

Delay Sensitive Routing for High Speed Packet-Switching Networks*

Yao-Nan Lien and Yu-Sheng Huang, National Chengchi University, Taiwan, lien@cs.nccu.edu.tw

Abstract—In a high speed packet switching network, the processing time in a node becomes an important delay component that must be taken into account in designing routing algorithms for time sensitive services. To illustrate this hypothesis, we propose a new delay sensitive routing algorithm that takes both link and node delay time into consideration. Our simulation study shows that our algorithm could easily outperform traditional routing algorithms that only consider link delay time.

I. INTRODUCTION

Providing time sensitive services becomes an essential task for some packet-switching networks such as All-IP networks [1], which have to carry all types of traffics currently supported by both circuit-switching and packet-switching networks. Since routing is a critical task to select path to deliver a packet in a packet-switching network, such a network requires a delay sensitive routing mechanism to provide time sensitive services with QoS guarantee. However, most traditional routing algorithms do not take delay time as a major concern. Only a few are designed for time sensitive services [7,9,10,11]. These time sensitive routing algorithms were designed at the time when networks were slow and link bandwidth was the scarcest resource. As link bandwidth grows rapidly in recent years due to the advance of optical communication technologies, link bandwidth is no longer the only scarce resource. The processing time in a node, e.g. router, becomes another critical source of time delay. It must be taken into account in designing adequate routing algorithms for time sensitive services.

In this paper, we designed a new flow-based routing algorithm, KLONE, that takes average delay time as its minimization objective and both node and link as delay components. Through an intensive evaluation using simulation method, we demonstrate that the KLONE algorithm outperforms the traditional OSPF algorithm [6] by about 30%.

A. Path and Link Delay Time

The delay time of a packet traveled along a path, referred to as *path delay time* in this paper, can be divided into three components: link delay, node delay, and end-host delay. Among them, end-host delay, which is the delay occurred at the hosts of both ends, can be ignored in the design of routing algorithms since it is independent of the path it travels.

Link delay time contains three components, the queuing delay, the propagation delay and the transmission delay. Transmission delay is the time of a data unit being transmitted along a specific link with the propagation delay time ignored. For example, it will take 82ms to transmit 128k bits data along a T1 link.

B. Node Delay Time

Node delay is the delay time occurs in a node (e.g., a router). It contains two components: the processing delay and the queuing delay for a processor. The tasks performed in a router and the processing power of the router determine its delay time.

As link bandwidth grows rapidly and is getting closer to node processing capacity, node delay is increasing its share in a path delay time making itself a notable component in network performance.

C. The Myth of Bandwidth

In the beginning of router algorithm development, the tasks performed by a router are simple and the power of the processor within a router is much faster than the links in terms of processing or transmitting packets. The link delay was the major concern in designing a delay sensitive routing algorithm. In recently years, network operators start to deploy fiber optic networks with DWDM technique making a dramatic increase in network bandwidth. This fast growth in bandwidth makes link bandwidth closer to the node processing capacity and results in the increase of relative weigh of node delay. One way to reduce the delay time cause by routers is to embed higher layer protocols to lower layer equipments such as Layer 3 switch or MPLS [2]. Another way is to choose better routing algorithms that take node delay time into consideration.

D. An Illustration Example

In the following example, discovering a minimal delay path with and without taking node delay into account will be compared. As shown in Figure 1(a), a simple network is composed of eight nodes, A, B, C, ..., H, and eight directed links. The weight of link AD is 2, and all the other links is 1. We assume the delay time caused by a node is proportional to the number of traffic flows passing that node. There is a unit traffic demand from A to F, from B to G and from C to H, respectively. Without considering node delay, the best routing algorithm will route all three traffic flows through node E, as shown in Figure 1(b). The delay time of each path will be 6.

* This research is sponsored by the NSC grant NSC 91-2219-E-004-002.

Figure 1(c) shows the result of another possible routing that takes node delay time into account. One traffic flow will pass node D, instead of node E. The delay time is then 5 for each flow.

