

**ROUTING AND CAPACITY ASSIGNMENT IN A  
NETWORK WITH DIFFERENT CLASSES OF MESSAGES**

by

**Irina Neuman**

Information Systems Department  
Leonard N. Stern School of Business  
New York University  
90 Trinity Place  
New York, NY 10006

January 1989

Center for Research on Information Systems  
Information Systems Department  
Leonard N. Stern School of Business  
New York University

**Working Paper Series**

CRIS #198  
GBA #89-10

# Routing and Capacity Assignment in a Network with Different Classes of Messages

Irina Neuman

Information Systems Area  
Graduate School of Business Administration  
New York University, New York, NY

## Abstract

A mathematical model is presented for the problem of jointly assigning routes to the communicating pairs of nodes and capacities to the links in a packet switched network. It is assumed that several classes of flow are using the network, different service requirements and message characteristics being associated with each class. An algorithm that generates good feasible solutions to the model, together with tight lower bounds on the value of the objective function, is presented. Results of numerical experiments using several network topologies are reported.

## 1 Introduction

A starting point for most of the existing research in the area of backbone network design is the implicit assumption that all messages in the network have similar characteristics and requirements. As a result, a uniform treatment is adopted for all messages, with no distinction being made among different types of applications, each with their own specific characteristics, nor between different user requirements. Such an approach greatly reduces the complexity of the analysis, but in most cases the assumption does not correspond to the real world environment. Explicitly taking into account the characteristics and requirements of different classes of messages not only leads to a solution that is preferable from a global perspective (e.g. depending on the performance criterion, the average delay in the network may be reduced, or an overall less costly design may be achieved), but also the solution is better tailored to the individual user needs.

The practical relevance of the issue is suggested by the fact that routing strategies that differentiate among messages in accordance with their various characteristics and requirements are commonly implemented by many operational networks, SNA [2], DATAPAC [17], and SITA [4] being just some of the examples.

This paper addresses the following problem: how to *simultaneously* select the link capacities and the routes to be used by the communicating nodes in a network that supports several classes of messages with different priority levels. The model is a generalization of earlier models introduced in [16] and [6].

The capacity and flow assignment aspects of network design are usually dealt with separately in the literature. In most cases, such an approach is not appropriate. The close interaction that exists between the capacity value of a link, and the delay incurred by a given flow on that same link, makes it difficult to claim that a truly good solution has been found for either of the two problems when considered independently.

A growing body of research literature deals with the performance *analysis* of computer networks. Significant efforts have been made to tailor the general models for networks of priority queues to the specific characteristics of computer communication systems [3,11,15,17,18,19]. Even when simplified networks are considered, the complexity of the underlying phenomena is such that finding optimal or near optimal solutions to these models is a difficult task.

The comparative results in the area of performance evaluation of computer communication systems supporting several classes of service strongly suggest that the overall performance is significantly improved when messages are prioritized. These theoretical indications, together with the experience gained from the operational networks that chose to implement similar methods, are powerful arguments in favor of such schemes. Nevertheless, the literature dealing with the related *design* issues is very limited. To our knowledge, the only authors who incorporated this important aspect into their design methodology, are K. Maruyama and D. T. Tang. A sequence of their papers deals with increasingly complex aspects of the problem. In [14] only discrete link capacity assignment is considered. Messages are classified according to their processing and delay characteristics, and known priority levels are associated with each message class. The heuristic procedure suggested for the solution of the model attempts to minimize the total link cost, while satisfying the delay requirement constraints specific to each class. An interesting refinement is introduced in [13]. This time the priority levels are no longer assumed to be known in advance, i.e. they are no longer user-assigned, and are instead determined by the system. Thus, an additional reduction in the cost of the capacity assignment can be achieved, by determining the best mapping of  $n$  message classes into  $r$  different priority levels. The heuristic is a composite procedure, that alternates between two separate algorithms, for capacity and priority assignment respectively, until a local minimum is found. Finally, in [12] the scope of the proce-

cedure is further broadened, by also including an algorithm that handles static flow assignment. The global algorithm starts by determining the flow on each link, based on the maximum available capacities. The initial flow assignment satisfies the throughput requirement, but ignores the capacity constraints. Next, the procedure iterates between the capacity and priority assignment, and the flow assignment algorithms, until no further improvement is possible.

The complexity of the issues involved in network design renders an attempt to find optimal solutions an illusory goal for all but the most trivial cases (e.g. very small networks and/or models based on highly unrealistic assumptions). As a result, the majority of the solution methods suggested in the literature are of a heuristic nature. They do not provide for a way to evaluate the quality of the feasible solution generated, a fact which may significantly hamper their usefulness for real life applications.

The remainder of the paper is organized as follows: in section 2, the problem is defined, and a mathematical model is developed; section 3 describes the solution procedure, while methods for obtaining good upper and lower bounds on the optimal value are outlined in section 4; finally, section 5 contains the results and the analysis of the computational experiments conducted with the model, as well as some concluding remarks.

