

An efficient cost scaling algorithm for the assignment problem

Andrew V. Goldberg ^{*,1}, Robert Kennedy ²

Computer Science Department, Stanford University, Stanford, CA 94305-2140, United States

Received 18 October 1993; revised manuscript received 19 January 1995

Abstract

The cost scaling push-relabel method has been shown to be efficient for solving minimum-cost flow problems. In this paper we apply the method to the assignment problem and investigate implementations of the method that take advantage of assignment's special structure. The results show that the method is very promising for practical use.

Keywords: Network optimization; Assignment problem; Algorithms; Experimental evaluation; Cost scaling

1. Introduction

Significant progress has been made in the last decade on the theory of algorithms for network flow problems. Some of the algorithms that came out of this research have been shown to have practical impact as well. In particular, the push-relabel method [11, 16] is the best currently known way for solving the maximum flow problem [2, 8, 23]. This method extends to the minimum-cost flow problem using cost scaling [11, 17]. Earlier implementations of this method [5, 14] performed well on some problems but relatively poorly on others. A later implementation [12] has been shown very competitive on a wide class of problems. In this paper we study efficient implementations of the cost scaling push-relabel method for the assignment problem.

* Corresponding author. Present address: NEC Research Institute, 4 Independence Way, Princeton, NJ 08540, United States, e-mail: avg@research.nj.nec.com.

¹ This author's research was supported in part by ONR Young Investigator Award N00014-91-J-1855, NSF Presidential Young Investigator Grant CCR-8858097 with matching funds from AT&T, DEC and 3M, and a grant from the Powell Foundation.

² This author's research was supported by the above-mentioned ONR and NSF grants.

The most relevant theoretical results on the assignment problem are as follows. The best currently known strongly polynomial time bound of $O(n(m + n \log n))$ is achieved by the classical Hungarian method of Kuhn [21]. Here n denotes the number of nodes in the input network and m denotes the number of edges. (Implementations of the Hungarian method and related algorithms are discussed in [7].) Under the assumption that the input costs are integers in the range $[-C, \dots, C]$, Gabow and Tarjan [10] use cost scaling and blocking flow techniques to obtain an $O(\sqrt{n} m \log(nC))$ time algorithm. Algorithms with the same running time bound based on the push-relabel method appeared in [15, 24].

In this paper we study implementations of the scaling push-relabel method in the context of the assignment problem. We use the ideas behind the minimum-cost flow codes [5, 12, 14], the assignment codes [3, 4, 6], and the ideas of theoretical work on the push-relabel method for the assignment problem [15], as well as new techniques aimed at improving practical performance of the method. We develop several CSA (Cost Scaling Assignment) codes based on different heuristics which improve the code performance on many problem classes. The “basic” code CSA-B does not use any heuristics, the CSA-Q code uses a “quick-minima” heuristic, and the CSA-S code uses a “speculative arc fixing” heuristic. Another outcome of our research is a better understanding of cost scaling algorithm implementations, which may lead in turn to improved cost scaling codes for the minimum-cost flow problem.

We compare the performance of the CSA codes to four other codes: SFR10, an implementation of the auction method for the assignment problem [6]; SJV and DJV, implementations of Jonker and Volgenant’s shortest augmenting path method [19] tuned for sparse and dense graphs respectively; and ADP/A, an implementation of the interior-point method specialized for the assignment problem [25]. We make the comparison over classes of problems from generators developed for the First DIMACS Implementation Challenge [18]³ and on problems obtained from digital images as suggested by Knuth [20]. Of our codes, CSA-Q is best overall. This code outperforms ADP/A on all problem instances in our tests, outperforms SFR10 on all except one class, and outperforms SJV and DJV on large instances in every class. Although our second-best code, CSA-S, is somewhat slower than CSA-Q on many problem classes, it is usually not much slower and it outperforms CSA-Q on three problem classes, always outperforms ADP/A, is worse than SFR10 by only a slight margin on one problem class and by a noticeable margin on only one problem class, and loses to the Jonker–Volgenant codes only on one class and on small instances from two other classes. While we use the CSA-B code primarily to gauge the effect of heuristics on performance, it is worth noting that it outperforms ADP/A in all our tests, the Jonker–Volgenant codes on large instances from all but one class, and SFR10 on all but one class of problems we tested.

This paper is organized as follows. Section 2 gives the relevant definitions. Section 3 outlines the scaling push-relabel method for the assignment problem. Section 4

³The DIMACS benchmark codes, problem generators, and other information we refer to are available by anonymous ftp from dimacs.rutgers.edu.

discusses our implementation, in particular the techniques used to improve our code's practical performance. Section 5 describes the experimental setup. Section 6 gives the experimental results. In Section 7, we give concluding remarks.

2. Background

Let $\bar{G} = (\bar{V} = X \cup Y, \bar{E})$ be an undirected bipartite graph and let $n = |\bar{V}|$ and $m = |\bar{E}|$. A *matching* in \bar{G} is a subset of edges $M \subseteq \bar{E}$ that have no node in common. The *cardinality* of the matching is $|M|$. Given a weight function $\bar{c} : \bar{E} \rightarrow \mathbb{R}$, we define the weight of M to be the sum of weights of edges in M . The *assignment problem* is to find a maximum cardinality matching of maximum weight. We assume that the weights are integers in the range $[-C, \dots, C]$. To simplify the presentation, we assume that $|X| = |Y|$, \bar{G} has a perfect matching (i.e., a matching of cardinality $|X|$), and every node degree in \bar{G} is at least two. We can dispense with these last assumptions without any significant decrease in performance by using a slightly more complicated reduction to the transportation problem than the one described below.

Our implementation reduces the assignment problem to the *transportation problem* defined as follows. Let $G = (V, E)$ be a digraph with a real-valued *capacity* $u(a)$ and a real-valued *cost* $c(a)$ associated with each arc⁴ a and a real-valued *supply* $d(v)$ associated with each node v . We assume that $\sum_v d(v) = 0$. A *pseudoflow* is a function $f : E \rightarrow \mathbb{R}_+$ satisfying the *capacity* constraints: for each $a \in E$, $f(a) \leq u(a)$. For a pseudoflow f and a node v , the *excess flow into* v , $e_f(v)$, is defined by $e_f(v) = d(v) + \sum_{(u,v) \in E} f(u,v) - \sum_{(v,w) \in E} f(v,w)$. A node v with $e_f(v) > 0$ is called *active*. Note that $\sum_{v \in V} e_f(v) = 0$.

A *flow* is a pseudoflow f such that, for each node v , $e_f(v) = 0$. Observe that a pseudoflow f is a flow if and only if there are no active nodes. The *cost* of a pseudoflow f is given by $\text{cost}(f) = \sum_{a \in E} c(a)f(a)$. The transportation problem is to find a flow of minimum cost.