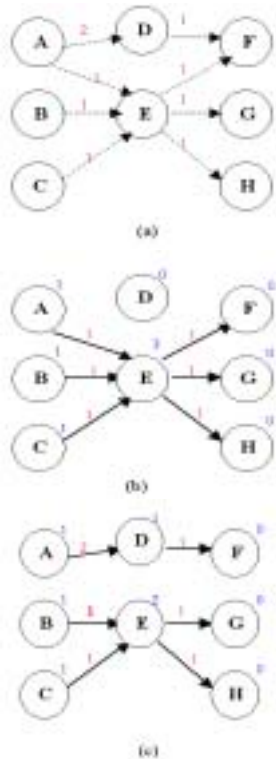


Fig. 1. Example routing with node delay considered.

Above example shows that, for high-speed networks, a routing algorithm that takes node delay time into account may obtain a better result as compared to the traditional routing algorithms that consider link delay time only.

II. RELATED WORK

A. Current Routing Approaches

Shortest path routing is to build up a shortest path tree to present the network topology, so then routes each of request traffics to its destination. The Dijkstra's shortest path algorithm[4,8] is a very famous example of shortest path routing, and it discovers a node's shortest paths to other nodes in $O(n \log n)$. Another routing algorithm is *flooding* [8], in which every incoming packet is sent out on every outgoing line except the one it arrived on. It generates vast numbers of duplicate packets. It is not practical in most applications, but it does have some uses. *Flow-based routing* considers both the network topology and load [8]. Link delay time is, instead of a given constant, a load dependent variable. *Distance vector routing* estimates the distance from source to destination by certain approaches [8,11]. They are often referred to as "Bellman-Ford" protocols because they are based on a shortest path computation algorithm. Distance vector routing protocols periodically send information to its neighbor nodes. Each node

could estimate the distance to other nodes according to the number of intermediate passed nodes. *Link state routing* focuses on the states of links [11]. Link state routing transforms the link states into some mathematical expressions to choose proper paths. Open Shortest Path First (OSPF[6]) is a widely used example of link state routing, and is parallelly used with RIP, which is a distance vector routing. *QoS routing* selects routes based on flow QoS requirements and network resource availability. QoS routing determines feasible paths satisfying QoS requirements, while optimizing resource usage and degrading gracefully during periods of heavy load [9,11].

B. Our Approach

The networks that are to support time sensitive services could choose to use a proper QoS routing. However, existing routing algorithms only consider link delays such that they may not be adequate for high speed networks as explained in Section I. In Section III, we model the problem as a flow-based routing problem with link and node delay dependent on their loads. An iterative approach is taken to cope with the difficulty of estimating load dependent delay time on links and nodes. A transformation is then taken to convert node delay into link delay such that the intermediate problem of each iteration can be solved using a traditional routing algorithm. The proposed algorithm is evaluated by comparing with the OSPF algorithm in Section IV using average path delay time and goodput ratio as performance metrics.

III. ROUTING WITH NODE DELAY

A. Routing Problem Model

Given a directed graph $G(V, E)$, with $|V|$ nodes and $|E|$ links, the propagation delay time and bandwidth of each link, and the processing capacity of each node, the problem is to find a set of paths for a given set of traffic demands and the delay bound such that the total delay time is minimized. Given and derived parameters are listed in Table 1 and 2 respectively.

TABLE I
NOTATION OF INPUT PARAMETERS.

$G(V, E)$	a directed graph, G , with set of nodes V and set of directed link E
v_i	a node; $v_i \in V$
e_k	a directed link $e_k = (v_x, v_y) \in E$, v_x is the start node, v_y is the end node of link e_k ; also denoted as e_{xy}
ij	volume of traffic requests from v_i to v_j
k	volume of traffic requests from v_k to all other nodes. $k = \{k_i i=1, \dots, V \}$
D	allowable delay time to transmit a packet from source to destination
$b(e_k)$	bandwidth of link e_k
$t(e_k)$	propagation delay time of link e_k
$p(v_k)$	processing capacity of node v_k

TABLE 2
NOTATION OF DERIVED PARAMETERS AND ROUTING RESULTS.

