidates for new QoS flows, but this causes the already congested best-effort traffic to become even more congested.

Rerouting — There are a number of situations where rerouting is desired. First, the routes of the connections are typically selected based on network resource availability at the times when the requests arrive. Long paths are often assigned when resource contention occurs. However, as new connections are established and existing connections are torn down upon completion, the network state changes locally and globally, which makes the routes of the remaining connections less optimal [60]. Routes with light (heavy) traffic at the beginning may become congested (lightly loaded) later. Shorter paths for some existing connections may become available. Rerouting helps to balance network traffic on the fly and improve the resource efficiency. Second, when a feasible path for a new connection does not exist, instead of rejecting the connection, it is often possible to reroute some existing connections in order to make room for the new one. Rerouting is especially useful when connections have different priorities. A new connection with a higher priority will preempt the resources held by existing connections. Instead of throwing the preempted connections out of the network, we can reroute them to other paths. Rerouting should not be done too frequently in order to avoid excessive overhead and the oscillation of shifting the traffic from one part of the network to another.

Integration with Other Network Components — Routing must work with other network components in order to provide guaranteed services. The design of routing algorithms must consider how global state is maintained, resources reserved, and data packets scheduled. Different scheduling policies make different requirements on routing algorithms, and often provide special properties to simplify the routing problems [16]. For example, when rate-based scheduling policies are

used, the end-to-end delay constraint can be transformed into a bandwidth constraint. The following properties are desired for the routing component in an integrated network system.

Generality — Multimedia applications tend to have diverse QoS requirements on bandwidth, delay, delay jitter, cost, and so on. From a network designer's point of view, it would be beneficial to develop a *generic* routing algorithm instead of implementing different routing algorithms for different types of QoS requirements independently. The generic algorithm captures the common messaging and computational structure. Various concrete algorithms are derived from the generic algorithm by specifying the QoS-dependent open components [25].

Extensibility — As the network infrastructure evolves and capacity increases, new applications are made possible. It requires the routing algorithms to adapt in order to accommodate new service types. It is important to design extensible algorithms and make them adapt to new applications, because the networks become increasingly complex and the deployment of new routing algorithms is very costly and problem-prone.

Simplicity — The simplicity of a routing algorithm in terms of time/logical complexity often allows efficient implementation, debugging and evaluation. It also makes the algorithm easier to understand, maintain, and upgrade.

Summary

QoS routing is a key network function for the transmission and distribution of digitized audio/video across future high-speed networks. It has two objectives: finding routes that satisfy the QoS constraints and making efficient use of network resources. Based on the way the state information is maintained, the existing unicast/multicast routing algorithms can be divided into three classes: source routing, distributed routing, and hierarchical routing algorithms. Source routing algorithms are most thoroughly investigated. They simplify the path selection problem by locally computing a feasible path based on a global state that is maintained at every node. The responsibility of path selection is shared by intermediate nodes in the distributed routing. Most existing distributed routing algorithms also require the maintenance of global state. Limited work has been done on the hierarchical routing, especially for the NP-complete routing problems.

The polynomial-complexity routing problems (Fig. 3) were well solved by the shortest-path-based algorithms. Heuristics were proposed for the NP-complete routing problems with close to optimal results. However, problems still remain. Most source and distributed algorithms do not scale well due to the need to maintain global state. It is difficult to keep the global state up-to-date in large networks with dynamic data traffic. In addition, the source heuristic algorithms often have prohibitively high time complexities, which limits their practical value. Hierarchical routing provides a scalable solution because path selection is based on the aggregated state information whose size is much reduced. However, information imprecision is an issue of concern due to state aggregation. Furthermore, it is an unsolved problem to aggregate the state of a subnet with multiple QoS metrics. Future research should focus on efficient heuristic algorithms for the NP-complete routing problems, state aggregation with multiple QoS metrics, hierarchical routing with imprecise information, multipath routing, integration of QoS routing and best-effort routing, rerouting of dynamic traffic load, and efficient routing algorithms based on specific network models such as the rate-based scheduling network.

