

# The Auction Algorithm: A Distributed Relaxation Method for the Assignment Problem

D.P. Bertsekas

*Laboratory for Information and Decision Systems  
Massachusetts Institute of Technology  
Cambridge, Mass. 02139*

We propose a massively parallelizable algorithm for the classical assignment problem. The algorithm operates like an auction whereby unassigned persons bid simultaneously for objects thereby raising their prices. Once all bids are in, objects are awarded to the highest bidder. The algorithm can also be interpreted as a Jacobi - like relaxation method for solving a dual problem. Its (sequential) worst - case complexity, for a particular implementation that uses scaling, is  $O(N \log(NC))$ , where  $N$  is the number of persons,  $A$  is the number of pairs of persons and objects that can be assigned to each other, and  $C$  is the maximum absolute object value. Computational results show that, for large problems, the algorithm is competitive with existing methods even without the benefit of parallelism. When executed on a parallel machine, the algorithm exhibits substantial speedup.

## 1. INTRODUCTION

Relaxation methods for optimal network flow problems resemble classical coordinate descent, Jacobi, and Gauss-Seidel methods for solving unconstrained nonlinear optimization problems or systems of nonlinear equations. They operate on a dual problem which is unconstrained and involves a dual variable for every node (also called a node price). In their pure form they modify the node prices one at a time using only local node information while aiming to improve the dual cost. They are well suited for distributed implementation on massively parallel machines.

For problems with strictly convex arc costs, relaxation methods can be shown to converge to the correct solution under mild additional assumptions. An important property of strictly convex cost problems is that the dual cost function is differentiable. This facilitates the use of the relaxation idea because at any nonoptimal dual vector one can find a single node price coordinate along which the dual cost can be improved. However, the convergence properties of the method do not follow from classical results of unconstrained optimization because the dual cost function is not strictly convex and does not have bounded level sets. Nonetheless, there is sufficient structure to guarantee

\* Work supported by Grant NSF-ECS-8217668. Thanks are due to J. Kennington and L. Hatay of Southern Methodist Univ. for contributing some of their computational experience.

convergence, which can also be shown for a totally asynchronous implementation where price updates at each node are carried out asynchronously with out-of-date price information from neighboring nodes [15], [18]. Computational experiments with parallel relaxation methods applied to strictly convex arc cost problems have been very encouraging [47], [48].

By contrast, for problems where the arc costs are linear the dual cost function is nondifferentiable (piecewise linear). As a result the relaxation idea may encounter difficulty at points where the dual cost cannot be improved in any coordinate direction as seen in Fig. 1.

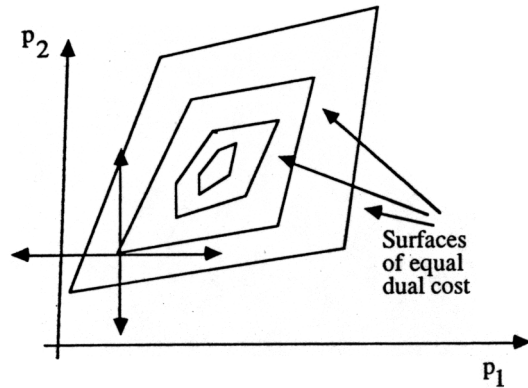


FIGURE 1: Illustration of the difficulty of the relaxation approach for linear arc costs. At the indicated point it is impossible to improve the dual cost by changing any single price.

One approach to overcome the difficulty has been pursued by the author and his coworkers for the past several years [7], [9], [16], [17]. In this approach, in addition to the single coordinate relaxation steps, one occasionally changes the prices of several nodes as a group as illustrated in Figure 2. Unfortunately, however, these multiple node price changes need *global* node price information, thereby substantially complicating the distributed implementation of the algorithm. Nonetheless, public domain sequential codes based on this approach have proved highly successful, and have outperformed the classical primal simplex and primal-dual methods by a substantial margin on standard benchmark problems [9], [16], [17] and by an overwhelming margin on large problems [16], [17].

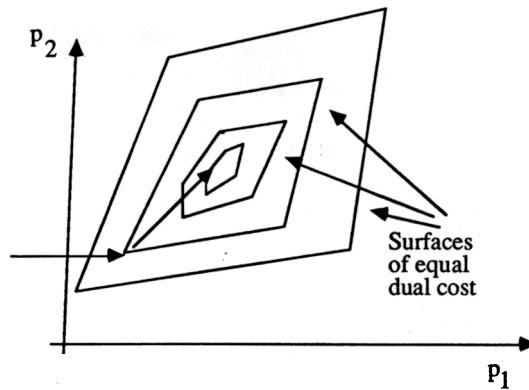


FIGURE 2: Illustration of the first approach for overcoming the difficulty of Fig. 1. The idea is to reduce the dual cost by changing simultaneously the prices of several nodes.

A second approach for overcoming the difficulty shown in Fig. 1, that is more amenable to parallelization, is based on performing single node price changes exclusively. Some of these price changes may worsen the dual cost by small amounts, but Fig. 3 illustrates that if the price changes are properly regulated, the algorithm can eventually approach the optimal solution.

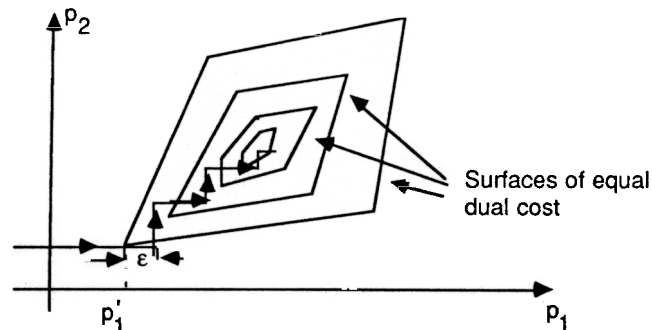


FIGURE 3: Relaxation approach whereby single node price changes are performed exclusively even if this leads to deterioration of the dual cost. If the size of the price changes is properly regulated, the method can approach the optimal solution.

This second approach was first proposed for the classical assignment problem in 1979 by the author in an unpublished report [6], and was also briefly described in the conference paper [8]. The corresponding method, called the *auction algorithm*, is the subject of the present paper. Among the many methods for the assignment problem [2] - [4], [7], [24], [26], [27], [28], [32] - [34] [37], [38], [46] the auction algorithm seems to be the only one that has a naturally parallel character and is well suited for implementation on a massively parallel machine. It also admits an intuitively appealing economic interpretation, although we make no claim about its utility as a model for market price adjustment. (See [23], [43], [44] for work on the assignment problem with an

economic modeling orientation. In particular, [23] describes an algorithm which involves ideas of competitive bidding, but differs substantially from ours; for example it does not use the idea of  $\epsilon$ -complementary slackness described later in this paper.)