## 2 Problem Formulation

The limited capacity of network components gives rise to queuing phenomena. These are modeled by associating a server with each link, whose service rate is determined by the link capacity and by the message length. Messages are viewed as customers competing for the link server. Unlimited buffering space and no processing delays at the network nodes are assumed for ease of exposition. Propagation delays, which are negligible for terrestrial links, are also ignored. Messages from each class arrive on the network boundaries according to Poisson processes with known average interarrival times. Message lengths are exponentially distributed for each class. The independence assumption, first introduced in [9], is also used. The resulting model is that of a Jacksonian network of queues, in which average delay measures are easily computable.

Each message class is associated with a known priority level. A head-of-the-line non-preemptive discipline is imposed on the messages waiting at each link. Static routing is assumed in the model. Such routing mechanisms are used in many operational networks (e.g. [2], [4],[20]), and are known to perform well, mostly due to their simplicity and stability.

Two distinct types of costs, which reflect the unified way in which the model deals with the flow and capacity assignment issues, are considered:

1. *capacity costs*, comprised of a fixed setup cost, and a variable cost, which is a function of the traffic on the line; and

2. *queuing costs*, associated with the delay incurred by messages in the network.

The model requires the following notation:

$L$  =set of links in the network

$J$  =total number of priority classes

$1/\mu_j$  =average message length for class  $j \in J$

$I_l$  =set of line types available for link  $l, l \in L$

$Q_{lk}$  =the capacity of line type  $k, k \in I_l$

$S_{lk}$  =the fixed cost of line type  $k, k \in I_l$

$C_{lk}$  =the variable cost of line type  $k, k \in I_l$

$D_j$  =unit cost of delay for messages in class  $j \in J$

$R$  =the set of candidate routes

$\Pi$  =the set of communicating nodes in the network

$S_{pj}$  =set of candidate routes for class  $j$  messages associated with origin-destination pair  $p \in \Pi$ .  $S_p, p \in \Pi$  is defined as  $\cup_{j \in J} S_{pj}$ . The sets of candidate routes for different classes of messages are not necessarily disjoint, i.e.  $\cap_{j \in J} S_{pj}$  is not necessarily empty.

$\lambda_{rj}$  =the class  $j$  message arrival rate for the unique origin-destination pair associated with route  $r \in R$ . Also,  $\lambda_{pj} = \lambda_{rj}, \forall r \in S_{pj}$ .

$\phi_{lj}$  =the class  $j$  message rate on link  $l$

$F_{lj} = \phi_{lj}/\mu_j$  =the class  $j$  bit rate on link  $l$

$T_{lj}$  =the average delay incurred on link  $l$  by a class  $j$  message.

$\delta_{rl}$  =an indicator function, taking the value one if link  $l$  is used in route  $r$ , and zero otherwise

$x_{rj}$  =a *decision variable*, taking the value one if route  $r$  is chosen to carry the class  $j$  flow of its associated origin-destination pair, and zero otherwise.

$y_{lk}$  =a *decision variable*, which is one if line type  $k$  is assigned to link  $l$ , and zero otherwise

$T_{lj}$ , the average delay on link  $l$  for class  $j$  messages, can be computed as (see [10]):

$$T_{lj} = \frac{s_{lj}(1 - \sigma_j) + \sum_{i=j}^J \phi_{li}(s_{li})^2}{(1 - \sigma_j)(1 - \sigma_{j+1})} \quad (1)$$

where  $s_{lj} = 1/\mu_j \sum_{k \in I_l} Q_{lk} y_{lk}$  is the average 'service time' for class  $j$  messages on link  $l$ , and  $\sigma_j = \sum_{i=j}^J \phi_{li} s_{li}$ .  $T_{lj}$  includes both the queuing delay incurred by a message while waiting in the buffers of a network switch before transmission, as well as the transmission time.

As a result, the average end-to-end delay in the network for class  $j$  messages can be expressed as:

$$T_j = \frac{1}{\gamma_j} \sum_{i \in L} \phi_{ij} T_{ij}$$

where  $\gamma_j = \sum_{p \in \Pi} \lambda_{pj}$  is the total external arrival rate for class  $j$ .

The above expression becomes untractably complex as the number of message classes increases. We will therefore concentrate in the following on the case of a network supporting just two classes of messages and, without loss of generality, assume that the higher priority is associated with the second class. It is an important case, as indicated by the sharp distinction, both in terms of their requirements as well as of their processing characteristics, between traffic generated by interactive computation, with its tight delay requirement, on one hand, and such applications as file transfer and remote job entry, for which response time is less of a critical factor, and which as a result may be associated with a lower priority, on the other.