We use a slight variation of the standard reduction from the assignment problem to the minimum-cost flow problem (see, e.g., [22]). Given an instance of the assignment problem (\bar{G}, \bar{c}) , we construct a transportation problem instance $(G = (V, E), c, u)$ as follows. We define $V = \bar{V} = X \cup Y$. For every edge $\{v, w\} \in \bar{E}$ such that $v \in X$ and $w \in Y$, we add the arc (v, w) to E and define $c(v, w) = -\bar{c}(v, w)$ and $u(v, w) = 1$. Finally we define $d(v) = 1$ for all $v \in X$ and $d(w) = -1$ for all $w \in Y$. Note that the graph G is bipartite.

For a given pseudoflow f , the *residual capacity* of an arc $a \in E$ is $u_f(a) = u(a) - f(a)$. The set of *residual arcs* E_f contains the arcs $a \in E$ with $f(a) < u(a)$ and the reverse arcs a^R , for every $a \in E$ with $f(a) > 0$. The *residual graph* $G_f = (V, E_f)$ is the graph induced by the residual arcs. For $a \in E$, we define $c(a^R) = -c(a)$. Note that if

⁴ Sometimes we refer to an arc a by its end points, e.g., (v, w) . This is ambiguous if there are multiple arcs from v to w . An alternative is to refer to v as the tail of a and to w as the head of a , which is precise but inconvenient.

```

procedure MIN-COST( $\forall E, u, c$ );
  [initialization]
   $\epsilon \leftarrow C$ ;  $\forall v, p(v) \leftarrow 0$ ;
  [loop]
  while  $\epsilon \geq 1/n$  do
     $(\epsilon, f, p) \leftarrow \text{REFINE}(\epsilon, p)$ ;
  return( $f$ );
end.

```

Fig. 1. The cost scaling algorithm.

G is obtained by the above reduction, then for any integral pseudoflow f and for any arc $a \in E$, $u(a), f(a) \in \{0, 1\}$.

A price function is a function $p : V \rightarrow \mathbb{R}$. For a given price function p , the reduced cost of an arc (v, w) is $c_p(v, w) = c(v, w) + p(v) - p(w)$ and the partial reduced cost is $c'_p(v, w) = c(v, w) - p(w)$.

For a constant $\epsilon \geq 0$, a pseudoflow f is said to be ϵ -optimal with respect to a price function p if, for every residual arc $a \in E_f$, we have

$$a \in E \Rightarrow c_p(a) \geq 0, \quad a^R \in E \Rightarrow c_p(a) \geq -\epsilon.$$

A pseudoflow f is ϵ -optimal if f is ϵ -optimal with respect to some price function p . If the arc costs and capacities are integers and $\epsilon < 1/n$, any ϵ -optimal flow is optimal.

For a given f and p , an arc $a \in E_f$ is admissible iff

$$a \in E \text{ and } c_p(a) < \frac{1}{2}\epsilon \quad \text{or} \quad a^R \in E \text{ and } c_p(a) < -\frac{1}{2}\epsilon.$$

The admissible graph $G_A = (V, E_A)$ is the graph induced by the admissible arcs.

3. The method

First we give a high-level description of the successive approximation algorithm (see Fig. 1). For a detailed presentation of the successive approximation framework and the associated proofs, see [17]. The algorithm starts with $\epsilon = C$ and $p(v) = 0$ for all $v \in V$. At the beginning of every iteration, the algorithm divides ϵ by a constant factor α and sets f to the zero pseudoflow. The iteration modifies f and p so that f is an (ϵ/α) -optimal flow with respect to p . When $\epsilon < 1/n$, f is optimal and the algorithm terminates. The number of iterations of the algorithm is $1 + \lfloor \log_\alpha(nC) \rfloor$.

Reducing ϵ is the task of the subroutine *refine*. The input to *refine* is ϵ and p such that there exists a flow f that is ϵ -optimal with respect to p . The output from *refine* is $\epsilon' = \epsilon/\alpha$, a flow f , and a price function p such that f is ϵ' -optimal with respect to p .

The generic *refine* subroutine (described in Fig. 2) begins by decreasing the value of ϵ , setting f to the zero pseudoflow (thus creating some excesses and making some nodes active), and setting $p(v) = -\min_{(v,w) \in E} \{c'_p(v, w)\}$ for every $v \in X$. This converts the f into an ϵ -optimal pseudoflow (indeed, into a 0-optimal pseudoflow). Then the subroutine converts f into an ϵ -optimal flow by applying a sequence of *push* and *relabel*

```

procedure REFINE( $\epsilon, p$ );
  [initialization]
   $\epsilon \leftarrow \epsilon/\alpha$ ;
   $\forall a \in E, f(a) \leftarrow 0$ ;
   $\forall v \in X, p(v) \leftarrow -\min_{(v,w) \in E} c'_p(v, w)$ ;
  [loop]
  while  $f$  is not a flow
    apply a push or a relabel operation;
  return( $\epsilon, f, p$ );
end.

```

Fig. 2. The generic *refine* subroutine.

operations, each of which preserves ϵ -optimality. The generic algorithm does not specify the order in which these operations are applied. Next, we describe the *push* and *relabel* operations for the unit-capacity case (see Fig. 3).

A *push* operation applies to an admissible arc (v, w) whose tail node v is active. It consists of pushing one unit of flow from v to w , thereby decreasing $e_f(v)$ by one, increasing $e_f(w)$, and either increasing $f(v, w)$ by one if $(v, w) \in E$ or decreasing $f(w, v)$ by one if $(w, v) \in E$. A *relabel* operation applies to a node v . The operation sets $p(v)$ to the smallest value allowed by the ϵ -optimality constraints, namely $\max_{(v,w) \in E_f} \{p(w) - c(v, w)\}$ if $v \in X$, or $\max_{(v,w) \in E_f} \{p(w) - c(v, w) - \epsilon\}$ otherwise.

The analysis of cost scaling push-relabel algorithms is based on the following facts [15, 17]. During a scaling iteration,

- the node prices monotonically decrease;
- for any $v \in V$, $p(v)$ decreases by $O(n\epsilon)$.

4. Implementation and heuristics

In this section we discuss implementation issues and heuristics aimed at speeding up the method.

The efficiency of a scaling implementation depends on the choice of scale factor α . Although an earlier study [6] suggests that the performance of scaling codes for the assignment problem may be quite sensitive to the choice of scale factor, our observations are to the contrary. Spot checks on instances from several problem classes indicated that

```

PUSH( $v, w$ ).
  send a unit of flow from  $v$  to  $w$ .
end.

RELABEL( $v$ ).
  if  $v \in X$ 
    then replace  $p(v)$  by  $\max_{(v,w) \in E_f} \{p(w) - c(v, w)\}$ 
    else replace  $p(v)$  by  $\max_{(u,v) \in E_f} \{p(u) + c(u, v) - \epsilon\}$ 
  end.