In addition to the idea of independent single node price changes which is responsible for the parallel character of the method, the auction algorithm embodies several ideas that were new at the time of its original proposal and have recently found considerable extension and use in network algorithms and computational complexity analysis. These include the notion of  $\epsilon$ -complementary slackness ( $\epsilon$ -CS) and the idea of  $\epsilon$ -scaling.

In Section 2, we define  $\epsilon$ -CS as a condition whereby the usual complementary slackness relation for assignment of a person to an object can be violated by a small ( $\epsilon$ ) amount. A key fact is that a feasible assignment satisfying  $\epsilon$ -CS is optimal if  $\epsilon$  is sufficiently small. This idea was generalized and was used algorithmically for linear network flow problems in [16], and for nonlinear network flow problems in [18]. A related idea was also used in a different context in the analysis of [45].  $\epsilon$ -CS is central in the extension of the auction algorithm to general linear network flow problem, first given in [11], [12], and further analyzed later in [14], [30], [31].

The idea of  $\epsilon$ -scaling is similar to the idea of gradually reducing the penalty parameter in penalty function methods in order to alleviate ill-conditioning, and is essential for obtaining computational efficiency and polynomial complexity. In the pure form of the auction algorithm,  $\epsilon$  is kept constant, and the complexity (with the use of some simple data structures [19]) can be shown to be  $O(NAC/\epsilon)$ . It is necessary to take  $\epsilon$  sufficiently small (less than  $1/N$  assuming object values are integer), in order that the assignment obtained by the algorithm be optimal (see Proposition 1 in Section 3). As a result, the pure form of the algorithm has pseudopolynomial complexity, and its computational performance can be poor as shown by an example given in Section 3.  $\epsilon$ -scaling corrects this by applying the auction algorithm with successively smaller values of  $\epsilon$  starting with a large value and ending with a value less than  $1/N$ . The final prices and assignment corresponding to each value of  $\epsilon$  are used to obtain "good" starting prices and assignment for the next smaller value.  $\epsilon$ -scaling was noted in [6] as a method for improving the performance of the auction algorithm, and was tested computationally in 1979 and again in 1985 with encouraging results (unpublished work). It was first analyzed in the context of the more general min-cost flow problem in [29] where polynomial complexity results were given that were more fully established in [30], [31]. Another scaling procedure that is closer in spirit to traditional scaling methods in network flow problems [21], [25], [41], [39], and can serve as an alternative to  $\epsilon$ -scaling was later proposed and analyzed in [14], using some of the ideas of [29] and also of [21]. The method of analysis of [14], [30], [31] can be used to show that the auction algorithm with  $\epsilon$ -scaling as well as the more traditional cost scaling method has complexity  $O(N \log(NC))$ ; see [19]. This is relatively straightforward, but requires replication of large portions of the analysis in these references, and so will not be given. The computational results of the present paper

suggest that the auction algorithm with  $\epsilon$ -scaling is competitive with the best known implementations of assignment methods even without the benefit of parallelism. When implemented in a parallel machine, the auction algorithm should converge much faster.

The ideas of  $\epsilon$ -CS and  $\epsilon$ -scaling are likely to prove useful in other contexts. For example they play an important role in an extension of the auction algorithm that solves transportation problems [13], and also in the recent assignment algorithm of [27]. The latter algorithm, which is not naturally parallelizable, has complexity  $O(N^{1/2} \text{Alog}(NC))$  where  $N$  is the number of persons to be assigned,  $A$  is the number of person - object pairs that can be assigned to each other, and  $C$  is the maximum absolute object value. For reasonably small values of  $C$ , this is the best complexity bound for assignment problems. We note also that the auction algorithm can be modified to obtain an  $O(N^{1/2} \text{Alog}(NC))$  complexity. This modification, due to [1], employs  $\epsilon$ -scaling, and within a scaling phase, it first uses the auction algorithm to assign  $N - O(N^{1/2})$  persons, and then uses a Hungarian-like method to assign the remaining persons.

The purpose of the present paper is to describe the auction algorithm, to interpret it as a dual Jacobi-like relaxation method, and to provide computational results. While a version of the auction algorithm can be derived from its min-cost flow generalization of [11], [12], this derivation is neither obvious nor straightforward; indeed one must introduce modifications to the basic form of the  $\epsilon$ -relaxation method of [11], [12] in order to be able to derive the auction algorithm as a special case. (The reverse is also true. The min-cost flow problem can be converted into a transportation problem ([40], p. 149), which can be converted to an assignment problem by replacing single sources and single sinks of the transportation problem with multiple copies of persons and objects respectively in the assignment problem. By applying the auction algorithm to this assignment problem, one can derive a generic form of the  $\epsilon$ -relaxation method of [11], [12].) The conceptually useful interpretation as a Jacobi-like relaxation method is also nontrivial to obtain starting from the general network flow framework. It is therefore worthwhile to derive the auction algorithm from basic principles, and Section 2 is devoted to this purpose. The validity of the algorithm is established in Section 3, and the computational results are given in Section 5. The auction algorithm can also be implemented in a totally asynchronous (chaotic) distributed environment; this appears to be a generic characteristic of relaxation methods for both linear and nonlinear network flow problems. We describe briefly and somewhat informally in Section 4 how this can be done, and we refer to [11], [12], [19] for discussion and analysis of totally asynchronous implementations in the context of the general min-cost flow problem.

## 2. ASSIGNMENT BY MEANS OF AN AUCTION

Consider  $N$  persons wishing to decide among themselves  $N$  objects. We number persons and objects as  $1, 2, \dots, N$ . For each person  $i$  there is a nonempty subset  $A(i)$  of objects that can be assigned to  $i$ . An *assignment*  $S$  is a (possibly empty) set of person-object pairs  $(i, j)$  such that  $j \in A(i)$  for all  $(i, j) \in S$ , for each person  $i$  there is at most one pair  $(i, j) \in S$ , and for each object  $j$  there is at most one pair  $(i, j) \in S$ . In the context of a given assignment  $S$ , we say that person  $i$  is *assigned* if there exists an object  $j$  such that  $(i, j) \in S$ ; otherwise we say that  $i$  is *unassigned*. We use similar terminology for objects. A *complete assignment* is an assignment containing  $N$  pairs (*i.e.* every person is assigned to a distinct object). There is a given integer value  $a_{ij}$  that a person  $i$  associates with an object  $j \in A(i)$ . We want to find a complete assignment that maximizes

$$\sum_{(i,j) \in S} a_{ij}$$

over all complete assignments  $S$ . We call this the *primal assignment problem* and note its well-known equivalence to a linear programming (linear network flow) problem as shown in Fig. 4.