M_k	$\sum_1^{ \mathcal{K} } \lambda_{ki}$, total volume of requested traffics starting from node v_k ; the traffic volume of $ \mathcal{K} $
$S_k^{(n)}$	the selected routing path set(slice) of iteration n , corresponding to the request k
$P^{(n)}$	the selected routing path set(pasta) of iteration n , set of $S_k^{(n)}$
μ_h	volume of traffics passing through link e_h , starting from v_x , and ending at v_y , also denoted as μ_{xy}
μ_k	volume of traffic passing through node v_k
U	set of μ_h , $U = \{ \mu_h \}$
ϕ_{ij}	the selected path for ij , by the routing algorithm; $\phi_{ij} = v_i, e_{i+1}, v_{i+1}, e_{i+2}, \dots, e_{j-1}, v_j$
Φ	set of ϕ_{ij}
$d(v_k)$	delay time caused by node v_k
$d(e_h)$	delay time caused by link e_h
$d(\phi_{ij})$	total delay time along path set ij , $d(\phi_{ij}) = \sum_{e \in \phi_{ij}} d(e) + \sum_{v \in \phi_{ij}} d(v)$

The problem is then formulated as follows:

Find Φ

$$\begin{aligned} & \ni \min \sum_{\phi_{ij} \in \Phi} d(\phi_{ij}), \\ & \text{s.t. } d(\phi_{ij}) < D, \forall \phi_{ij} \in \Phi. \end{aligned} \quad (1)$$

The delay bound of each traffic flow, D , could be variant without incurring a significant impact to the model. $d(\phi_{ij})$ is the total delay time for a traffic flow; it is an accumulation of the delay time on all links and nodes along all the selected paths. The delay time occurred in the components of a real network is actually dependent on the stochastic behavior of the traffic and routing process. In reality, it is extremely difficult to solve a routing problem that takes stochastic behavior into account. Therefore, we take a compromised approach by relaxing the stochastic property of traffic and routing process in estimating of delay time on links and nodes.

We assume all traffics are of Constant Bit Rate (CBR) type and all resources (processing and link bandwidth) are proportionally shared by all traffic flows passing through. In this way, the load of each resource can be computed based on the total amount of traffics passing through that resource. Although this is a compromised model, it is more realistic as compared to the traditional fixed weight model.

The problem can be easily proved NP-hard by reducing into a 0-1 knapsack problem. Furthermore, both objective and constraints are not simple functions of given parameters (delay time). Instead, they are result dependent variables. This makes the problem much more complicated.

B. Iterative Solution Approach

An iterative approach is taken to cope with the difficulty of load dependent delay time on nodes and links. In the first

iteration, the delay time of every node is set to zero and the delay time of every link is set to its propagation delay time. In other iterations, the delay time of nodes and links are computed based on the result obtained in the previous iteration.

For convenience, the result obtained in an iteration is called a *pasta*. In each iteration, the problem is further divided into some number of sub-problems. Each routing solution (pasta) can be decomposed into a number of single root flow trees, named a *slice*. In such a tree, the root node is any node and the tree presents the flows generated from that root node and are forwarded to all other nodes. In each incremental step within an iteration, a slice corresponding to a request set $|\mathcal{K}|$ is extracted from the pasta, recomputed using an algorithm similar to Dijkstra's, and superimposed back to the pasta. Thus, a pasta, the result of an iteration, is recomputed incrementally slice by slice.

After some number of iterations, hopefully, the delay time of each network component will be stabilized, and the routing solution obtained will be a stable solution. Termination is triggered by two conditions: when average path delay of two consecutive iterations is smaller than the predetermined value ϵ ; or the number of iterations exceeds a given number. In the first condition, ϵ is defined as the difference in total delay time of two consecutive iterations divided by the total path delay time of previous iteration as shown in Eq. 2:

$$\epsilon = \left| \frac{\sum d(\phi_{ij}) - \sum d(\phi_{im})}{\sum d(\phi_{ij})} \right| \Big|_{\phi_{ij} \in S_k^{(n)}, \phi_{im} \in S_{k+1}^{(n)}} \quad (2)$$

We denote the result (pasta) obtained in the n -th iteration as $P^{(n)}$, the single root path tree (slice) corresponding to the k in the n -th iteration as $S_k^{(n)}$, and $P^{(n)} = \{S_1^{(n)} \oplus S_2^{(n)} \dots \oplus S_k^{(n)}\}$, where \oplus denotes a superposition. The iterative procedure is summarized in the followings:

(I) Initial condition

for all nodes and links, set $\mu = 0$, $\mu = 0$, $d(v) = 0$, $d(e) = t(e)$;

$$S_1^{(0)} = S_2^{(0)} = S_3^{(0)}, \dots, = S_{|\mathcal{V}|}^{(0)} = \{\}; \quad // \text{empty set}$$

$$P^{(0)} = S_1^{(0)} \oplus S_2^{(0)} \oplus \dots \oplus S_{|\mathcal{V}|}^{(0)};$$

// denotes superimposing a slice into a pasta

// denotes removing a slice from a pasta

(II) First iteration

$$P^{(1)} = \{\};$$

route 1 based on $(P^{(0)}, S_1^{(0)})$, to obtain $S_1^{(1)}$;

$$P^{(1)} = P^{(1)} \oplus S_1^{(1)};$$

route 2 based on $(P^{(0)}, S_1^{(0)}, S_2^{(0)}, S_1^{(1)})$, to obtain

$$S_2^{(1)};$$

$$P^{(1)} = P^{(1)} \quad S_2^{(1)};$$

route $|V|$ based on $(P^{(0)} \quad S_1^{(0)} \quad S_{|V|}^{(0)} \quad S_{|V|-1}^{(1)})$,
to obtain $S_{|V|}^{(0)}$; $P^{(1)} = P^{(1)} \quad S_{|V|}^{(1)}$;

(III) On the k-th iteration:

$$S^{(k)} = \{\};$$

for $j \leftarrow 1$ to $|V|$ {

route j based on $(P^{(k-1)} \quad S_2^{(k-1)} \quad S_j^{(k-1)} \quad S_{j-1}^{(k)})$, to obtain $S_j^{(k)}$; $P^{(k)} = P^{(k)} \quad S_j^{(k)}$;

$$S_{j-1}^{(k)});$$

}

(IV) Termination Conditions

When $P^{(M)} \approx P^{(M+1)}$ or the number of iteration exceeds a given number, terminate;

C. Estimation of Path Delay Time

The delay time of path ϕ_{ij} , $d(\phi_{ij})$, consists of the delay time on all nodes and links in a path, which is $d(v_i) + d(e_{i,i+1}) + d(v_{i+1}) + d(e_{i+1,i+2}) + \dots + d(e_{j-1,j}) + d(v_j)$. μ_h and σ_k are defined as the total volume of traffic flows passing through a link e_h and a node v_k , respectively and can be computed as follows:

$$\mu_h = \sum_{\substack{\phi_{ij} \in \Phi \\ e_h \in \phi_{ij}}} \lambda_{ij}, \text{ and} \quad (3)$$

$$\sigma_k = \sum_{\substack{\phi_{ij} \in \Phi \\ v_k \in \phi_{ij} \\ v_k \neq v_j}} \lambda_{ij}. \quad (4)$$

1) Link Delay Time

$d(e_h)$ is the delay time of a flow of packets passing through link e_h . μ_h is the total flows passing through e_h . As mentioned in Section III.A, we assume every traffic flow is a CBR and the bandwidth of a link is shared by all the traffic flows passing through that link. The queuing delay on the link, thus, can be ignored. Therefore, the delay time on a link for a flow is approximately the propagation delay time plus the total traffic flows divided by the bandwidth of that link, as shown in Eq. 5,

$$d(e_h) = \mu_h / b(e_h) + t(e_h) = \left(\sum_{\substack{\phi_{ij} \in \Phi \\ e_h \in \phi_{ij}}} \lambda_{ij} \right) / b(e_h) + t(e_h) \quad (5)$$

Notice that the delay time of a link, $d(e_h)$, is independent of the size of the flow passing that link. All traffic flows passing the same link are delayed by the same amount of time.