```

Fig. 3. The *push* and *relabel* operations.

the running times seem to vary by a factor of no more than 2 for values of α between 4 and 40. We chose $\alpha = 10$ for our tests; different values of α would yield running times that are somewhat worse on some problem classes and somewhat better on others, but the difference is not drastic. We believe the lack of robustness alluded to in [6] may be due to a characteristic of the implementation of SFR10 and related codes. In particular, SFR10 contains an “optimization” that seems to terminate early scaling phases prematurely. Our codes run every scaling phase to completion as suggested by the theory.

The efficiency of an implementation of *refine* depends on the number of operations performed by the method and on the implementation details. We discuss the operation ordering first.

The implementation maintains the price function p and the flow f . For each node $w \in Y$ with $e_f(w) = 0$, we maintain a pointer to the unique node $v = \mu(w)$ such that $f(v, w) = 1$.

Our implementation maintains the invariant that only the nodes in X are active, except possibly in the middle of the double-push operation described below. The implementation picks an active node and applies the double-push operation to it.

The performance of the implementation depends on the strategy for selecting the next active node to process. We experimented with several operation orderings, including those suggested in [13, 17]. Our implementation uses the LIFO ordering, i.e., the set of active nodes is maintained as a stack. This ordering worked best in our tests; the FIFO ordering usually worked somewhat worse, although the difference was never drastic.

4.1. The double-push operation

The *double-push* operation is similar to a sequential version of the match-and-push procedure from [15]. This operation was independently discovered in [1]. The operation applies to an active node v . Recall that at the beginning of a double-push, all active nodes are in X , so $v \in X$.

First the double-push operation processes v by relabeling v , pushing flow from v along an admissible arc (v, w) , and then relabeling v again. If $e_f(w)$ becomes positive, the operation pushes flow from w to $\mu(w)$ and sets $\mu(w) = v$. Finally, double-push relabels w .

Lemma 4.1. *The double-push operation is correct.*

Proof. We only need to show that double-push applies the pushing operation correctly. Since immediately before the flow is pushed out of v the node is relabeled, there is an admissible arc out of v and the push is correct. If this push makes w active, then there is a second push from w to $\mu(w)$.

Consider the last double-push into w which set $\mu(w)$ to its current value. Because the network is obtained via a reduction described in Section 2, $(w, \mu(w))$ is the only residual arc out of w . So when the double-push relabeled w , $c_p(\mu(w), w)$ became

ϵ . From this double-push to the current one, w and $\mu(w)$ have not been relabeled; (the latter holds because $(w, \mu(w))$ was the only residual arc into $\mu(w)$ during that time period). Thus during the current push from w , $c_p(\mu(w), w) = \epsilon$, so the push is valid. \square

Lemma 4.2. *A double-push operation decreases the price of a node $w \in Y$ by at least ϵ .*

Proof. Just before the double-push, w is either unmatched or matched.

In the first case, the flow is pushed into w and at this point the only residual arc out of w is the arc (w, v) . Just before that the double-push relabeled v and $c_p(v, w) = 0$. Next, double-push relabels w and $p(w)$ decreases by ϵ .

In the second case, the flow is pushed to w and at this point w has two outgoing residual arcs, (w, v) and $(w, \mu(w))$. As we have seen, $c_p(v, w) = 0$ and $c_p(\mu(w), w) = \epsilon$. After the second relabeling of v , double-push pushes flow from w to $\mu(w)$ and relabels w , reducing $p(w)$ by ϵ . \square

Corollary 4.3. *There are $O(n^2)$ double-push operations per refine.*

4.2. Efficient implementation

Suppose we apply double-push to a node v . Let (v, w) and (v, z) be the arcs out of v with the smallest and the second-smallest reduced costs, respectively. These arcs can be found by scanning the adjacency list of v once. The effects of double-push on v are equivalent to pushing flow along (v, w) and setting $p(v) = -c'_p(v, z)$. To relabel w , we set $p(w) = p(v) + c(v, w) - \epsilon$. This implementation of double-push is summarized in Fig. 4.

It is not necessary to maintain the prices of nodes in X explicitly; for $v \in X$, we can define $p(v)$ implicitly by $p(v) = -\min_{(v,w) \in E} \{c'_p(v, w)\}$ if $e_f(v) = 1$ and $p(v) = c'(v, w) + \epsilon$ if $e_f(v) = 0$ and (v, w) is the unique arc with $f(v, w) = 1$. One can easily verify that using implicit prices is equivalent to using explicit prices in the above implementation. The only time we need to know the value of $p(v)$ is when we relabel w in double-push, and at that time $p(v) = -c'_p(v, z)$ which we compute during the previous relabel of v . Maintaining the prices implicitly saves memory and time. The

```

DOUBLE-PUSH( $v$ ).
  let  $(v, w)$  and  $(v, z)$  be the arcs with the smallest and the second-smallest reduced costs;
  PUSH( $v, w$ );
   $p(v) = -c'_p(v, z)$ ;
  if  $e_f(w) > 0$ 
    PUSH( $w, \mu(w)$ );
   $\mu(w) = v$ ;
   $p(w) = p(v) + c(v, w) - \epsilon$ ;
end.

```

Fig. 4. Efficient implementation of double-push.

implementation of the double-push operation with implicit prices is similar to the basic step of the auction algorithm of [3].

Our code CSA-B implements the scaling push-relabel algorithm using stack ordering of active nodes and the implementation of double-push with implicit prices mentioned above.

4.3. Heuristics

In this section we describe two heuristics that often improve the algorithm's performance.

The k th-best heuristic [3] is aimed at reducing the number of scans of arc lists of nodes in X . The idea of the k th-best heuristic is as follows. Recall that we scan the list of v to find the arcs (v, w) and (v, z) with the smallest and second-smallest values of the partial reduced cost. Let $k \geq 3$ be an integer. When we scan the list of $v \in X$, we compute the k th-smallest value K of the partial reduced costs of the outgoing arcs and store the $k - 1$ arcs with the $k - 1$ smallest partial reduced costs. The node prices monotonically decrease during refine, hence during the subsequent double-push operations we can first look for the smallest and the second-smallest arcs among the stored arcs whose current partial reduced cost is at most K . We need to scan the list of v again only when all except possibly one of the saved arcs have partial reduced costs greater than K .

Our code CSA-Q is a variation of CSA-B that uses the fourth-best heuristic.

The idea of the *speculative arc fixing* heuristic [9, 12] is to move arcs with reduced costs of large magnitude to a special list. These arcs are not examined by the double-push procedure but are examined as follows at a (relatively large) periodic interval. When the arc (v, w) is examined, if the ϵ -optimality condition is violated on (v, w) , $f(v, w)$ is modified to restore ϵ -optimality and (v, w) is moved back to the adjacency list of v ; if ϵ -optimality holds for (v, w) but $|c_p(v, w)|$ is no longer large, (v, w) is simply moved back to the adjacency list. This heuristic takes advantage of the fact that the flow is *fixed* on arcs of high reduced cost [17].

Our code CSA-S is a variation of CSA-B that uses the speculative arc fixing heuristic.