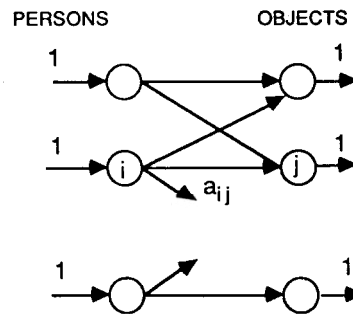


FIGURE 4: Formulation of the primal assignment problem as a linear programming/network flow problem. The nodes are the persons and the objects, and an arc  $(i, j)$  connects every person  $i$  and object  $j$  such that  $j$  belong to  $A(i)$ . Denoting the flow of arc  $(i, j)$  by  $f_{ij}$  the corresponding linear programming problem is:

$$\begin{aligned} & \text{maximize} \sum_{i=1}^N \sum_{j \in A(i)} a_{ij} f_{ij} \\ & \text{subject to} \\ & \sum_{j \in A(i)} f_{ij} = 1, \quad i = 1, \dots, N \\ & \sum_{(i|j \in A(i))} f_{ij} = 1, \quad j = 1, \dots, N \\ & 0 \leq f_{ij} \end{aligned}$$

For each object  $j$ , it is useful to introduce a dual variable  $p_j$  called the *price* of  $j$ . We call the vector with coordinates  $p_j, j = 1, \dots, N$  a *price vector*. For a

given price vector  $p$  the scalar

$$\pi_i = \max_{j \in A(i)} \{a_{ij} - p_j\} \quad (1)$$

is called the *profit margin* of person  $i$  corresponding to  $p$ . It is helpful to think of  $p_j$  as the amount of money that a person must pay when assigned to  $j$ . Therefore, for a given price vector  $p$ ,  $a_{ij} - p_j$  may be thought of as the benefit person  $i$  associates with being assigned to object  $j$ . In this context the name “profit margin” for  $\pi_i$  as given by (1) becomes meaningful.

A dual problem to the assignment problem is

$$\text{minimize } \sum_{i=1}^N \pi_i + \sum_{j=1}^N p_j$$

$$\text{subject to } \pi_i + p_j \geq a_{ij}, \quad \forall i, \text{ and } j \in A(i)$$

For a given price vector  $p$ , the cost of this problem is minimized when  $\pi_i$  equals the maximum value of  $a_{ij} - p_j$  over  $j \in A(i)$ , i.e. the profit margin of  $i$  corresponding to  $p$  as given by (1). Therefore an equivalent form of the dual problem is the *unconstrained* minimization problem

$$\text{minimize } q(p) \quad (2)$$

subject to no constraints on  $p$

where  $p$  is the vector of object prices  $p_j$ , and

$$q(p) = \sum_i \max_{j \in A(i)} \{a_{ij} - p_j\} + \sum_j p_j \quad (3)$$

From linear programming theory it is known that a complete assignment  $S = \{(i, j_i) | i = 1, \dots, N\}$  and a price vector  $p$  are simultaneously primal and dual optimal respectively if and only if

$$\pi_i = \max_{k \in A(i)} \{a_{ik} - p_k\} = a_{ij} - p_j, \quad \text{for each } (i, j) \in S.$$

This is known as the *complementary slackness condition* and states that at an optimum each person is assigned to a “best” object, i.e. to one attaining the maximum in the definition of profit margin (1).

A relaxation of the complementary slackness condition is to allow persons to be assigned to objects that come within  $\epsilon$  of attaining the maximum in (1). Formally we say that an assignment  $S$  (not necessarily complete) and a price vector  $p$  satisfy  $\epsilon$ -complementary slackness ( $\epsilon$ -CS) if

$$\pi_i - \epsilon = \max_{k \in A(i)} \{a_{ik} - p_k\} - \epsilon \leq a_{ij} - p_j, \quad \text{for each } (i, j) \in S, \quad (4)$$

where  $\pi_i$  is given by (1), and  $\epsilon$  is a nonnegative constant.

We now describe formally the auction algorithm. We fix  $\epsilon > 0$ , and we start with some assignment (possibly empty) and price vector satisfying  $\epsilon$ -CS. The algorithm proceeds iteratively and terminates when a complete assignment is obtained. At the start of the generic iteration we have an assignment  $S$  and a price vector  $p$  satisfying  $\epsilon$ -CS. At the end of the iteration,  $S$  and  $p$  are updated while maintaining the  $\epsilon$ -CS condition. There are two phases in each iteration, the *bidding phase*, and the *assignment phase* described below:

**Bidding Phase:** For each person  $i$  that is unassigned under the assignment  $S$ : Compute the “current value” of each object  $j \in A(i)$  given by

$$v_{ij} = a_{ij} - p_j \quad (5)$$

Find a “best” object  $j^*$  having maximum value

$$v_{ij^*} = \max_{j \in A(i)} v_{ij},$$

and find the best value offered by objects other than  $j^*$

$$w_{ij^*} = \max_{j \in A(i), j \neq j^*} v_{ij} \quad (6)$$

(If  $j^*$  is the only object in  $A(i)$  we define  $w_{ij^*}$  to be  $\infty$ , or, for computational purposes, a number that is much smaller than  $v_{ij^*}$ .)

Compute the “bid” of person  $i$  for object  $j^*$  given by

$$b_{ij^*} = p_{j^*} + v_{ij^*} - w_{ij^*} + \epsilon = a_{ij^*} - w_{ij^*} + \epsilon \quad (7)$$

(We characterize this situation by saying that person  $i$  bid for object  $j^*$ , and that object  $j^*$  received a bid from person  $i$ . The algorithm works if the bid has any value between  $p_{j^*} + \epsilon$  and  $p_{j^*} + v_{ij^*} - w_{ij^*} + \epsilon$ , but it tends to work fastest for the maximal choice (7). Note that from (5)-(7) we have  $b_{ij^*} \geq p_{j^*} + \epsilon$ .)