From (1), the following expressions are obtained for the average delay on link  $l$  for class 1 and class 2 messages, respectively:

$$T_{l1} = \frac{(Q_l - F_{l2})/\mu_1 + F_{l2}/\mu_2}{(Q_l - F_{l1} - F_{l2})(Q_l - F_{l2})} \quad (2)$$

$$T_{l2} = \frac{1}{\mu_2(Q_l - F_{l2})} \quad (3)$$

where the average class  $j$  bit flow on link  $l$  can be expressed in terms of the decision variables  $x_{rj}$  as:

$$F_{lj} = \sum_{r \in R} \lambda_{rj} \delta_{rl} x_{rj} / \mu_j \quad j = 1, 2 \quad (4)$$

The problem of optimally assigning primary routes and link capacities in a network supporting two classes of messages is then equivalent to finding the binary variables  $x_{rj}$  and  $y_{lk}$  values that satisfy:

#### Problem P

$$Z_P = \min \left\{ D_1 \sum_{l \in L} \frac{F_{l1}(Q_l - F_{l2}) + a F_{l1} F_{l2}}{(Q_l - F_{l1} - F_{l2})(Q_l - F_{l2})} + D_2 \sum_{l \in L} \frac{F_{l2}}{Q_l - F_{l2}} \right. \\ \left. + \sum_{k \in I_l} S_{lk} y_{lk} + C_{lk} y_{lk} (F_{l1} + F_{l2}) \right\} \quad (5)$$

subject to:

$$F_{l1} + F_{l2} \leq \sum_{k \in I_l} Q_{lk} y_{lk} \quad \forall l \in L \quad (6)$$

$$\sum_{r \in S_{pj}} x_{rj} = 1 \quad \forall p \in \Pi \quad j = 1, 2 \quad (7)$$

$$\sum_{k \in I_l} y_{lk} = 1 \quad \forall l \in L \quad (8)$$

$$x_{rj} = 0, 1 \quad \forall r \in R \quad j = 1, 2 \quad (9)$$

$$y_{lk} = 0, 1 \quad \forall l \in L, k \in I_l \quad (10)$$

where  $F_{lj}$ ,  $j = 1, 2$  are defined by (4), and  $a = \mu_1/\mu_2$ .

The first two objective function terms capture the total cost of delay for the lower priority and higher priority message classes, respectively. The third term corresponds

to the total fixed capacity cost, while the fourth represents the total variable cost. The constraints in (6) ensure that the chosen capacity is feasible in terms of the flow assigned to the link. They are equivalent to the constraint set of the NP-complete multiconstrained knapsack problem. The problem studied here is therefore of at least the same complexity. Constraints in (7) and (8) guarantee that only one route is chosen for each origin-destination pair, and only one line type for each link, respectively. Notice that, since  $S_{p1}$  and  $S_{p2}$  are not necessarily disjoint, the formulation allows for the two types of flow to be directed either along the same or along different routes.

The nature of the problem imposes certain restrictions upon the characteristics of the higher priority messages. On an average, they must be shorter than the low priority messages, and they have to pay for the increase in performance they require. As a result, the following relations must hold among the problem parameters:

1.  $a \leq 1$ , i.e. the average length of class 2 messages cannot exceed that of class 1 messages, and
2.  $D_2 \geq D_1$ , i.e. the unit cost of delay is at least as high for class 2 messages as for class 1 messages.

The unit costs of delay are estimates based on user requirements. The model implicitly takes into account the different delay requirements of the two message classes. Priority messages, with their tighter response time requirement, are associated with a higher cost of delay, which lowers the average delay they incur in the final solution. As a result, it is no longer necessary to introduce delay bounds specific to each message class (such as, for instance, the approach used in [12] and [14]), and the structure of the constraint set is significantly simplified.

To better evidence the underlying structure of the problem, a set of derived decision variables is next introduced. The  $f_{lj}$ ,  $j = 1, 2$  variables are defined as the portion of the utilization of link  $l$  attributable to type  $j$  flow:

$$f_{lj} = \frac{\sum_{r \in R} \lambda_{rj} \delta_{rl} x_{rj} / \mu_j}{\sum_{k \in I_l} Q_{lk} y_{lk}} \quad \forall l \in L, j = 1, 2$$

In terms of the augmented set of decision variables, the problem becomes that of finding the  $f_{lj}$ ,  $x_{rj}$  and  $y_{lk}$  variables that satisfy:

#### Problem P<sub>1</sub>

$$Z_P = \min \left\{ \sum_{l \in L} D_1 \frac{f_{l1}(1 - f_{l2}) + a f_{l1} f_{l2}}{(1 - f_{l2})(1 - f_{l1} - f_{l2})} + D_2 \frac{f_{l2}}{1 - f_{l2}} \right. \\ \left. + C_{lk} Q_{lk} (f_{l1} + f_{l2}) y_{lk} \right\} \quad (11)$$

subject to:

$$\sum_{r \in R} \lambda_{rj} \delta_{rl} x_{rj} / \mu_j \leq f_{lj} \sum_{k \in I_l} Q_{lk} y_{lk} \quad \forall l \in L, j = 1, 2 \quad (12)$$

$$f_{l1} + f_{l2} \leq 1 \quad \forall l \in L \quad (13)$$

$$f_{lk} \geq 0 \quad \forall l \in L, k \in I_l \quad (14)$$

and: (7)-(10).

### 3 Solving the Model

A Lagrangian relaxation to Problem P1 is obtained by multiplying the capacity constraints in (12) by a vector of non-positive Lagrange multipliers  $\{\alpha_{lj}, l \in L, j = 1, 2\}$ , and adding them to the objective function. With the coupling constraints no longer present, the Lagrangian problem can be decomposed into a problem depending only on the link decision variables  $f_{lj}$  and  $y_{lk}$ , and a second problem over the routing variables  $x_{rj}$ . Each of these problems, in turn, can be further decomposed over the links in the network, and over the origin-destination pairs and message classes, respectively.