We implemented a number of other heuristics that are known to improve performance of cost scaling code for the minimum-cost flow problem [12]. Among these are: *global price updates* which periodically ensure, via a specialized shortest-paths computation, that the admissible graph contains a path from every node with flow excess to some node with flow deficit; and *price refinement* which determines at each iteration whether the current assignment is actually ϵ' -optimal for some $\epsilon' < \epsilon$, and hence avoids unnecessary executions of *refine*. Our best implementation uses neither of these strategies, however, since even taking advantage of the assignment problem's structure to simplify and speed up these heuristics, a typical price refinement iteration used more time than simply executing *refine* in our tests. The double-push operation seems to maintain a sufficiently "aggressive" price function and global price updates cannot reduce the number of *push* and *relabel* operations enough to improve the running time.

Table 1
DIMACS benchmark times

C benchmarks		FORTRAN benchmarks	
Test 1	Test 2	Test 1	Test 2
2.7 sec	24.0 sec	1.2 sec	2.2 sec

5. Experimental setup

All the test runs were executed on a Sun SparcStation 2 with a clock rate of 40 MHz and 96 Megabytes of main memory. We compiled the SFR10 code supplied by Castañon with the Sun Fortran-77 compiler, release 2.0.1 using the `-O4` optimization switch.⁵ We compiled the DJV and SJV codes supplied by Hao with the Sun C compiler release 1.0, using the `-O2` optimization option. We compiled our CSA codes with the Sun C compiler release 1.0, using the `-fast` optimization option; each choice seemed to yield the fastest execution times for the code where we used it. Times reported here are UNIX *user* CPU-times, and were measured using the `times()` library function. During each run, the programs collect time usage information after reading the input problem and initializing all data structures and again after computing the optimum assignment; we take the difference between the two figures to indicate the CPU-time actually spent solving the assignment problem.

To give a baseline for comparison of our machine's speed to others, we ran the DIMACS benchmarks `wmatch` (to benchmark C performance) and `netflo` (to benchmark FORTRAN performance) on our machines, with the timing results given in Table 1. It is interesting (though neither surprising nor critical to our conclusions) to note that the DIMACS benchmarks do not precisely reflect the mix of operations in the codes we developed. Of two C compilers available on our system, the one that consistently ran our code faster by a few percent also ran the benchmarks more slowly by a few percent; (the C benchmark times in Table 1 are for code generated by the same compiler we used for our experiments). But even though they should not be taken as the basis for very precise comparison, the benchmarks provide a useful way to estimate relative speeds of different machines on the sort of operations typically performed by combinatorial optimization codes.

We did not run the ADP/A code on our machine, but because the benchmark times reported in [25] differ only slightly from the times we obtained on our machine, we conclude that the running times reported for ADP/A in [25] form a reasonable basis for comparison with our codes. Therefore, we report running times directly from [25]. As the reader will see, even if this benchmark comparison introduces a significant amount of error, our conclusions about the codes' relative performance are justified by the large differences in performance between ADP/A and the other codes we tested.

⁵ Castañon [6] recommends setting the initial "bidding increment" in SFR10 to a special value for problems of high density; we found this advice appropriate for the dense problem class, but discovered that it hurt performance on the geometric class. We followed Castañon's recommendation only on the class where it seemed to improve SFR10's performance.

The DJV code is designed for dense problems and uses an adjacency-matrix data structure. The memory requirements for this code would be prohibitive on sparse problems with many nodes. For this reason, we included it only in experiments on problem classes that are dense. On these problems, DJV is faster than SJV by a factor of about 1.5. It is likely that our codes and the SFR10 code would enjoy a similar improvement in performance if they were modified to use the adjacency-matrix data structure.

We collected performance data on a variety of problem classes, many of which we took from the First DIMACS Implementation Challenge. Following is a brief description of each class; details of the generator inputs that produced each set of instances are included in Appendix A.

5.1. *The high-cost class*

Each $v \in X$ is connected by an edge to $2 \log_2 |V|$ randomly-selected nodes of Y , with integer edge costs uniformly distributed in the interval $[0, 10^8]$.

5.2. *The low-cost class*

Each $v \in X$ is connected by an edge to $2 \log_2 |V|$ randomly-selected nodes of Y , with integer edge costs uniformly distributed in the interval $[0, 100]$.

5.3. *The two-cost class*

Each $v \in X$ is connected by an edge to $2 \log_2 |V|$ randomly-selected nodes of Y , each edge having cost 100 with probability $\frac{1}{2}$, or cost 10^8 with probability $\frac{1}{2}$.

5.4. *The fixed-cost class*

For problems in this class, we view X as a copy of the set $\{1, 2, \dots, \frac{1}{2}|V|\}$, and Y as a copy of $\{\frac{1}{2}|V| + 1, \frac{1}{2}|V| + 2, \dots, |V|\}$. Each $v \in X$ is connected by an edge to $\frac{1}{16}|V|$ randomly-selected nodes of Y , with edge (x, y) , if present, having cost $100xy$.

5.5. *The geometric class*

Geometric problems are generated by placing a collection of integer-coordinate points uniformly at random in the square $[0, 10^6] \times [0, 10^6]$, coloring half the points blue and the other half red, and introducing an edge between every red point r and every blue point b with cost equal to the floor of the distance between r and b .

5.6. *The dense class*

Like instances of the geometric class, dense problems are complete, but edge costs are distributed uniformly at random in the range $[0, 10^7]$.

5.7. Picture problems

Picture problems, suggested by Knuth [20], are generated from photographs scanned at various resolutions, with 256 greyscale values. The set V is the set of pixels; the pixel at row r , column c is a member of X if $r + c$ is odd, and is a member of Y otherwise. Each pixel has edges to its vertical and horizontal neighbors in the image, and the cost of each edge is the absolute value of the greyscale difference between its two end points. Note that picture problems are extremely sparse, with an average degree always below 4. Although picture problems are an abstract construct with no practical motivation, the solution to a picture problem can be viewed as a tiling of the picture with dominos, where we would like each domino to cover greyscale values that are as different as possible.

For our problems, we used two scanned photographs, one of each author of this paper.

6. Experimental observations and discussion

In the following tables and graphs, we present performance data for the codes. Note that problem instances are characterized by the number of nodes on a single side, i.e., *half* the number of nodes in the graph.

We report times on the test runs we conducted, along with performance data for the ADP/A code taken from [25]. The instances on which ADP/A was timed in [25] are identical to those we used in our tests. We give mean running times computed over three instances for each problem size in each class; in the two-cost and geometric classes we also give mean running times computed over fifteen instances and sample deviations for each sample size. We computed sample deviations for each problem class and size, and observed that in most cases they were less than ten percent of the mean (often much less). The two exceptions were the two-cost and geometric classes, where we observed larger sample deviations in the running times for some of the codes. For these two classes we also collected data on fifteen instances for each problem size. The sample statistics taken over fifteen instances seem to validate those we observed for three instances. All statistics are reported in seconds.