**Assignment Phase:** For each object  $j$ : let  $P(j)$  be the set of persons from which  $j$  received a bid in the bidding phase of the iteration. If  $P(j)$  is nonempty increase  $p_j$  to the highest bid

$$p_j := \max_{i \in P(j)} b_{ij}, \quad (8)$$

remove from the assignment  $S$  any pair  $(i, j)$  (if one exists), and add to  $S$  the pair  $(i^*, j)$  where  $i^*$  is some person in  $P(j)$  attaining the maximum above.

It is seen that at the end of the iteration we have a new price vector that differs from the preceding vector in the prices of the objects that received a bid during the iteration. We also have a new assignment that differs from the preceding one in that each object that received a bid is now assigned to some person that was unassigned at the start of the iteration. However, the assignment at the end of the iteration need not have more pairs than the one at the start of the iteration because it is possible that all objects that received a bid were assigned at the start of the iteration.

An important fact is that *the algorithm preserves  $\epsilon$ -CS throughout its execution*, i.e. if the assignment and price vector available at the start of an iteration satisfy  $\epsilon$ -CS, the same is true for the assignment and price vector obtained at the end of the iteration. To see this suppose that object  $j^*$  received a bid from person  $i$  and was assigned to  $i$  during the iteration. Let  $p_j$  and  $p'_j$  be the object prices before and after the assignment phase. Then we have [cf. (7), (8)]

$$p'_{j^*} = b_{ij^*} - w_{ij^*} + \epsilon \quad (9)$$

Using this equation and the fact  $p'_j \geq p_j$  for all  $j$ , it follows that



$$a_{ij^*} - p'_{j^*} = a_{ij^*} - b_{ij^*} = w_{ij^*} - \epsilon = \max_{j \in A(i), j \neq j^*} \{a_{ij} - p_j\} - \epsilon$$

$$\geq \max_{j \in A(i), j \neq j^*} \{a_{ij} - p'_j\} - \epsilon$$

This equation implies that

$$a_{ij^*} - p'_{j^*} \geq \max_{j \in A(i)} \{a_{ij} - p'_j\} - \epsilon \quad (10)$$

which shows that (4) continues to hold after the assignment phase of an iteration for a pair  $(i, j^*)$  that entered the assignment during the iteration. Consider also any pair  $(i, j^*)$  that belonged to the assignment just before an iteration, and also belongs to the assignment after the iteration. Then  $j^*$  must not have received a bid during the iteration, so  $p'_{j^*} = p_{j^*}$ . Therefore, (10) holds in view of the  $\epsilon$ -CS condition that holds prior to the iteration, and the fact  $p'_j \geq p_j$  for all  $j$ .

Note that both the bidding and the assignment phases are highly parallelizable. In the extreme case of a fine grain parallel computing environment where there is a processor associated with each person and a processor associated with each object, all unassigned persons/processors can compute their bids simultaneously and communicate them to the relevant objects/processors. Those object/processors that received at least one bid can determine the highest bidder simultaneously and communicate the changes in the current assignment and price vector to the relevant persons/processors.

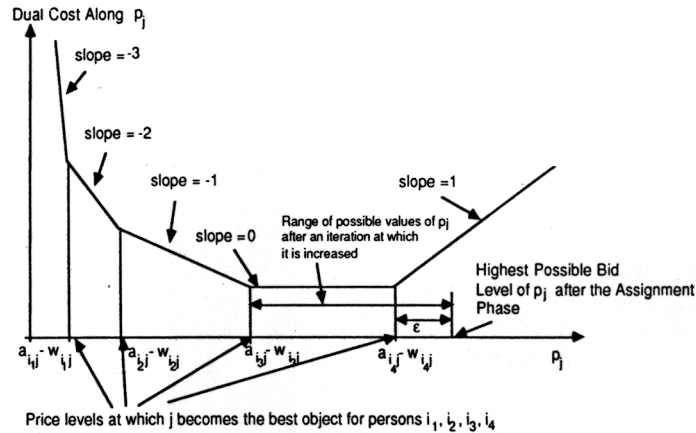


FIGURE 5: Form of the dual cost along the price coordinate  $p_j$ . From (3) the right directional derivative of  $q$  along  $p_j$  is

$$d_j^+ = 1 - (\text{number of persons } i \text{ with } j \in A(i) \text{ and } p_j < a_{ij} - w_{ij})$$

where  $w_{ij}$  is given by (6). The break points are  $a_{ij} - w_{ij}$  for all  $i$  such that  $j \in A(i)$ . If  $p_j < a_{ij} - w_{ij}$ , then an unassigned person  $i$  bids for object  $j$  the amount  $a_{ij} - w_{ij} + \epsilon$ . The price  $p_j$  after the assignment phase is increased to  $\epsilon$  plus the highest level  $a_{ij} - w_{ij}$  over all unassigned persons  $i$  with  $j \in A(i)$ .

Figure 5 indicates how each bidding and subsequent assignment phase can

be interpreted as a Jacobi - like relaxation step for minimizing the dual function  $q(p)$  of (3). In particular, *the price  $p_j$  of each object that received a bid during the assignment phase is increased to either a value that minimizes  $q(p)$  when all other prices are kept constant, or else exceeds the largest such value by no more than  $\epsilon$ .* To see this suppose that at some iteration object  $j$  received a bid and its price was raised from  $p_j$  to  $p'_j$ . Then

$$p'_j = \max\{a_{ij} - w_{ij} | i \text{ was unassigned, } j \in A(i), \text{ and } j \text{ received a bid from } i\} + \epsilon \quad (11)$$

$$p_j \geq \max\{a_{ij} - w_{ij} | i \text{ was unassigned, } j \in A(i), \text{ and } j \text{ did not receive a bid from } i\} \quad (12)$$

Since whenever an object receives a bid, its price increases by at least  $\epsilon$ , we have

$$p'_j \geq p_j + \epsilon,$$

so from (11) and (12) we obtain

$$p'_j \geq \max\{a_{ij} - w_{ij} | i \text{ was unassigned, and } j \in A(i)\} + \epsilon \quad (13)$$

Since the algorithm maintains property (4) throughout, we have that if at the start of the iteration person  $i$  was assigned to some  $k \neq j$  and  $j \in A(i)$  then

$$a_{ij} - p_j \leq \pi_i \leq a_{ik} - p_k + \epsilon \leq w_{ij} + \epsilon$$

Using this relation and the fact  $p'_j \geq p_j + \epsilon$  we obtain

$$p'_j \geq p_j + \epsilon \geq a_{ij} - w_{ij}$$

and

$$p'_j \geq \max\{a_{ij} - w_{ij} | i \text{ was assigned to some } k \neq j, \text{ and } j \in A(i)\} \quad (14)$$

Combining (13) and (14) we obtain

$$p'_j \geq \max\{a_{ij} - w_{ij} | i \text{ was not assigned to } j, \text{ and } j \in A(i)\}$$

Since there can be at most one person assigned to  $j$ , it follows from the form of the dual cost shown in Figure 5, that  $p'_j$  is no less than the smallest value of  $p_j$  that minimizes  $q(p)$ . Combining this fact with (11), we see that  $p'_j$  has the property stated in the beginning of this paragraph.