The  $|J| \times |\Pi|$  subproblems associated with a given traffic type for each of the communicating pairs, have a simple structure:

**Problem P1( $\alpha, p, j$ )**

$$L(\alpha, p, j) = \min \left\{ \sum_{r \in S_{pj}} a_{rj} x_{rj} \right\}$$

subject to:

$$\begin{aligned} \sum_{r \in S_{pj}} x_{rj} &= 1 \\ x_{rj} &= 0, 1 \quad r \in S_{pj} \end{aligned}$$

where:  $a_{rj} = \sum_{l \in L} -\alpha_{lj} \lambda_{lj} \delta_{rl} / \mu_j Q_l$ .

The subproblems are readily solved by setting to one that  $x_{rj}$  variable that has the lowest coefficient in the objective function, i.e.

$$a_{bj} = \min_{r \in S_{pj}} \{a_{rj}\} \Rightarrow x_{bj} = 1$$

The  $|L|$  link subproblems resulting from the decomposition are more complex:

**Problem P1( $\alpha, l$ )**

$$L(\alpha, l) = \min \left\{ D_1 \frac{f_{l1}(1-f_{l2}) + a f_{l1} f_{l2}}{(1-f_{l2})(1-f_{l1}-f_{l2})} + D_2 \frac{f_{l2}}{1-f_{l2}} + \sum_{k \in I_l} S_{lk} y_{lk} \right. \\ \left. + f_{l1} \sum_{k \in I_l} Q_{lk} y_{lk} (C_{lk} + \alpha_{l1}) + f_{l2} \sum_{k \in I_l} Q_{lk} y_{lk} (C_{lk} + \alpha_{l2}) \right\}$$

subject to:

$$f_{l1} + f_{l2} \leq 1 \quad (15)$$

$$\sum_{k \in I_l} y_{lk} = 1 \quad (16)$$

$$f_{l1}, f_{l2} \geq 0 \quad (17)$$

$$y_{lk} = 0, 1 \quad (18)$$

The set of candidate capacities is likely to be of small cardinality. It is therefore possible to simplify the above problem by successively fixing the  $y_{lk}$  variables to all the possible values that satisfy the constraints in (16) and (18).

The subproblem becomes:

**Problem P1( $\alpha, l, k$ )**

$$L(\alpha, l, k) = \min \left\{ D_1 \frac{f_{l1}(1-f_{l2}) + a f_{l1} f_{l2}}{(1-f_{l2})(1-f_{l1}-f_{l2})} + D_2 \frac{f_{l2}}{1-f_{l2}} \right. \\ \left. + f_{l1} y_{lk} Q_{lk} (C_{lk} + \alpha_{l1}) + f_{l2} y_{lk} Q_{lk} (C_{lk} + \alpha_{l2}) + S_{lk} \right\}$$

subject to: (15) and (17), where the  $k$  index corresponds to the  $y_{lk}$  variable set to one.

The numerical solution to the subproblem, is based on the following theorem:

**Theorem 1** The objective function of Problem P1( $\alpha, l, k$ ) is unimodal over  $\Omega = \{f_{l1}, f_{l2} : f_{l1} + f_{l2} \leq 1, f_{l1}, f_{l2} \geq 0\}$ .

The theorem proof contains a lengthy argument that is left out for the sake of brevity. The interested reader is referred to [16] for further details.

The result in theorem 1 implies that any algorithm that numerically searches for the minimum within the domain over which the function is defined is guaranteed to converge to a global optimum. Initial experiments showed that a simple successive substitution method has a good convergence rate. The procedure alternately optimizes the function for fixed  $f_{l2}$  and fixed  $f_{l1} + f_{l2}$  values until no further improvement is obtained in two subsequent iterations. Theorem 1 ensures that the unique minimum is reached at this point.

The objective function value for Problem P1( $\alpha, l$ ) is computed as:

$$L(\alpha, l) = \min_{k \in I_l} L(\alpha, l, k)$$

Once all the subproblems are solved, the Lagrangian value is given by:

$$L(\alpha) = \sum_{p \in \Pi} L(\alpha, p) + \sum_{l \in L} L(\alpha, l)$$

It is a known result in optimization theory [7] that the best lower bound is provided by the vector  $\alpha^*$  that corresponds to:

$$L(\alpha^*) = \max_{\alpha \leq 0} L(\alpha) \leq Z_P$$

The following theorem states the relationship that exists between  $L(\alpha^*)$  and the continuous relaxation of Problem P1.

**Theorem 2**

$$L(\alpha^*) = \bar{Z}$$

where:  $\bar{Z}$  is the objective function value of the continuous relaxation of Problem P1.

The proof, based on duality theory, is similar to the one presented in [6] for the no priority case.

### 4 Subgradient Optimization and Heuristic Procedures

A subgradient procedure, an iterative method successfully applied to a variety of combinatorial problems (e.g. [1], [6]).