6.1. The high-cost class

Fig. 5 and Table 2 summarize the timings on DIMACS high-cost instances. The k th-best heuristic yields a clear advantage in running time on these instances. CSA-Q beats CSA-B, its nearest competitor, by a factor of nearly 2 on large instances, and CSA-Q seems to have an asymptotic advantage over the other codes, as well. The overhead of speculative arc fixing is too great on high-cost instances; the running times of CSA-S for large graphs are essentially the same as those of SFR10. SJV has the worst asymptotic behavior.

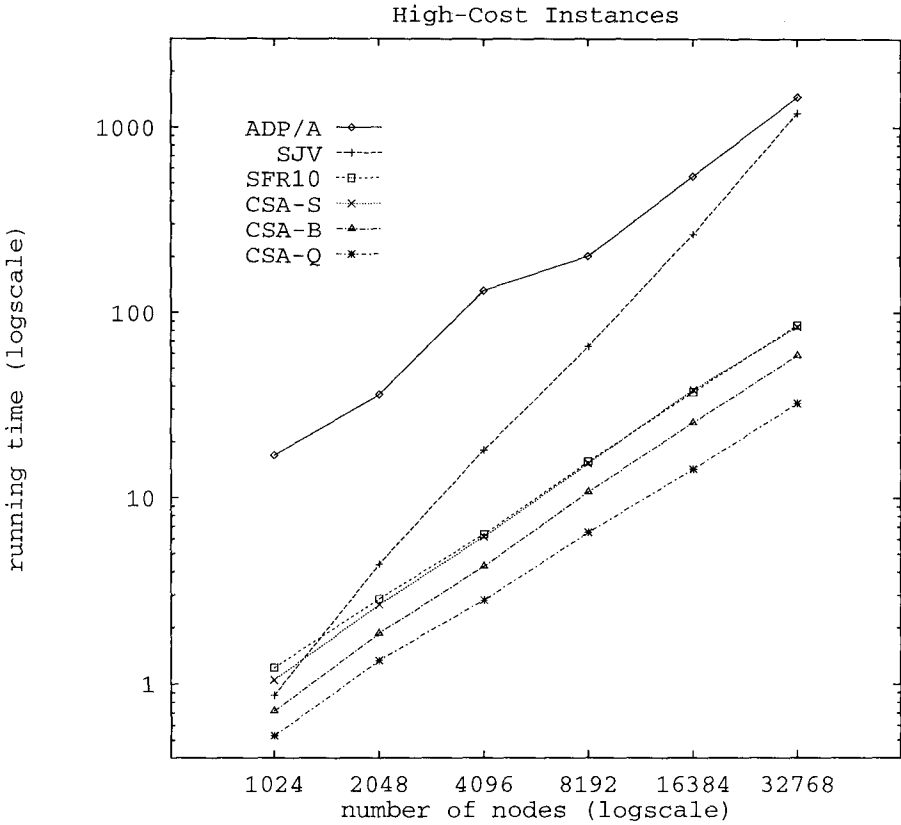


Fig. 5. Running times for the high-cost class.

Table 2
Running times for the high-cost class

Nodes X	ADP/A	SFR10	SJV	CSA-B	CSA-S	CSA-Q
1024	17	1.2	0.87	0.7	1.1	0.5
2048	36	2.9	4.40	1.9	2.7	1.3
4096	132	6.4	18.1	4.3	6.2	2.8
8192	202	15.7	65.6	10.8	15.3	6.5
16384	545	37.3	266	25.5	38.3	14.3
32768	1463	85.7	1197	58.7	84.0	32.4

6.2. The low-cost class

The situation here is very similar to the high-cost case: CSA-Q enjoys a slight asymptotic advantage as well as a clear constant-factor advantage over the competing codes. SJV has worse asymptotic behavior than the other codes on the low-cost class, just as it does on high-cost instances. See Fig. 6 and Table 3.

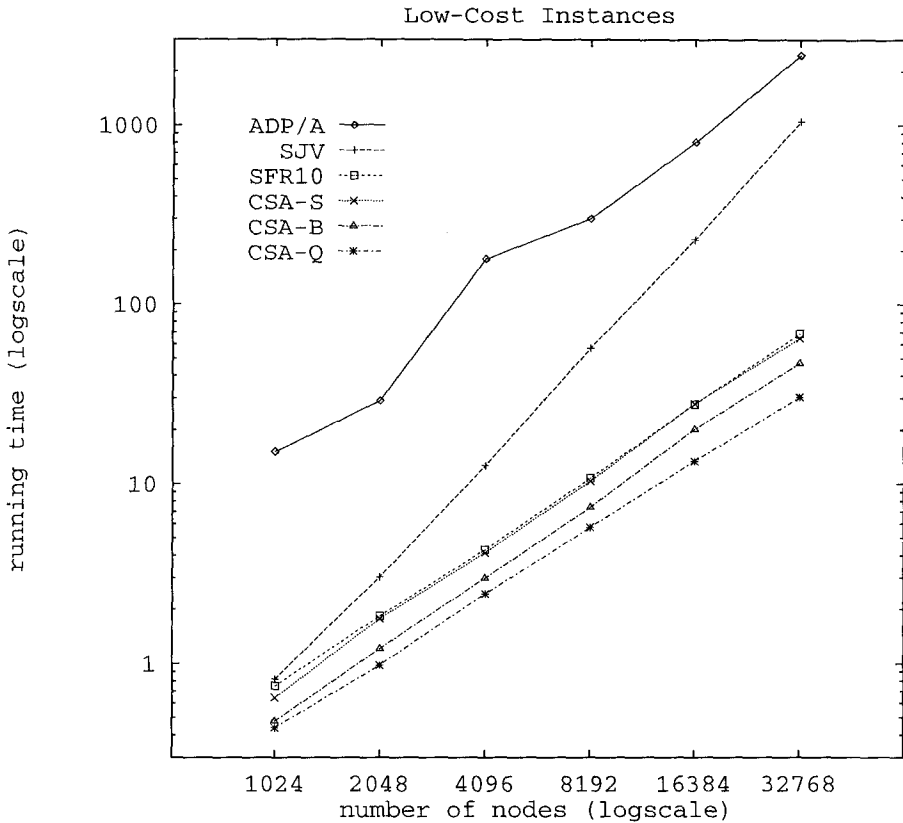


Fig. 6. Running times for the low-cost class.