Note that our method is not quite equivalent to a coordinate descent method which reduces the dual cost along each price coordinate, since the dual cost (3) may deteriorate strictly after a price increase. However, the cost deterioration is at most  $\epsilon$ . It will be seen shortly that, for  $\epsilon$  small enough, an optimal solution can still be obtained thanks to the integer nature of the problem data, and the fact that  $\epsilon$ -CS holds at termination [cf. (4)].

Figure 5 also suggests a variation of the algorithm whereby, in addition to all unassigned persons, each *assigned* person  $i$  bids for his *own* assigned object  $j$  the amount  $a_{ij} - w_{ij} + \epsilon$ . This variation has not been tested, but may accelerate convergence in some parallel computing environments.

The above algorithm may be viewed as a *Jacobi version* of the relaxation idea since the bids of unassigned persons are calculated simultaneously, and the prices of objects that receive a bid are raised simultaneously. An alternative is a *Gauss-Seidel version* whereby a single unassigned person bids for an object, and the price rise of the object is taken into account when the next bid by an unassigned person takes place. This version is just as valid as the Jacobi version and in fact tends to converge a little faster, but is generally much less parallelizable.

Another variation of the auction algorithm applies to asymmetric assignment problems where the number of objects is larger than the number of persons, and there is a requirement that all persons be assigned. To solve this problem, the auction algorithm of this section need only be modified in the choice of the initial conditions. It is sufficient to require that the initial assignment be empty, and that all the initial object prices be zero. The algorithm terminates when all persons have been assigned, and it can be shown that the appropriate form of  $\epsilon$ -CS condition holds at termination. The details are left for the reader. Additional variations of the auction algorithm are discussed in [13].

### 3. PROPERTIES OF THE ALGORITHM

Suppose that the algorithm terminates with the final (complete) assignment  $\{j_i | i = 1, \dots, N\}$ , the object prices  $p_j$ , and the profit margins  $\pi_i$  given by (1). Then, by adding (4) over  $i$ ,

$$\sum_i a_{ij_i} \geq \sum_i (\pi_i + p_{j_i}) - N\epsilon$$

If  $A^*$  is the optimal primal value and the (equal) optimal dual value we have using the relation above

$$A^* \geq \sum_i a_{ij_i} \geq \sum_i (\pi_i + p_{j_i}) - N\epsilon = q(p) - N\epsilon \geq A^* - N\epsilon$$

Therefore the assignment  $\{j_i | i = 1, \dots, N\}$  is within  $N\epsilon$  of being optimal. Since  $a_{ij}$  are integer, an optimal assignment is obtained when  $\epsilon < 1/N$ . Thus we have shown:

**PROPOSITION 1:** *A complete assignment obtained upon termination of the rithm is within  $N\epsilon$  of being optimal, and is optimal if  $\epsilon < 1/N$ .*

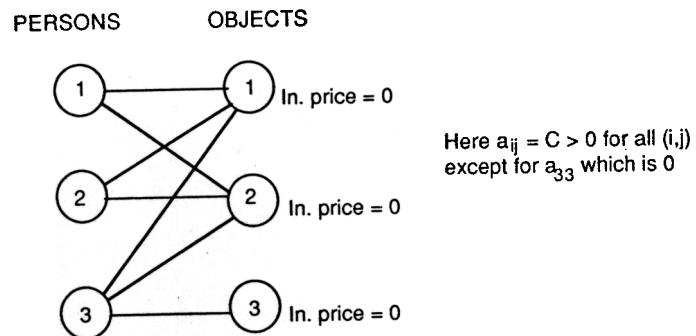
The next result asserts that the algorithm terminates assuming existence of at least one complete assignment. The proof relies on the following facts:

- a) Once an object is assigned, it remains assigned throughout the remainder of the algorithm's duration. Furthermore, except at termination, there will always exist at least one object that has never been assigned, and has a price equal to its initial price. This is due to the fact that a bidding and assignment phase can result in a reassignment of an already assigned object to a different person, but cannot result in the object becoming unassigned.

- b) Each time an object receives a bid its price increases by at least  $\epsilon$  [cf. (7), (8)]. Therefore if the object receives a bid an infinite number of times, its price increases to  $\infty$ .
- c) Each time a person  $i$  bids for an object a number of times at most equal to the cardinality of  $A(i)$ , his profit margin  $\pi_i$  as defined by (1) must decrease by at least  $\epsilon$ . This is due to the fact that a bid by person  $i$  either decreases  $\pi_i$  by at least  $\epsilon$ , or else leaves  $\pi_i$  unchanged because there is more than one object  $j$  attaining the maximum in (1). However, in the latter case the price of the object  $j^*$  receiving the bid will increase by at least  $\epsilon$ , and object  $j^*$  will not receive a bid again from person  $i$  until  $\pi_i$  decreases by at least  $\epsilon$ . The conclusion is that if a person bids an infinite number of times his profit margin must decrease to  $-\infty$ .

**PROPOSITION 2:** *If at least one complete assignment exists, the algorithm terminates in a finite number of iterations.*

**PROOF:** If the algorithm continues indefinitely, the prices of a proper [cf. a) above] subset  $J^\infty$  of objects increases to  $\infty$ , while the profit margins  $\pi_i$  of a subset  $I^\infty$  of persons decrease to  $-\infty$ , [cf. c) above]. Furthermore, eventually, in view of (4), at any given time each object in  $J^\infty$  can only be assigned to a person from  $I^\infty$ , and a person from  $I^\infty$  will either be assigned to an object in  $J^\infty$  or be unassigned. Also, in view of c) above, eventually only persons from  $I^\infty$  will be unassigned. Therefore the cardinality of  $I^\infty$  is greater than the cardinality of  $J^\infty$  while, in view of (1), we have  $J^\infty \supset A(i)$  for all  $i$  in  $I^\infty$ . This contradicts the existence of a complete assignment. **Q.E.D.**



**FIGURE 6:** Example where the number of bidding phases is proportional to  $C/\epsilon$ . Here at each bidding phase, one of the persons 1, 2, or 3 bids the price of object 1 or 2 up by an increment  $\epsilon$  until the time that these prices reach the level  $C$ .