2) Node Delay Time

$d(v_k)$ is the delay time caused by a node, v_k . Again, to simplify the delay time model, we assume all traffic flows

passing a node are processed in time-sharing fashion, such that the $d(v_k)$ can be estimated as the total volume of traffics divided by the processing capacity of that node, as shown in Eq. 6:

$$d(v_k) = \sigma_k / p(v_k) = \left(\sum_{\substack{\phi_{ij} \in \Phi \\ v_k \in \phi_{ij} \\ v_k \neq v_j}} \lambda_{ij} \right) / p(v_k). \quad (6)$$

3) Path Delay Time

Thus, the delay time of a path ϕ_{ij} is

$$d(\phi_{ij}) = \sum_{e_h \in \phi_{ij}} d(e_h) + \sum_{v_k \in \phi_{ij}} d(v_k). \quad (7)$$

4) Node Delay to Link Delay Conversion

We need an efficient algorithm to solve a single-source shortest path routing problem to obtain a *slice*. Unfortunately, current shortest path algorithms all assume zero weight on nodes such that they are not adequate to solve this problem even though the delay time of network components are all constant within each iteration. There are two approaches to solve this problem: either to develop a new algorithm that considers both node and link delay together or to convert node delays into link delays, and then apply a conventional shortest path algorithm to solve this problem. We choose the second approach for simplicity. In the transformation, node delay time is added to the propagation delay time of each incoming link to the node. By doing so, we obtain another graph that has weights on its links only and is equivalent to the original graph with respect to the path delay time, as depicted in Fig. 2.

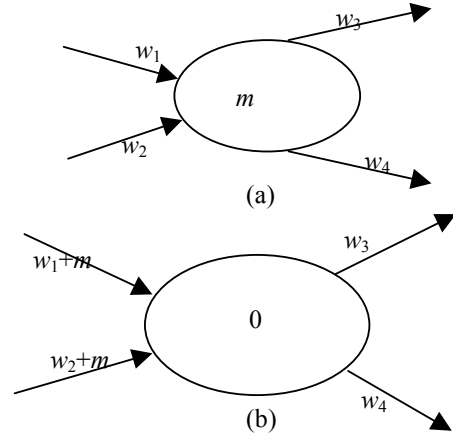


Fig. 2. Transformation of node delay to link delay.

IV. PERFORMANCE EVALUATION

It can be easily proved that the complexity of the KLONE algorithm is $N^4 \log N$ in the worst case. It is further evaluated by comparing with the traditional OSPF routing algorithm using a numerical simulation. Performance metrics are convergence speed, average path delay time, and goodput ratio. Convergence speed is evaluated by two different values: the

number of iterations/slices when the convergence occurs (the average path delay time of two consecutive pastas differ by a predefined value, ϵ) divided by the total number of nodes (K_1/N and K_2/N). Goodput ratio is the ratio of total satisfied traffic requests to the total traffic requests. Average path delay time is the average time for a unit of request traffic passing through the network. It is computed as summation of the size of a traffic multiplied by the delay time on the selected path and then divided by the size of total traffic,

$$\frac{\sum [|\lambda_{ij}| * d(\phi_{ij})]}{\sum |\lambda_{ij}|} \quad (8)$$

A. Design of Experiments

We compared KLONE algorithm and OSPF algorithm in 64,000 different test instances, in the combinations of different number of nodes, network connectivity ratio, and different link bandwidth/processing capacity ratio. The range of link bandwidth was set from 0 to 400 Gbps, and propagation delays stayed below 20 ms. The number of nodes was set from 10 to 100 with a processing capacity in the range of Gbps.

Connectivity is defined as $\frac{P}{N * (N - 1)}$, where P is the number of links, and N is the number of nodes. The range of connectivity was set from 0 to 100 percents. The *BP ratio* is defined as $b(e)/p(v)$, where $b(e)$ is the link bandwidth and $p(v)$ is the node processing capacity. We varied it from 1/300 to 1/1. The traffic coming into an edge node is assumed in an aggregated form. For a graph of N nodes, there are $N*(N-1)$ requests, one from each node to every other node. The upper bound of delay time of all paths is set to D and D varies from 100 to 2000 ms.

B. Experiments and Results

The experiments and results are presented in this section. Most of the figures shown in this section are for the networks of size 50.

1) Exp-1: Convergence Test

We adjusted the following three parameters in the experiment to study their impact to the convergence speed: the BP ratio, the number of nodes and the connectivity. The results show that neither BP ratio nor the number of nodes has any impact to the convergence speed. On the other hand, we found that the convergence speed is dependent on the connectivity. This may be caused by two different reasons. First, higher connectivity may make a request easier to find a very good satisfied path, and then there is a higher probability to choose the same path in the succeeding iteration. On the other hand, lower connectivity may make a request having fewer paths to choose, so that the solution domain is much smaller and thus the convergence speed is faster. In more than 90% of the test instances, we found that the lines of both average path delay

time and goodput ratio become stable after the $K_1=2/N$ and $K_2=2$.