Table 3

Running times for the low-cost class

Nodes $ X $	ADP/A	SFR10	SJV	CSA-B	CSA-S	CSA-Q
1024	15	0.75	0.82	0.48	0.64	0.44
2048	29	1.83	3.03	1.21	1.77	0.98
4096	178	4.31	12.6	2.99	4.13	2.43
8192	301	10.7	57.0	7.39	10.3	5.72
16384	803	27.7	229	20.1	27.8	13.4
32768	2464	68.5	1052	46.9	64.6	30.3

6.3. The two-cost class

The two-cost data appear in Fig. 7 and Tables 4 and 5. It is difficult for robust scaling algorithms to exploit the special structure of two-cost instances; the assignment problem for most of the graphs in this class amounts to finding a perfect matching on the high-cost edges, and none of the scaling codes we tested is able to take special advantage of this observation. Because SJV does not use scaling, it would seem a good

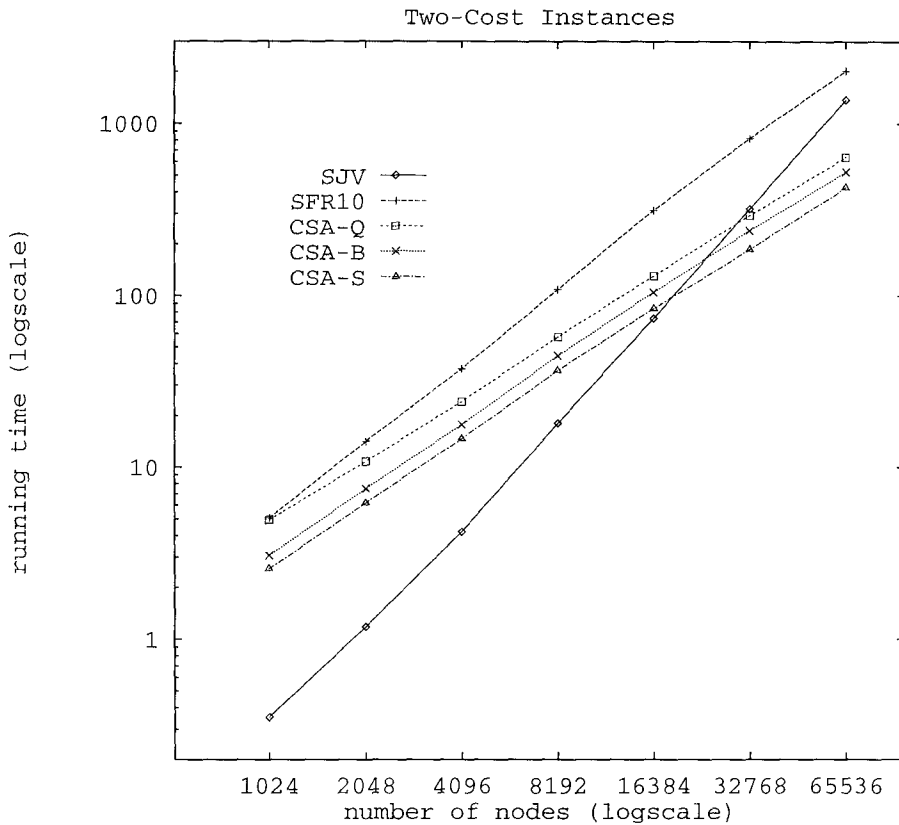


Fig. 7. Running times (three-instance samples) for the two-cost class.

Table 4
Running times (three-instance samples) for the two-cost class

Nodes $ X $	SFR10		SJV		CSA-B		CSA-S		CSA-Q	
	time	s	time	s	time	s	time	s	time	s
1024	5.13	0.09	0.35	0.00	3.09	0.24	2.58	0.15	5.21	0.33
2048	14.0	1.1	1.16	0.01	7.72	0.28	6.19	0.18	11.1	1.0
4096	37.3	1.1	4.21	0.16	17.7	1.1	14.2	1.6	23.3	2.4
8192	107	12	18.2	0.43	43.4	3.5	36.7	2.1	58.6	3.3
16384	366	81	73.6	0.58	102	2.8	85.4	3.2	133	7.8
32768	894	180	320	1.2	240	6.0	185	6.8	299	6.4
65536	1782	60	1370	5.8	531	15	417	11	628	25

candidate to perform especially well on this class, and indeed it does well on small two-cost instances. For large instances, however, SJV uses a great deal of time in its shortest augmenting path phase, and performs poorly for this reason. Speculative arc fixing improves significantly upon the performance of the basic CSA implementation, and the k th-best heuristic hurts performance on this class of problems. It seems that the

Table 5
Running times (fifteen-instance samples) for the two-cost class

Nodes $ X $	SFR10		SJV		CSA-B		CSA-S		CSA-Q	
	time	s	time	s	time	s	time	s	time	s
1024	5.05	0.32	0.35	0.02	3.07	0.21	2.56	0.10	4.93	0.40
2048	14.1	1.1	1.18	0.04	7.49	0.37	6.18	0.26	10.8	0.73
4096	37.4	2.7	4.22	0.14	17.7	1.0	14.7	0.84	24.1	1.6
8192	109	9.8	18.0	0.37	44.6	2.5	36.5	1.5	57.5	3.2
16384	314	50	73.7	0.57	105	4.2	84.1	2.9	130	8.4
32768	822	194	320	2.1	239	8.5	186	4.8	293	15
65536	2021	342	1376	7.5	524	25	426	16	637	27

k th-best heuristic tends to speed up the last few iterations of *refine*, but it hurts in the early iterations. Like k th-best, the speculative arc fixing heuristic is able to capitalize on the fact that later iterations of *refine* can afford to ignore many of the arcs incident to each node, but by keeping all arcs of similar cost under consideration in the beginning, speculative arc fixing allows early iterations to run relatively fast. On this class, CSA-S is the winner, although for applications limited to this sort of strongly bimodal cost distribution, an unscaled push-relabel or blocking flow algorithm might perform better than any of the codes we tested. No running times are given in [25] for ADP/A on this problem class, but the authors suggest that their program performs very well on two-cost problems. Relative to those of the other codes, the running times of SFR10 are comparatively scattered at each problem size in this class; we believe this phenomenon results from the premature termination of early scaling phases in SFR10 (see Section 4).

The relatively large sample deviations shown in Fig. 7 and Table 4 motivated our experiments with fifteen instances of each problem size. The sample means and deviations of the fifteen-instance data are shown in Table 5, and they are consistent with and very similar to the three-instance data shown in Fig. 7 and Table 4.

6.4. The fixed-cost class

Fig. 8 and Table 6 give the data for the fixed-cost problem class. On smaller instances of this class, CSA-B and CSA-Q have nearly the same performance. On instances with $|X| = 1024$ and $|X| = 2048$, CSA-Q is faster on fixed-cost problems than CSA-B, or indeed any of the other codes. On smaller instances, speculative arc fixing does not pay for itself; when $|X| = 2048$, the overhead is just paid for. Perhaps on larger instances, speculative arc fixing would pay off. It is doubtful, though, that CSA-S would beat CSA-Q on any instances of reasonable size. SJV exhibits the worst asymptotic behavior among the codes we tested on this problem class.