It can be shown (see Fig. 6) that the computational complexity of the algorithm is sensitive to the maximal absolute object value

$$C = \max\{|a_{ij}| \mid i = 1, \dots, N, j \in A(i)\}$$

It was found experimentally that this sensitivity manifests itself often in the practical performance of the algorithm. To improve this performance, scaling the value of  $\epsilon$  suggests itself. Here we multiply all costs  $a_{ij}$  with  $(N + 1)$  and apply the algorithm with progressively lower value of  $\epsilon$  up to the point where  $\epsilon$  becomes 1 or smaller. The sequence of  $\epsilon$  values used is

$$\epsilon(k) = \max\{1, \Delta/\theta^k\}, \quad k = 0, 1, \dots$$

where  $\theta$  and  $\Delta$  are parameters with  $\Delta > 0$  and  $\theta > 1$ . Scaling procedures are discussed in more detail in [14], [30], [31]. The scaling procedure given above, called  $\epsilon$ -scaling, was mentioned in [6] and was tested by the author in 1979 and 1985. The computational results of Section 5 show that, with proper choice of the parameters  $\delta$  and  $\theta$ , the algorithm has excellent performance even without the benefit of parallelism.

#### 4. ASYNCHRONOUS IMPLEMENTATION

The auction algorithm given in Section 2 may be described as *synchronous* since the bidding and assignment are carried out simultaneously for all persons and objects respectively. A totally asynchronous (chaotic) algorithm (see [5], [10], [19], [20], [22] for computation models and convergence analysis) results if each unassigned person makes a bid at arbitrary times on the basis of object price information that may be outdated (because of additional bidding of which the person is not informed). Furthermore assignment of objects may be decided even if some potential bidders have not been heard from. We can similarly show the same termination properties as for the synchronous version of the algorithm subject to two conditions.

- a) An unassigned person will bid for some object within finite time, and cannot bid twice (*i.e.* cannot bid for a second object while waiting for a reply regarding the disposition of an earlier bid for another object).
- b) Whenever one or more bids are received that raise the price of an object then, within finite time, that price must be updated, and its value must be communicated (not necessarily simultaneously) to all persons. Furthermore the new bidder assigned to the object must be informed that he has been assigned simultaneously with receiving the new price.

A more careful formulation of the asynchronous model, and a proof of validity for the more general case of the minimum cost flow problem is given in [11], [12], and [14].

#### 5. COMPUTATIONAL RESULTS

In this section we provide the results of computational experimentation on both a serial machine (VAX 11-750), and a parallel machine (sequent Balance 21000). A more extensive computational study is currently under way.

A code called AUCTION was used to test the performance of the algorithm on a serial machine. The code implements  $\epsilon$ -scaling as described at the end of Section 3 for various values of the parameters  $\theta$  and  $\Delta$ . The initial object prices were set to  $p_j = \min_i a_{ij}$  for all  $j$ ; this is a common choice for dual type of assignment algorithms. At the end of the  $k$ th subproblem (*i.e.* the problem

corresponding to  $\epsilon(k-1)$ ) the assignment obtained is checked to see if condition (4) is satisfied for  $\epsilon = \epsilon(k)$ . If (4) is not satisfied for a person  $i$  and the corresponding object  $j_i$ , the pair  $(i, j_i)$  is removed from the assignment, so that condition (4) is enforced at the beginning of each subproblem.

The test problems were generated using the widely used public domain generator NETGEN [36] with the random seed number set at the value used in the original paper [36]. The AUCTION code was compared with a version of the relaxation code RELAX-II [16], [17] specially adapted to solve assignment problems. This code, called RELAX-IIA, outperformed a variety of specialized relaxation and primal-dual codes for the assignment problem based on existing methods. On the basis of this fact and in view of the excellent performance of the RELAX-II code for assignment problems, as documented in [16], [17], we believe that AUCTION has been given state-of-the-art competition.

The codes compared were written in standard FORTRAN on a VAX 11-750, and were compiled under VMS version 4.1. A rough indication of the potential speedup of AUCTION in a parallel machine is provided by the *degree of parallelism* defined as the number of unassigned persons in a bidding phase averaged over all bidding phases.

Our results are summarized in Figures 7-10. Figure 7 gives the solution times for five problems with equal density (1.5%). The parameters used were  $\theta = 8$  and  $\Delta = NC/2$ . Figure 8 gives the times for five problems where the set  $A(i)$  has average cardinality equal to 5. The parameters used were  $\theta = 6$  and  $\Delta = NC/2$ . Figure 9 gives the times for the five standard assignment problems (problems 11-15 in [36]). The parameters used were  $\theta = 4$  and  $\Delta = NC/2$ . The parallelism factor for these problems is roughly constant near 10 as Figure 10 shows. This indicates that the efficiency of the algorithm may be quite low when implemented on a parallel machine with a large number of processors.

The auction algorithm was also independently implemented and tested on dense assignment problems by Prof. J. Kennington and Mr. L. Hatay at Southern Methodist University using a Sequent Balance 21000 computer - a shared memory parallel machine. Figure 11 shows their results with a multiple processor implementation of the auction algorithm and, in particular, the speedup as a function of the number of CPUs employed. It is seen that the efficiency of the algorithm for a small number of processors is quite satisfactory. Note that the speedup with one processor is less than one because we are comparing a serial with a parallel code, and, even on a single processor, the parallel code is not as efficient as the serial code.

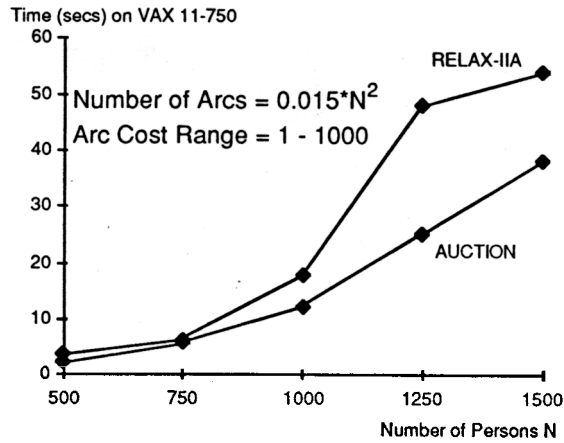


FIGURE 7: Solution times for AUCTION and RELAX-IIA on VAX 11-750. All problems generated by NETGEN.