References

- [1] W. C. Lee, M. G. Hluchiy, and P. A. Humblet, "Routing Subject to Quality of Service Constraints Integrated Communication Networks," *IEEE Network*, July/Aug. 1995.
- [2] M. Garey and D. Johnson, *Computers and Intractability: A Guide to the Theory of NP-Completeness*, New York: W. H. Freeman and Co., 1979.
- [3] J. Moy, "OSPF Version 2," Network Working Group Internet Draft, Nov. 1992.
- [4] ISO/IEC/JTC1/SC6 WG2 N323, "Intra-Domain IS-IS Routing Protocol," Sept. 1989.
- [5] M. Steenstrup, "Inter-Domain Policy Routing Protocol Specification: Version 1," Internet Draft, May 1992.
- [6] C. Hedrick, "Routing Information Protocol," RFC 1058, Network Info. Center, SRI Int'l., Menlo Park, CA, June 1988.
- [7] K. Lougheed and Y. Rekhter, "Border Gateway Protocol 3 (BGP-3)," RFC 1267, SRI Int'l., Menlo Park, CA, Oct. 1991.
- [8] ISO/IEC/JTC1/SC6 CD10747, "Protocol for Exchange of Inter-domain Routing Information among Intermediate Systems to Support Forwarding of ISO 8473 PDUs."
- [9] ATM Forum, *Private Network Network Interface (PNNI) v1.0 Specifications*, May 1996.
- [10] R. Guerin and A. Orda, "QoS-based Routing in Networks with Inaccurate Information: Theory and Algorithms," *IEEE INFOCOM '97*, Japan, Apr. 1997.
- [11] B. Awerbuch et al., "Routing Through Teranode Networks with Topology Aggregation," *ISCC '98*, Athens, Greece, June 1998.
- [12] Z. Wang and J. Crowcroft, "QoS Routing for Supporting Resource Reservation," *IEEE JSAC*, Sept. 1996.
- [13] E. Dijkstra, "A Note on Two Problems in Connexion with Graphs," *Numerische Mathematik*, vol. 1, 1959, pp. 269-71.
- [14] R. E. Bellman, *Dynamic Programming*, Princeton University Press, Princeton, NJ, 1957.
- [15] S. Chen and K. Nahrstedt, "On Finding Multi-Constrained Paths," *IEEE ICC '98*, June 1998.
- [16] Q. Ma and P. Steenkiste, "Quality-of-Service Routing with Performance Guarantees," *Proc. 4th Int'l. IFIP Wksp. QoS*, May 1997.
- [17] H. F. Salama, D. S. Reeves, and Y. Viniotis, "Evaluation of Multicast Routing Algorithms for Real-Time Communication on High-Speed Networks," *IEEE JSAC*, vol. 15, no. 3, Apr. 1997, pp. 332-45.
- [18] G. N. Rouskas and I. Baldine, "Multicast Routing with End-to-End Delay and Delay Variation Constraints," *IEEE JSAC*, vol. 15, Apr. 1997, pp. 346-56.
- [19] D. Ferrari and D. C. Verma, "A Scheme for Real-Time Channel Establishment in Wide-Area Networks," *IEEE JSAC*, vol. 8, no. 3, Apr. 1990, pp. 368-79.
- [20] L. Zhang et al., "RSVP: A New Resource ReSerVation Protocol," *IEEE Network*, Sept. 1993.
- [21] I. Cidon, R. Rom and Y. Shavitt, "Multi-Path Routing Combined with Resource Reservation," *IEEE INFOCOM '97*, Japan, Apr. 1997, pp. 92-100.
- [22] A. Shaikh, J. Rexford, and K. Shin, "Dynamics of quality-of-service routing with inaccurate link-state information," U. of MI Tech. rep. CSE-TR-350-97, Nov. 1997.
- [23] H. F. Salama, D. S. Reeves, and Y. Viniotis, "A Distributed Algorithm for Delay-Constrained Unicast Routing," *IEEE INFOCOM '97*, Japan, Apr. 1997.
- [24] Q. Sun and H. Langendorfer, "A New Distributed Routing Algorithm with End-to-End Delay Guarantee," Unpublished paper, 1997.
- [25] S. Chen and K. Nahrstedt, "Distributed Quality-of-Service Routing in High-Speed Networks Based on Selective Probing," Tech. rep., Univ. of IL at Urbana-Champaign, Dept. Comp. Sci., 1998; the extended abstract was accepted by LCN '98.