6.5. The geometric class

On geometric problems, both heuristics improve performance over the basic CSA-B code. The performance of CSA-S and CSA-Q is similar to and better than that of the

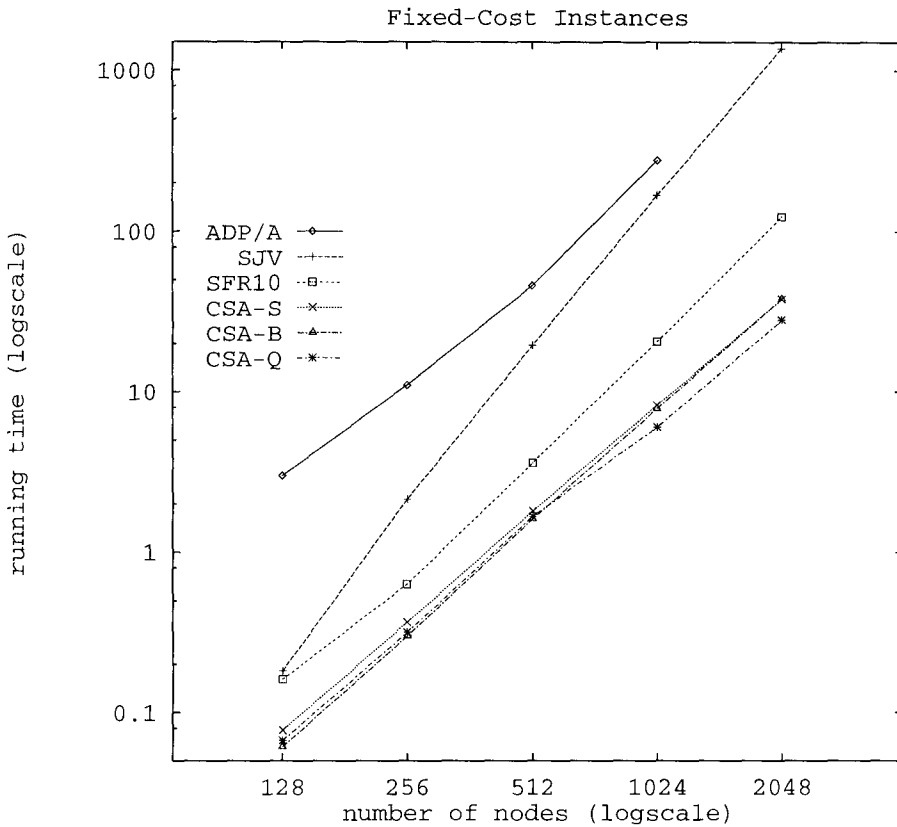


Fig. 8. Running times for the fixed-cost class.

Table 6
Running times for the fixed-cost class

Nodes $ X $	ADP/A	SFR10	SJV	CSA-B	CSA-S	CSA-Q
128	3	0.16	0.18	0.06	0.08	0.07
256	11	0.63	2.14	0.30	0.37	0.32
512	46	3.59	19.4	1.6	1.8	1.7
1024	276	20.5	168	7.8	8.2	6.0
2048	n.a.	123	1367	37.8	37.6	27.9

other codes. The Jonker-Volgenant codes seem to have asymptotic behavior similar to the other codes on this class. See Fig. 9 and Table 7.

Because the sample deviations shown in Fig. 9 and Table 7 are somewhat large compared to those we observed on most other problem classes, we ran experiments on fifteen instances as a check on the validity of the data. Statistics calculated over fifteen-instance samples are reported in Table 8, and they are very much like the three-instance data.

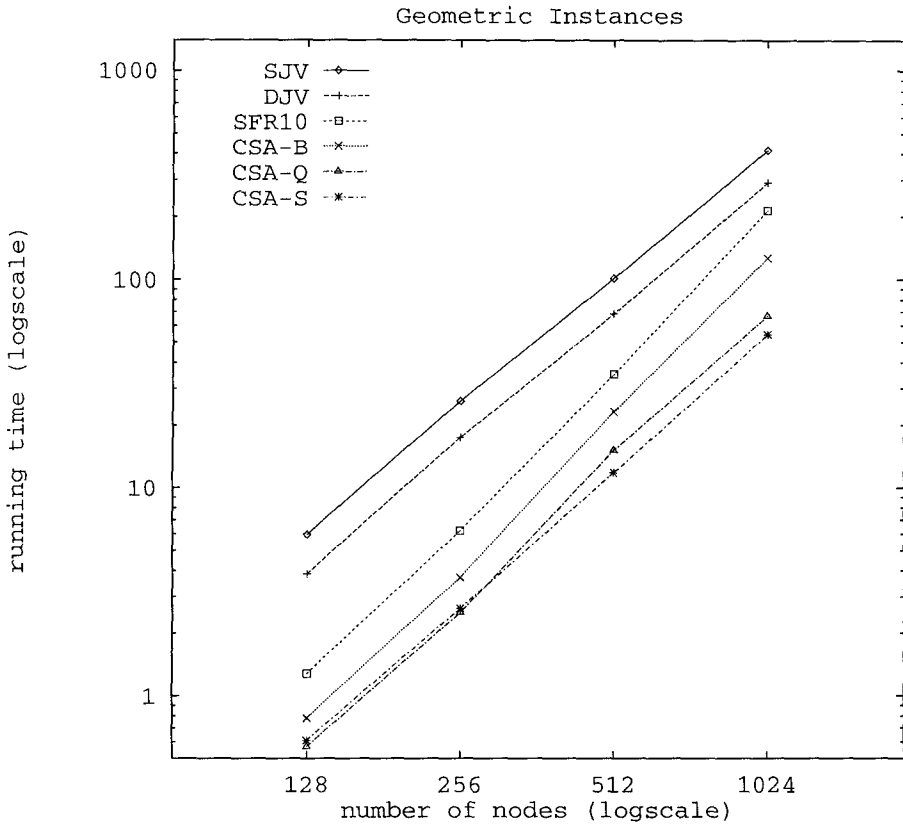


Fig. 9. Running times (three-instance samples) for the geometric class.

Table 7
Running times (three-instance samples) for the geometric class

Nodes $ X $	ADP/A		SFR10		SJV		DJV		CSA-B		CSA-S		CSA-Q	
	time	s	time	s	time	s	time	s	time	s	time	s	time	s
128	12	0.5	1.27	0.46	6.64	4.4	4.36	2.9	0.79	0.28	0.62	0.05	0.58	0.19
256	47	1	6.12	0.23	25.3	3.3	16.9	2.0	3.67	0.67	2.56	0.08	2.43	0.34
512	214	42	31.0	4.1	110	2.8	73.2	1.0	27.9	8.1	11.9	0.89	16.7	3.7
1024	1316	288	193	19	424	51	297	32	114	24	54.9	1.42	62.5	2.6

Table 8
Running times (fifteen-instance samples) for the geometric class

Nodes $ X $	SFR10		SJV		DJV		CSA-B		CSA-S		CSA-Q	
	time	s	time	s	time	s	time	s	time	s	time	s
128	1.28	0.21	5.96	2.0	3.85	1.3	0.78	0.16	0.61	0.03	0.57	0.11
256	6.21	0.82	26.1	4.7	17.5	2.9	3.72	0.51	2.63	0.09	2.50	0.27
512	35.0	6.0	101	11	68.2	7.4	23.2	4.9	11.8	0.67	15.1	2.4
1024	214	54	416	38	291	25	127	27	54.4	2.2	66.7	9.7