[8]), is used in order to obtain a good estimate of  $\alpha^*$ .

The domain of the original problem is only a subset of the one over which the Lagrangian problem is defined. It is possible to further tighten the lower bound by generating redundant constraints that will restrict the domain of the Lagrangian (though not that of the original Problem P1). These constraints aim at recapturing part of the problem description lost through relaxation. They quantify some of the implications that the choices available for route assignment have on the values that the flows on a link may be allowed to take, by making part of the structure of the set of candidate routes 'known' to the link subproblems. Specifically, upper and lower bounds on the values of the  $f_{ij}$  variables are computed as:

$$0 \leq L_{ij} \leq f_{ij} \leq \min\{U_{ij}, 1\} \leq 1 \quad (19)$$

where:

$$L_{ij} = \sum_{p \in A_{ij}} \lambda_{pj} / \mu_j \sum_{k \in I_l} Q_{lk} y_{lk}, j = 1, 2$$

$$U_{ij} = \sum_{p \in B_{ij}} \lambda_{pj} / \mu_j \sum_{k \in I_l} Q_{lk} y_{lk}, j = 1, 2$$

and  $A_{ij} = \{p : \delta_{r,l} = 1 \text{ for all } r \in S_{pj}\}$ , i.e. the set of all origin-destination pairs whose primary route for class  $j$  must use link  $l$ , and  $B_{ij} = \{p : \delta_{r,l} = 1 \text{ for some } r \in S_{pj}\}$ , i.e. the set of origin-destination pairs for which link  $l$  might belong to the primary route chosen for class  $j$  messages.

It is important to obtain good upper bounds, not only because they represent a benchmark against which, in the absence of the optimal solution, the quality of the lower bound provided by the Lagrangian can be measured, but foremost because they represent feasible solutions to the original problem. If the gap between the two bounds is reasonably small, the solution corresponding to the upper bound can confidently be used instead of the optimal one. The following ideas were incorporated into the procedure that searches for feasible solutions:

1. At each subgradient iteration, the Lagrangian solution is checked for feasibility in terms of the relaxed constraints.
2. *Randomization*: The following observation considerably increases the power of the algorithm to identify feasible solutions: whenever Problem P1( $\alpha, p, j$ ) is solved, it is often the case that the same reduced cost is associated with more than one route (where 'same' is meant to mean within an  $\epsilon \ll 0$  away from the minimum). A list of such routes is kept for each origin-destination pair, and each message class, out of which candidate solutions are later randomly selected and checked for feasibility in terms of the capacity assignment provided by the Lagrangian solution.
3. *Capacity improvement*: Taking advantage of the small cardinality of the  $I_l$  set, and of the fact that the objective function is decomposable over the links in the network, it is a simple task to determine, for any given flow assignment, what the least costly feasible capacity assignment is.

4. *Route improvement*: It is possible to determine the best flow assignment for a given capacity assignment by solving the following problem:

**Problem F**

$$Z_P = \min \left\{ \sum_{l \in L} D_1 \frac{f_{1l}(1-f_{12}) + a f_{1l} f_{12}}{(1-f_{12})(1-f_{1l}-f_{12})} + D_2 \frac{f_{12}}{1-f_{12}} \right. \\ \left. + C_{lk} Q_{lk} (f_{1l} + f_{12}) y_{lk} \right\}$$

subject to: (8), (9), (12), (14), and (19). The capacity constraints are once again relaxed, and Problem F is solved using a procedure similar to the one outlined earlier. The significant difference is that, with the capacity variables no longer present, not only the Lagrangian problem is more efficiently solved, but also the algorithm converges very fast, solution tolerances of under 5% being generally obtained in less than 40 subgradient iterations.

5. *Local optimum*: The route improvement procedure sometimes significantly alters the flow pattern, which may render the current capacity assignment no longer optimal. The algorithm therefore alternates between the capacity and the route improvement steps, until no further reduction in the overall cost can be achieved. Since this search for a local optimum may at times be quite time consuming, in the current implementation of the algorithm it is initiated only at the user's specific request.

## 5 Computational Results

The model and algorithm presented in the previous sections are implemented as an interactive system that allows the network topologies and model parameters to be easily defined and modified by the user. At the end of each major iteration (defined as a user specified number of subgradient iterations), a comprehensive output, corresponding to the current best feasible solution, is produced. In addition to the current value of the Lagrangian, the overestimate, and the corresponding average message delays, the output also gives a detailed description of the capacity assignment, specifying for each link its assigned capacity, its message rate and utilization, the associated fixed, variable, and queuing costs, and the percentage of the total cost attributable to it. Thus, the user is presented with a full picture of the current solution that can be used as a basis for gaining further insights into the characteristics of the problem under consideration.

Three different topologies were used in the experiments (fig. 1-3). The total average message traffic for all origin-destination pairs is of four messages for both directions, and is evenly divided between the two types of flow. The candidate routes were generated using the same hybrid procedure outlined in [5]. Table 1 shows the candidate line types used (the same for all the links in a network), and the associated

capacity costs.