- [26] K. G. Shin and C.-C. Chou, "A Distributed Route-Selection Scheme for Establishing Real-Time Channel," *6th IFIP Int'l. Conf. High Perf. Networking*, Sept. 1995, pp. 319-29.
- [27] J. Behrens and J. J. Garcia-Luna-Aceves, "Hierarchical Routing Using Link Vectors," *IEEE INFOCOM '98*, Mar. 1998.
- [28] J. Bennett and H. Zhang, "Hierarchical Packet Fair Queueing Algorithms," *ACM SIGCOMM '96*, Aug. 1996.
- [29] A. Demers, S. Keshav and S. Shenker, "Analysis and Simulation of a Fair Queuing Algorithm," *ACM SIGCOMM '89*, 1989, pp. 3-12.
- [30] S. Golestani, "A Self-Clocked Fair Queuing Scheme for Broadband Applications," *IEEE INFOCOM '94*, June 1994, pp. 636-46.
- [31] L. Zhang, "Virtual Clock: A New Traffic Control Algorithm for Packet Switching Networks," *ACM SIGCOMM '90*, Sept. 1990, pp. 19-29.
- [32] A. Orda, "Routing with End to End QoS Guarantees in Broadband Networks," *IEEE INFOCOM '98*, Mar. 1998.
- [33] D. H. Lorenz and A. Orda, "QoS Routing in Networks with Uncertain Parameters," *IEEE INFOCOM '98*, Mar. 1998.
- [34] J. Moy, "Multicast Extension to OSPF," Internet Draft, Sept. 1992.
- [35] L. Kou, G. Markowsky, and L. Berman, "A Fast Algorithm for Steiner Tree," *Acta Informatica* 15, 1981, pp. 141-45.
- [36] H. Takahashi and A. Matsuyama, "An Approximate Solution for the Steiner Tree Problem in Graphs," *Mathematica Japonica*, 1980.
- [37] V. P. Kompella, J. C. Pasquale, and G. C. Polyzos, "Multicast Routing for Multimedia Communication," *IEEE/ACM Trans. Networking*, June 1993.
- [38] Q. Sun and H. Langendorfer, "A New Distributed Routing Algorithm with End-to-End Delay Guarantee," *2nd Wksp. Protocols Multimedia Sys.*, Oct. 1995, pp. 452-58.
- [39] R. Widjono, "The Design and Evaluation of Routing Algorithms for Real-Time Channels," Technical rep., ICSI TR-94-024, Univ. CA at Berkeley Int'l. Comp. Sci. Inst., June 1994.
- [40] Q. Zhu, M. Parsa, and J. J. Garcia-Luna-Aceves, "A Source-Based Algorithm for Delay-Constrained Minimum-Cost Multicasting," *IEEE INFOCOM '95*, Boston, MA, Apr. 1995.
- [41] V. P. Kompella, J. C. Pasquale and G. C. Polyzos, "Two Distributed Algorithms for Multicasting Multimedia Information," *ICCCN '93*, San Diego, CA, June 1993, pp. 343-49.
- [42] B. Awerbuch et al., "Throughput-Competitive On-line Routing," *34th Annual Symp. Foundations of Comp. Sci.*, Palo Alto, CA, Nov. 1993.
- [43] B. Awerbuch et al., "Competitive Routing of Virtual Circuits with Unknown Duration," *5th ACM-SIAM Symp. Discrete Algorithms*, Jan. 1994.
- [44] S. Plotkin, "Competitive Routing of Virtual Circuits in ATM networks," *IEEE JSAC*, vol. 13, Aug. 1995, pp. 1128-36.
- [45] J.-Y. Le Boudec and T. Przygienda, "A Route Pre-Computation Algorithm for Integrated Services Networks," *J. Network and Sys. Mgmt.*, vol. 3, no. 4, 1995, pp. 427-49.
- [46] A. Iwata et al., "PNNI Routing Algorithms for Multimedia ATM Internet," *NEC R&D*, vol. 38, 1997.
- [47] M. Peyravian and A. D. Kshemkalyani, "Network Path Caching: Issues, Algorithms and A Simulation Study," *Comp. Commun. Rev.*, vol. 20, 1997, pp. 605-14.
- [48] A. Shaikh, J. Rexford, and K. Shin, "Efficient precomputation of quality-of-service routes," *Wksp. Network and Op. Sys. Support for Digital Audio and Video*, July 1998.