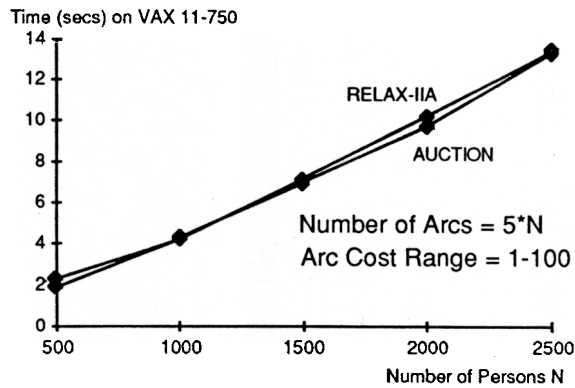


FIGURE 8: Solution Times for AUCTION and RELAX-IIA on VAX 11-750. All problems generated by NETGEN.

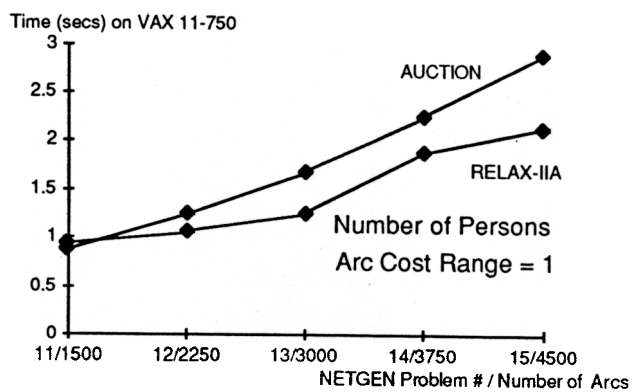


FIGURE 9: Solution times in secs for AUCTION and RELAX-IIA on VAX 11-750 for standard NETGEN assignment problems.

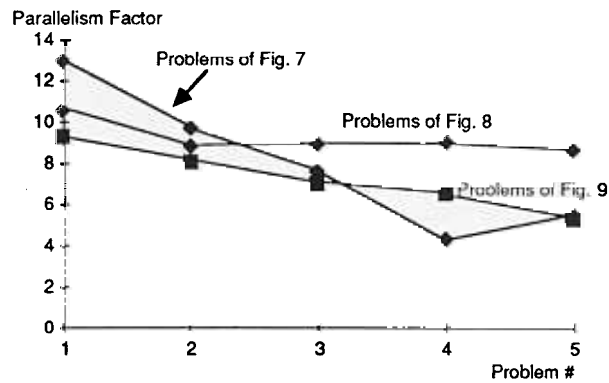


FIGURE 10: Parallelism factor for the problems of Figs. 7, 8, and 9.

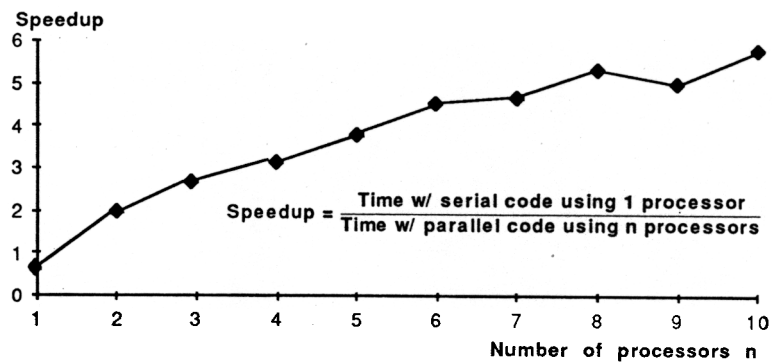


FIGURE 11: Speedup of a parallel implementation of the auction algorithm as a function of the number of processors used in a Sequent Balance 21000 computer. The problem solved is a randomly generated  $800 \times 800$  fully dense problem with arc cost range 1 - 1000. The time required by the serial auction code using a single processor is 336.13 secs.

#### REFERENCES

- [1] R. AHUJA and J.B. ORLIN, "An Improved Algorithm for Computing the Min Cycle Mean", Unpublished Manuscript, Nov. 1987.
- [2] M.L. BALINSKI, Signature methods for the assignment problem, *Oper. Res.* 33 (1985) 527-537.
- [3] M.L. BALINSKI, A competitive (dual) simplex method for the assignment problem, *Math. Progr.* 34 (1986) 125-141.
- [4] R. BARR, F. GLOVER and D. KLINGMAN, The alternating basis algorithm for assignment problems, *Math. progr.* 13 (1977) 1-13.
- [5] G.M. BAUDET, Asynchronous iterative methods for multiprocessors, *J. ACM* 15 (1978) 226-244.
- [6] D.P. BERTSEKAS, A distributed algorithm for the assignment problem, *Laboratory for Information and Decision Systems Working Paper* (M.I.T., March 1979).