The tolerance measure used for estimating the quality of the results generated by the algorithm is defined as:

$$\frac{(\text{Upper Bound} - \text{Lower Bound})}{\text{Upper Bound}}$$

The results summarized in the following tables show that the gaps between the lower and the upper bound the algorithm generates are between satisfactory, e.g. 15.5%, and very good, e.g. less than 1%. There does not seem to be an obvious correlation between the parameter values used and the behavior of the algorithm, but the matter requires further investigation.

The experiments were run on a VAX 780 machine. No exact values for the computation times were collected, but they were reasonably small, considering the off-line nature of the problem. The maximum number of subgradient iterations never exceeded 200, and was usually under 100. The amount of time used per iteration was negligible, even for the larger test problem.

From table 1, it can be observed that the average message delay for each message class is only marginally sensitive to the ratio between the two costs of delay, more so for the lower priority messages. On the other hand, as these costs are decreased, so that the objective function is even more dominated by the capacity costs, the average delay experienced by both classes of messages increases, the changes being more significant, once more, for the lower priority traffic.

Table 2 shows the results for fixed costs of delay, and varying average message lengths. As expected, as the message length increases/decreases, the average delay varies accordingly, for both message classes. The variations are rendered more significant by the changes in the capacity assignment.

Further testing of the model is required, before final conclusions can be reached. Nevertheless, the initial results seem to justify the introduction of the priority discipline, and are promising in terms of the quality of the solutions generated by the algorithm.

CAPACITY [bps]	SETUP COST [dollars/month]	DISTANCE COST [dollars/month/mile]	VARIABLE COST [dollars/month/bps]
4800	650	0.4	.360
9600	750	0.5	.252
19200	850	2.1	.126
50000	850	4.2	.030
108000	2400	4.2	.024
230000	1300	21.0	.020
460000	1300	60.0	.017