- [49] D. D. Kandlur, K. G. Shin and D. Ferrari, "Real-Time Communication in Multi-hop Networks," *Proc. 11th Int'l. Conf. Dist. Comp. Syst.*, 1991, pp. 300-7.
- [50] C. Hou, "Routing Virtual Circuits with Timing Requirements in Virtual Path Based ATM Networks," *IEEE INFOCOM '96*, 1996.
- [51] S. Chen and K. Nahrstedt, "Distributed QoS Routing with Imprecise State Information," *ICCCN '98*, Oct. 1998.
- [52] S. Deering and D. Cheriton, "Multicast Routing in Datagram Internetworks and Extended LANs," *ACM Trans. Comp. Sys.*, May 1990, pp. 85-111.
- [53] F. K. Hwang and D. S. Richards, "Steiner Tree Problems," *Networks*, vol. 22, 1992, pp. 55-89.
- [54] R. Prim, "Shortest Connection Networks and Some Generalizations," *Bell Systems Tech. J.*, vol. 36, 1957, pp. 1389-1401.
- [55] K. Carlberg and J. Crowcroft, "Building Shared Trees Using a One-to-Many Joining Mechanism," *Comp. Commun. Rev.*, Jan. 1997, pp. 5-11.
- [56] M. Faloutsos, A. Banerjee, and R. Pankaj, "QoSMIC: Quality of Service sensitive Multicast Internet protoCol," *ACM SIGCOMM '98*, Sept. 1998.
- [57] A. Banerjee, "Simulation Study of the Capacity Effects of Dispersy Routing for Fault Tolerant Realtime Channels," *ACM SIGCOMM '96*, Aug. 1996.
- [58] N. F. Maxemchuk, "Dispersy Routing," *IEEE ICC '75*, San Francisco, CA, June 1975.
- [59] N. S. V. Rao and S. G. Batsell, "QoS Routing Via Multiple Paths Using Bandwidth Reservation," *IEEE INFOCOM '98*, San Francisco, CA, Mar. 1998.
- [60] C. Parris, H. Zhang, and D. Ferrari, "A Mechanism for Dynamic Re-routing of Real-time Channels," *Tech. Rep. TR-92-053*, Int'l. Comp. Sci. Inst., Berkeley, CA, Apr. 1992.
- [61] Q. Ma, "Quality of Service Routing in Integrated Services Networks," Ph.D. thesis, Tech. rep. CMU-CS-98-138. Carnegie Mellon Univ., Dept. Comp. Sci., 1998.
- [62] Q. Ma and P. Steenkiste, "Supporting Dynamic Inter-Class Resource Sharing: A Multi-Class QoS Routing Algorithm," *IEEE INFOCOM '99*, New York, NY, Mar. 1999.

Biographies

SHIGANG CHEN (s-chen5@cs.uiuc.edu) is a Ph.D. candidate in the Department of Computer Science, University of Illinois at Urbana-Champaign. His research interests include quality of service, distributed computing, multimedia, and resource management. He received his B.A. in computer science from the University of Science and Technology of China, P.R.China, in 1993. He received his M.S. in computer science from the University of Illinois at Urbana-Champaign in 1997.

KIARA NAHRSTEDT (klara@cs.uiuc.edu) is an assistant professor at the University of Illinois at Urbana-Champaign, Computer Science Department. Her research interests are directed toward multimedia services, protocols, endpoint architectures for multimedia, quality of service provision and resource management in real-time multimedia distributed systems. She is co-author of the widely used multimedia book *Multimedia: Computing, Communications and Applications* (Prentice Hall), and the recipient of the Early NSF Career Grant. She received her B.A. in mathematics from Humboldt University, Berlin, in 1984, and her M.Sc. in numerical analysis from the same university in 1985. She was a research scientist in the Institut für Informatik in Berlin working on the X.400 mailing system from 1987 to 1990. In 1995 she received her Ph.D. from the University of Pennsylvania in the Department of Computer and Information Science. She is a member of the ACM and IEEE Computer Society.