2) Exp-2: Sensitivity to Connectivity

Intuitively, higher connectivity implies more available paths between nodes. We studied the dependency between the connectivity and the two performance metrics: average path delay time and goodput ratio. We varied connectivity from 0% to 100% to see how average path delay time and goodput ratio are influenced.

In this experiment, we found that, at the same number of nodes, the average path delay time becomes smaller as the connectivity increases, as shown in Figure 3(a). The average path delay time improvement is defined as $(T_2-T_1)/T_2$, where T_1 and T_2 are the average path delay time of KLONE algorithm and OSPF algorithm, respectively. The larger the value, the better the KLONE algorithm. In Figure 3(b), we show that at higher connectivity, KLONE algorithm has a higher goodput ratio than OSPF algorithm.

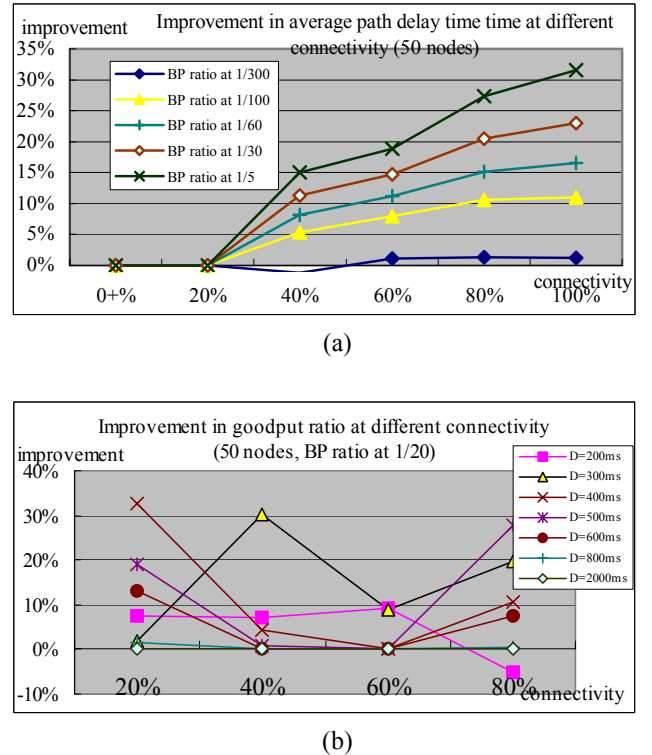
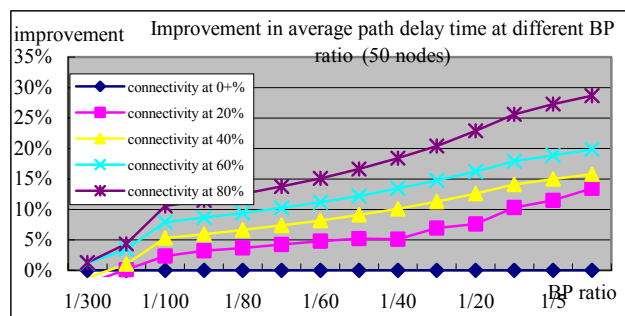


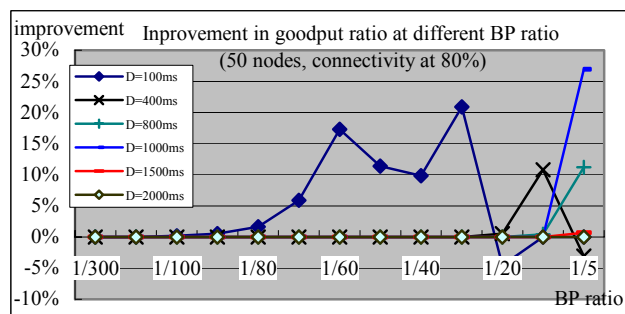
Fig. 3. Sensitivity to connectivity.