Table 1: Capacity set and base costs used in computational experiments

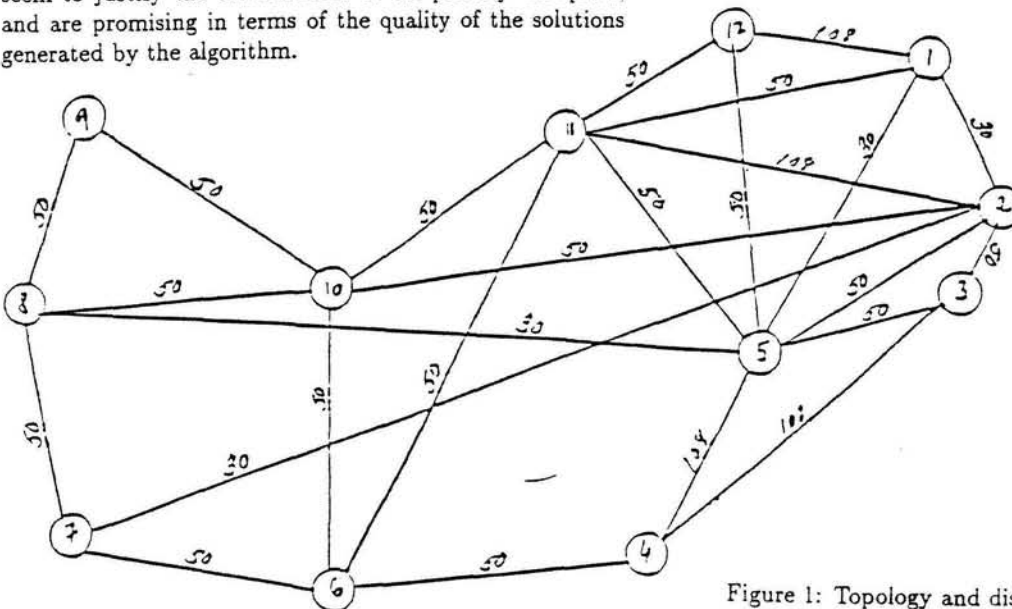


Figure 1: Topology and distances for the GTE network

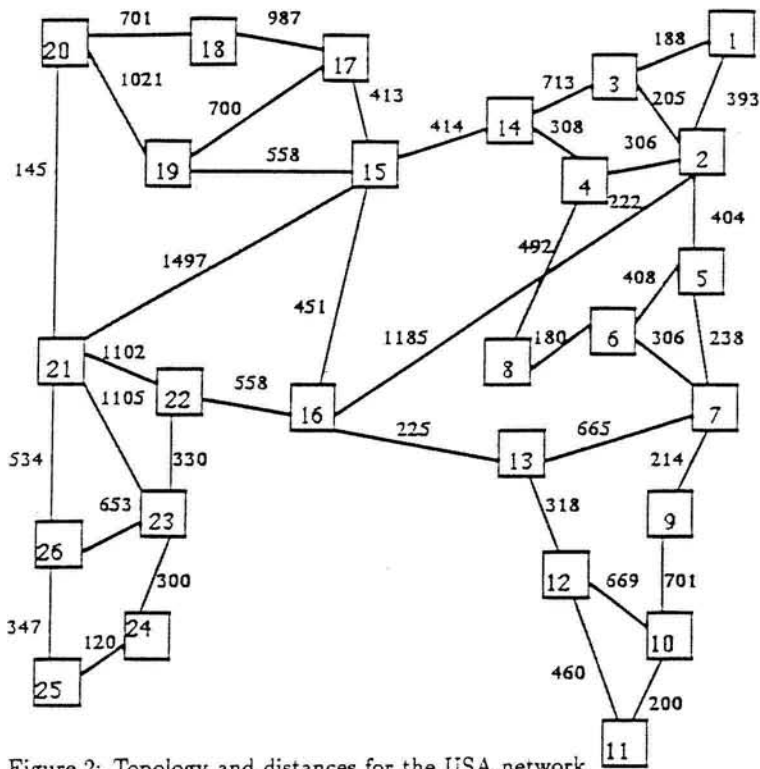


Figure 2: Topology and distances for the USA network

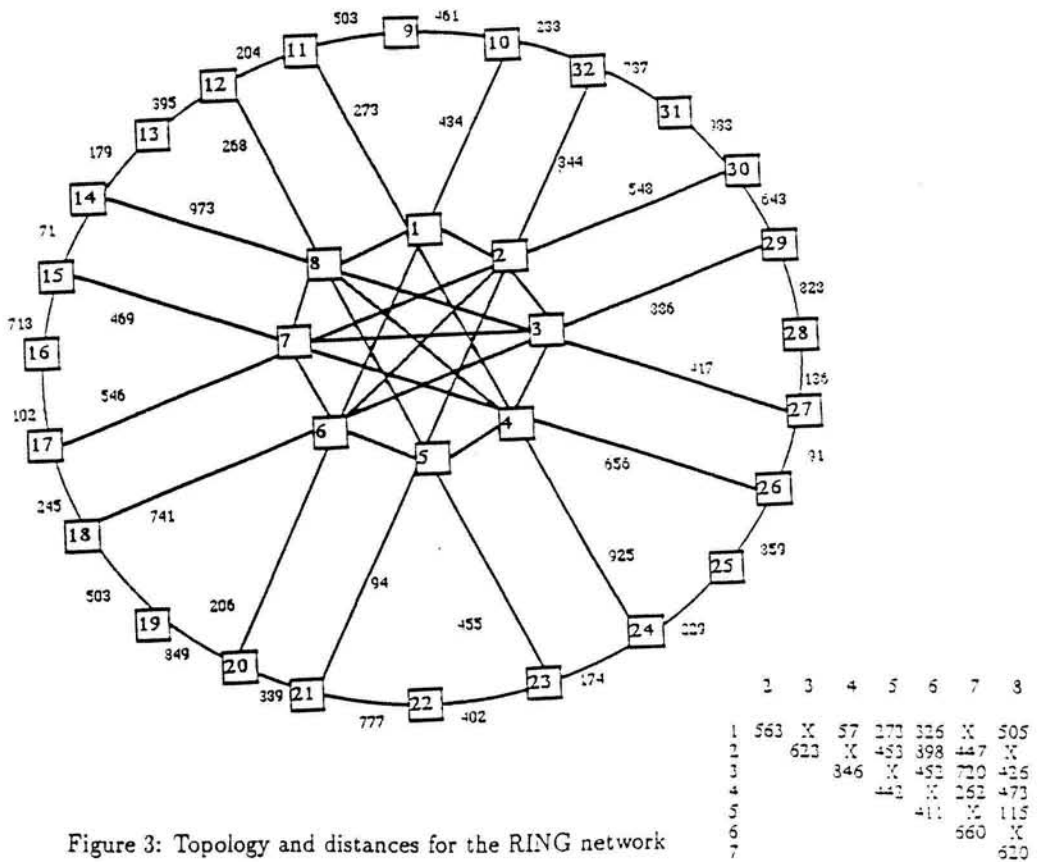


Figure 3: Topology and distances for the RING network



Network	Delay costs	Lower bound	Upper bound	Queuing costs	Fixed cost	Variable cost	Tolerance (%)	Average delays
GTE	cl.1=2000 cl.2=6000	131010	138077	18785 8080	82662	28549	5.12	33.3 13.0
	cl.1=2000 cl.2=2100	122609	132829	18896 3223	79904	30805	7.69	36.6 14.6
	cl.1=2000 cl.2=20000	155902	158340	19232 26172	85680	27255	1.54	27.8 11.3
	cl.1=20 cl.2=60	85270	101861	867 220	56103	44671	16.29	11.4 27.8
USA	cl.1=2000 cl.2=6000	1495758	1556373	201303 18897	1230818	105355	3.94	92.3 8.0
	cl.1=2000 cl.2=2100	1414826	1565693	260497 8236	1197921	99039	9.64	100.9 9.9
	cl.1=2000 cl.2=20000	1623522	1635090	260370 76433	1200230	98056	.71	78.1 9.0
	cl.1=20 cl.2=60	1110111	1276315	3558 366	1160268	112122	13.02	199.2 12.2
RING	cl.1=2000 cl.2=6000	982060	1083677	215312 32883	726937	108544	9.38	42.3 5.0
	cl.1=2000 cl.2=2100	1002690	1105734	276904 12921	709066	406875	9.32	45.2 5.0
	cl.1=2000 cl.2=20000	1008378	1125126	223809 87426	705125	108766	10.38	45.9 4.4
	cl.1=20 cl.2=60	669869	793296	4047 398	650622	138229	15.56	93.5 5.4