- [7] D.P. BERTSEKAS, A new algorithm for the assignment problem, *Math. Progr.* 21 (1981) 152-171.
- [8] D.P. BERTSEKAS, A Distributed Asynchronous Relaxation Algorithm for the Assignment Problem, *Proc. 24th IEEE Conf. Decision and Control* (Ft Lauderdale, Fla., Dec. 1985) pp. 1703-1704.
- [9] D.P. BERTSEKAS, A unified framework for primal-dual methods in minimum cost network flow problems, *Math. Progr.* 32 (1985) 125-145.
- [10] D.P. BERTSEKAS, Distributed asynchronous computation of fixed points, *Math. Progr.* 27 (1983) 107-120.
- [11] D.P. BERTSEKAS, *Distributed Asynchronous Relaxation Methods for Linear Network Flow Problems*, LIDS Report P-1606 (M.I.T., Sept. 1986, revised Nov. 1986).
- [12] D.P. BERTSEKAS, Distributed relaxation methods for linear network flow problems, in: *Proc. 25th IEEE Conf. Decision and Control, Athens, Greece* (1986) pp. 2101-2106.
- [13] D.P. BERTSEKAS and D. CASTANON, The Auction Algorithm for Transportation Problems, unpublished manuscript (Nov. 1987).
- [14] D.P. BERTSEKAS and J. ECKSTEIN, Distributed asynchronous relaxation methods for linear network flows problems, *Proc. IFAC '87*, Munich, Germany (July 1987).
- [15] D.P. BERTSEKAS and D. EL BAZ, Distributed asynchronous relaxation methods for convex network flow problems, *SIAM J. Control Optim.* 25 (1987) 74-85.
- [16] D.P. BERTSEKAS and P. TSENG, Relaxation methods for minimum cost ordinary and generalized network flow problems, LIDS Report P-1462, M.I.T., May 1985, *Oper. Res.* to appear.
- [17] D.P. BERTSEKAS and P. TSENG, RELAX: A Code for Linear Network Flow Problems, in *FORTTRAN Codes for Network Optimization*, (B. Simeone, ed.), Vol to be published by the *Annals of Operations Research*.
- [18] D.P. BERTSEKAS, P.A. HOSEIN and P. TSENG, Relaxation methods for network flow problems, *SIAM J. Control Optim.* 25 (1987) 1219-1243.
- [19] D.P. BERTSEKAS, and J.N. TSITSIKLIS, *Parallel and Distributed Computation: Numerical Methods* (Prentice-Hall, Englewood Cliffs, N.J., 1988) to appear.
- [20] D.P. BERTSEKAS, J.N. TSITSIKLIS and M. ATHANS, *Convergence Theories of Distributed Asynchronous Computation: A Survey*, LIDS Report P-1412, M.I.T., Oct. 1984; also in : *Stochastic programming*, eds. F. ARCHETTI, G. DI PILLO, and M. LUCERTINI, (Springer-Verlag, N.Y., 1986) pp. 107-120.
- [21] R.G. BLAND and D.L. JENSEN, *On the Computational Behaviour of a Polynomial-Time Network Flow Algorithm*, Tech. Report 661 (School of Operations Research and Industrial Engineering, Cornell University, June 1985).
- [22] D. CHAZAN and W. MIRANKER, Chaotic relaxation, *Linear Algebra Appl.* 2 (1969) 199-222.
- [23] V.P. CRAWFORD and E.M. KNOER, Job matching with heterogeneous firms and workers, *Econometrica* 49 (1981) 437-450.

- [24] E. DINIC and M. KRONROD, An algorithm for the solution of the assignment problem, *Soviet Math. Dokl.* 10 (1969).
- [25] J. EDMONDS and R.M. KARP, Theoretical improvements in algorithmic efficiency for network flow problems, *J. ACM* 19 (1972) 248-264.
- [26] M. ENGQUIST, A successive shortest path algorithm for the assignment problem, *INFOR* 20 (1982) 370-384.
- [27] H.N. GABOW and R.E. TARJAN, Faster scaling algorithms for graph matching, Revised Abstract, 1987.
- [28] F. GLOVER, R. GLOVER and D. KLINGMAN, *Threshold Assignment Algorithm*, Center for Business Decision Analysis report CBDA 107 (Graduate School of Business, Univ. of Texas at Austin, Sept. 1982).
- [29] A.V. GOLDBERG, Solving minimum-cost flow problems by successive approximations, extended abstract, submitted to *STOC 87* (Nov. 1986).
- [30] A.V. GOLDBERG, *Efficient Graph Algorithms for Sequential and Parallel Computers*, Tech. Report TR-374, Laboratory for Computer Science, M.I.T., Feb. 1987.
- [31] A.V. GOLDBERG, and R.E. TARJAN, Solving minimum cost flow problems by successive approximation, in: *Proc. 19th ACM STOC*, May 1987.
- [32] D. GOLDFARB, Efficient dual simplex methods for the assignment problem, *Math. Progr.* 33 (1985) 187-203.
- [33] M. HUNG, A polynomial simplex method for the assignment problem, *Oper. Res.* 31 (1983) 595-600.
- [34] M. HUNG and W. ROM, Solving the assignment problem by relaxation, *Oper. Res.* 28 (1980) 969-982.
- [35] J. KENNINGTON and R. HELGASON, *Algorithms for Network Programming* (Wiley, New York, 1980).
- [36] D. KLINGMAN, A. NAPIER and J. STUTZ, NETGEN - A program for generating large scale (un)capacitated assignment, transportation, and minimum cost flow network problems, *Management Sci.* 20 (1974) 814-822.
- [37] H.W. KUHN, The Hungarian method for the assignment problem, *Naval Res. Logistics Quarter.* 2 (1955) 83-97.
- [38] L.F. MCGINNIS, Implementation and testing of a primal-dual algorithm for the assignment problem, *Oper. Res.* 31 (1983) 277-291.
- [39] J.B. ORLIN, *Genuinely Polynomial Simplex and Non-Simplex Algorithms for the Minimum Cost Flow Problem*, Working Paper No. 1615-84 (Sloan School of Management, M.I.T., Dec. 1984).
- [40] C.H. PAPADIMITRIOU and K. STEIGLITZ, *Combinatorial Optimization: Algorithms and Complexity* (Prentice-Hall, Englewood Cliffs, N.J. 1982).
- [41] H. ROCK, Scaling techniques for minimal cost network flows, in: *Discrete Structures and Algorithms*, ed. V. PAGE Carl Hansen, Munich, 1980.
- [42] R.T. ROCKAFELLAR, *Network Flows and Monotropic programming* (J. Wiley, New York, 1984).
- [43] A.E. ROTH, The economics of matching: stability and incentives, *Math. Oper. Res.* 7 (1982) 617-628.
- [44] L.S. SHAPLEY and M. SHUBIK, The assignment game I: the core, *Intern. J. Game Theory* 1 (1972) 117-130.

- [45] E. TARDOS, A strongly polynomial minimum cost circulation algorithm, *Combinatorica* 5 (1985) 247-255.
- [46] A. WEINTRAUB and F. BARAHONA, *A Dual Algorithm for the Assignment Problem*, Publ. No. 79/02/C (Departamento de Industrias, Universidad de Chile-Sede Occidente, Santiago, Chile, 1979).
- [47] S.A. ZENIOS and J.M. MULVEY, Relaxation techniques for strictly convex network problems, *Ann. Oper. Res.* 5 (1985/86), 517-538.
- [48] S.A. ZENIOS and R.A. LASKEN, *Nonlinear Network Optimization on a Massively Parallel Connection Machine*, Report 87-08-03 (Decision Science Department, The Wharton School, Univ. of Pennsylvania, Aug. 1987).

NOTE ADDED IN PROOF: An improved version of the AUCTION code has been placed in the public domain and can be obtained from the author at no cost.