3) Exp-3: Sensitivity to BP Ratio

We varied the BP ratio from 1/300 to 1 to see the dependency between the BP ratio and the two performance metrics. We found that when the BP ratio increases, the improvement of average path delay time increases, as shown in Figure 4(a). This is consistent with our hypothesis that when the speed of links increases, an algorithm that concerns both link and node delay times might have a better performance than OSPF, which only concerns links delay times. We also compared the goodput ratio of KLONE algorithm and OSPF algorithm. We found that at different BP ratios, goodput ratio of KLONE algorithm is usually better than OSPF algorithm, as shown in Figure 4(b).



(a)



(b)

Fig. 4. Sensitivity to BP ratio.

4) Exp-4: Sensitivity to Number of Nodes

This experiment studied the dependency between the number of nodes and the delay time improvement. The number of nodes was varied from 20 to 70. The improvement increases as the number of nodes increases. Increasing the number of nodes will also increase the goodput ratio at the same delay bound, D . At different number of nodes, KLONE algorithm has a better goodput ratio than OSPF algorithm. Due to space limit, detailed results are not shown here.

5) Comparison with Optimal Solution

In order to estimate the absolute performance of KLONE algorithm, we compared both algorithms with the optimal solution in a very small scale test instance in which the number of nodes was set to 10, connectivity was set to 20%, BP ratio was set at 1/10. The average path delay time of KLONE and OSPF algorithm is approximately 30% and 60% higher than the optimal solution, respectively. While the goodput ratio is approximately 20% and 40% lower than the optimal solution, This toy-type study may not mean too much. However, it still gives us a sense of the discrepancy between the KLONE algorithm and the optimal solution.

V. CONCLUDING REMARK

With an intensive evaluation, we demonstrate the importance of the nodes delay in a time sensitive path in high-speed packet-switching networks. We hypothesized that a routing algorithm that considers the delay time of both nodes and links may have a better performance than that only considers link delay time in supporting delay sensitive services. We developed a flow-based routing algorithm, KLONE algorithm,

which considers both link and node delay time. The results of the evaluation show that KLONE algorithm could have a better performance than OSPF algorithm in most cases, with only a few exceptions. The hypothesis that considering node delay is important in high-speed packet-switching network is thus demonstrated.

This algorithm still has some weak points. First, KLONE algorithm may have worse goodput than OSPF algorithm when the delay bound is very low. Secondly, it does not support multi-path routing for the same traffic stream yet. Furthermore, a distributed version is needed in order to apply it onto real networks. In estimating the delay time of nodes and links, the traffic model should also be more realistic to include different traffic types, such as VBR, and in difference priorities.

REFERENCES

- [1] 3rd Generation Partnership Project, "Technical Specification Group Services and Systems Aspects; Architecture for an All IP network", 3GPP TR 23.922 version 1.0.0., October 1999.
- [2] D. Awduche, J. Malcolm, J. Agogbua, M. O'Dell and J. McManus, "Requirements for Traffic Engineering Over MPLS", RFC 2702, September 1999.
- [3] Christophe Beaujean, "Delay-Based Routing Issues in IP Networks", contact GRADIENT CR/98/148, May 2000.
- [4] Dijkstra, E.W., "A Note on Two Problems in Connection with Graphs", Numerische Math, vol. 1, March 1959, pp. 269-271.
- [5] C. Hedrick, "Routing Information Protocol", RFC 1058, June 1988.
- [6] J. Moy, "OSPF version 2", RFC 1583, March 1994.
- [7] Douglas S. Reeves and Hussein F. Salama, "A Distributed Algorithm for Delay-Constrained Unicast Routing", *IEEE Transaction on Network*, April 2000.
- [8] A. S. Tanenbaum, "Computer Networks, Third Edition", Prentice Hall, March 1996, pp. 345-366.
- [9] Z. Wang and J. Crowcroft, "Quality of Service Routing for Supporting Multimedia Applications", *IEEE J. on Selected Areas in Communications*, September 1996.
- [10] R. Wideyona, "The Design and Evaluation of Routing Algorithms for Real-Time Channels", *International Computer Science Institute, Univ. of California at Berkeley, Tech Rep. ISCI TR-94-024*, June 1994.
- [11] Ossama Younis and Sonia Fahmy, "Constraint-Based Routing in the Internet: Basic Principles and Recent Research", *IEEE Communications Society Surveys & Tutorials*, vol 5, no. 1, 3Q 2003.