Note:  $1/\mu_1 = 1000$  and  $1/\mu_2 = 400$  for the GTE and USA networks  
 $1/\mu_1 = 700$  and  $1/\mu_2 = 300$  for the RING network

Table 2: Results for different costs of delay

Network	Message length	Lower bound	Upper bound	Queuing costs	Fixed cost	Variable cost	Tolerance (%)	Average delays
GTE	cl.1=1000 cl.2=400	131010	138077	18785 8080	82662	28549	5.12	33.3 13.0
	cl.1=1400 cl.2=400	147636	160179	26223 8413	95389	30154	7.83	52.0 11.9
	cl.1=1000 cl.2=200	115000	127724	181456 4732	74337	30499	9.96	33.5 7.4
USA	cl.1=1000 cl.2=400	1495088	1556373	201303 18897	1230818	105354	3.94	92.3 8.0
	cl.1=1400 cl.2=400	1813868	1950383	224121 72291	1533981	119989	7.00	86.2 9.3
	cl.1=1000 cl.2=200	1239601	1420852	198750 10755	1108945	102401	12.76	78.8 4.1
RING	cl.1=700 cl.2=300	982060	1083677	2153112 32883	726937	108544	9.38	42.3 5.0
	cl.1=1000 cl.2=300	1822029	1918665	614213 109158	1044621	150672	5.04	77.4 4.6
	cl.1=700 cl.2=150	964094	977740	217547 16601	641289	102303	1.40	27.4 2.2

Table 3: Results for different average message lengths  
(D1=2000, D2=6000)

## References

- [1] P. Afentakis and B. Gavish. Computationally efficient optimal solutions to the lot sizing problem in multistage assembly problems. *Manag. Sc.*, 30:222-239, 1984.
- [2] V. Ahuja. Routing and flow control in systems network architecture. *IBM Syst. Jour.*, 18:298-314, 1979.
- [3] J.N. Daigle and C.E. Houstis. Analysis of a task oriented multipriority queuing system. *IEEE Trans. Commun.*, COM-29:1660-1677, 1981.
- [4] B. Dureste and J. Kravitz. The sita telecommunications network. In *IEEE Intern. Conf.*, pages 1067-1069, 1983.
- [5] B. Gavish and I. Neuman. Capacity and flow assignment in large computer networks. In *Proceedings, INFOCOM'86, Miami, Florida*, pages 275-284, april 1986.
- [6] B. Gavish and I. Neuman. A system for routing and capacity assignment in computer communication networks. *IEEE Trans. Commun.* (in print).
- [7] A.M. Geoffrion. Lagrangean relaxation and its uses in integer programming. *Mathematical Programming Study*, 2:82-114, 1974.
- [8] M. Held and R.M. Karp. The travelling salesman problem and minimum spanning trees, part ii. *Math. Program.*, 1:6-25, 1971.
- [9] L. Kleinrock. *Communication nets: stochastic message flow and delay*. McGraw-Hill, 1964.
- [10] L. Kleinrock. *Queuing systems, vol. 2: Computer applications*. John Wiley & Sons, 1976.
- [11] R.V. Laue. A versatile queuing model for data switching. In *Symp. on Data Commun.*, pages 118-128, 1981.
- [12] K. Maruyama, K. Fratta, and D.T. Tang. Heuristic design algorithm for computer communication networks with different classes of customers. *IBM Jour. Res. Develop.*, 360-369, 1977.
- [13] K. Maruyama and D. T. Tang. *Discrete link capacity and priority assignment in communication networks*. Technical Report RC 6121, IBM Thomas J. Watson Research Center, Yorktown Heights, New York, 1976.
- [14] K. Maruyama and D.T. Tang. Discrete link capacity assignment in communication networks. In *3rd ICC, Toronto*, pages 92-97, 1976.
- [15] R.J.T. Morris. Priority queuing networks. *Bell Sys. Tech. Jour.*, 1745-1769, 1981.
- [16] I. Neuman. *Routing in a network with different classes of messages*. Technical Report, Center for Research on Information Systems, Information Systems Area, Graduate School of Business Administration, New York University, 1986.
- [17] R.N. Pandya. Delay analysis for datapac: a packet switched network with two priority classes. In *5th Data Commun. Symp., Snowbird, Utah*, pages 314-121, 1977.
- [18] M. Reiser. A queuing network analysis of computer communication networks with window flow control. *IEEE Trans. Commun.*, COM-27:1199-1209, 1979.
- [19] E. Soueid and J.J. Metzner. Priority queuing models for congestion control and performance analysis of a distributed computer network. In *IEEE Computer networking symp.*, pages 105-112, 1981.
- [20] *TELENET Communications Corporation - Packet switching network*. Auerbach, New York, 1978.