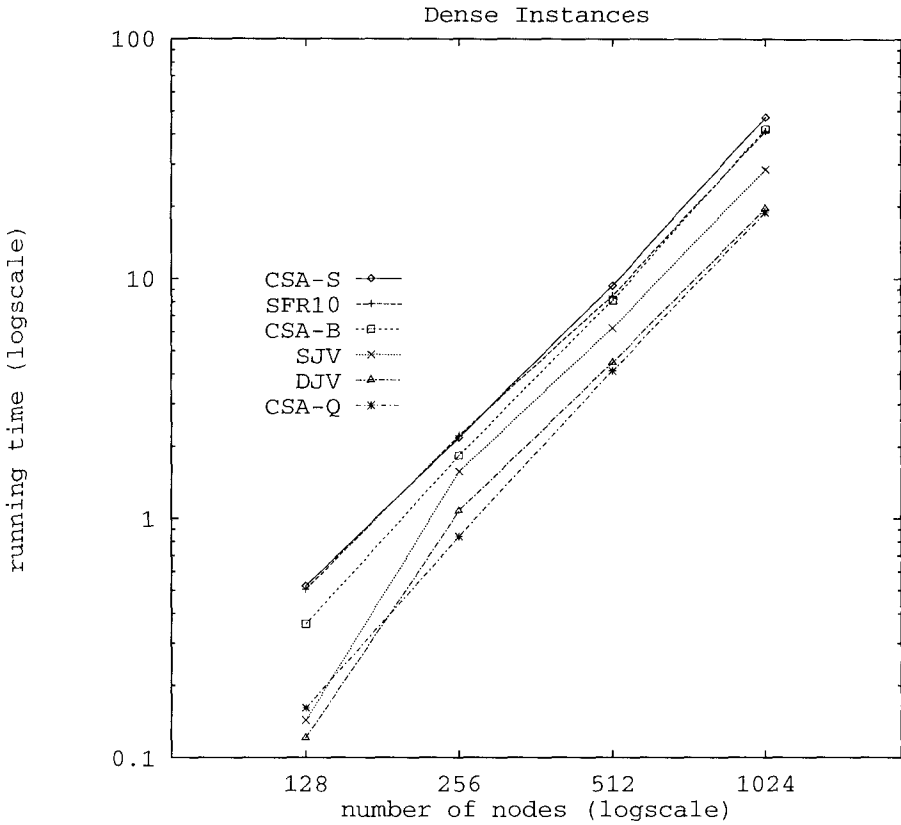


Fig. 10. Running times for the dense class.

Table 9
Running times for the dense class

Nodes $ X $	SFR10	SJV	DJV	CSA-B	CSA-S	CSA-Q
128	0.51	0.14	0.12	0.36	0.52	0.16
256	2.22	1.57	1.07	1.83	2.17	0.84
512	8.50	6.22	4.47	8.12	9.36	4.13
1024	41.2	28.5	19.6	42.0	47.1	18.9

6.6. The dense class

The difference between Fig. 10 and Tables 8 and 9 shows that the codes' relative performance is significantly affected by changes in cost distribution. Except on very small instances, CSA-Q is the winner in this class; DJV is its closest competitor, with SJV performing fairly well also. As in the case of geometric problems, SJV and DJV seem to have asymptotic performance similar to the scaling and interior-point codes on this class.

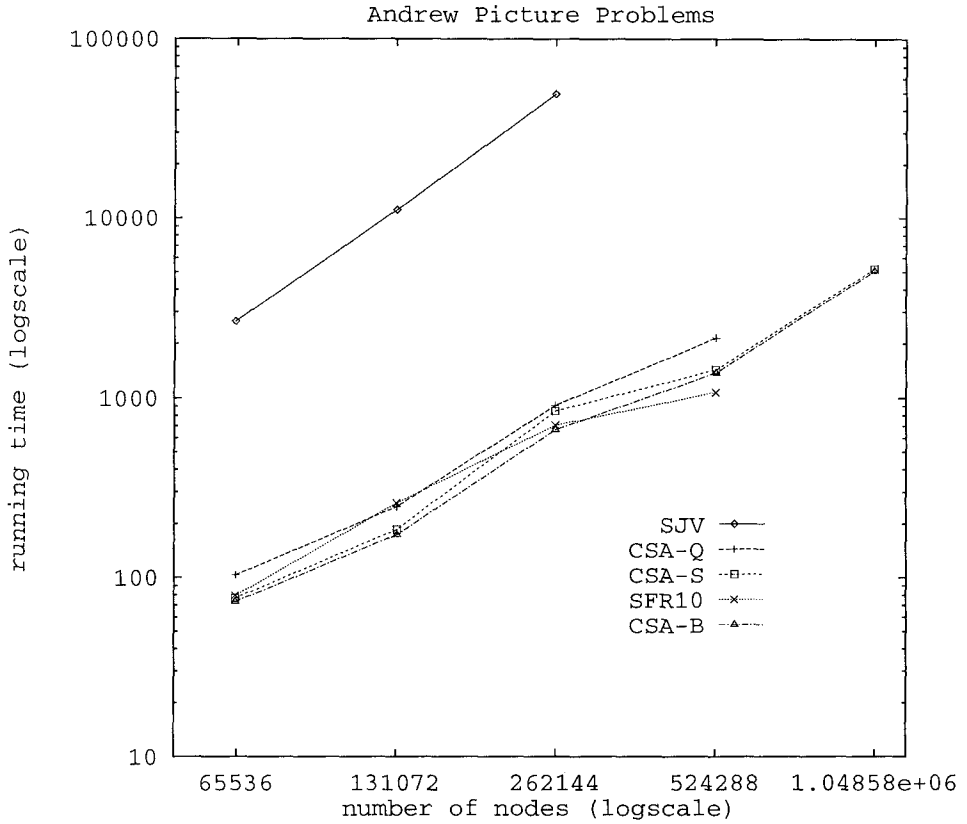


Fig. 11. Running times for problems from Andrew's picture.

Table 10
Running times for problems from Andrew's picture

Nodes X	SFR10	SJV	CSA-B	CSA-Q	CSA-S
65 158	79.20	2656	73.23	103.3	76.70
131 370	260.2	11 115	173.2	248.0	185.5
261 324	705.2	49 137	665.1	907.8	844.8
526 008	1073	n.a.	1375	2146	1432
1046 520	n.a.	n.a.	5061	n.a.	5204

Table 11
Running times for problems from Robert's picture

Nodes X	SFR10	SJV	CSA-B	CSA-Q	CSA-S
59 318	49.17	1580	50.13	68.10	51.82
119 132	153.1	6767	154.8	223.6	165.4
237 272	351.4	26 637	585.0	916.8	611.2
515 088	827.8	n.a.	2019	3095	3057
950 152	1865	n.a.	5764	n.a.	8215

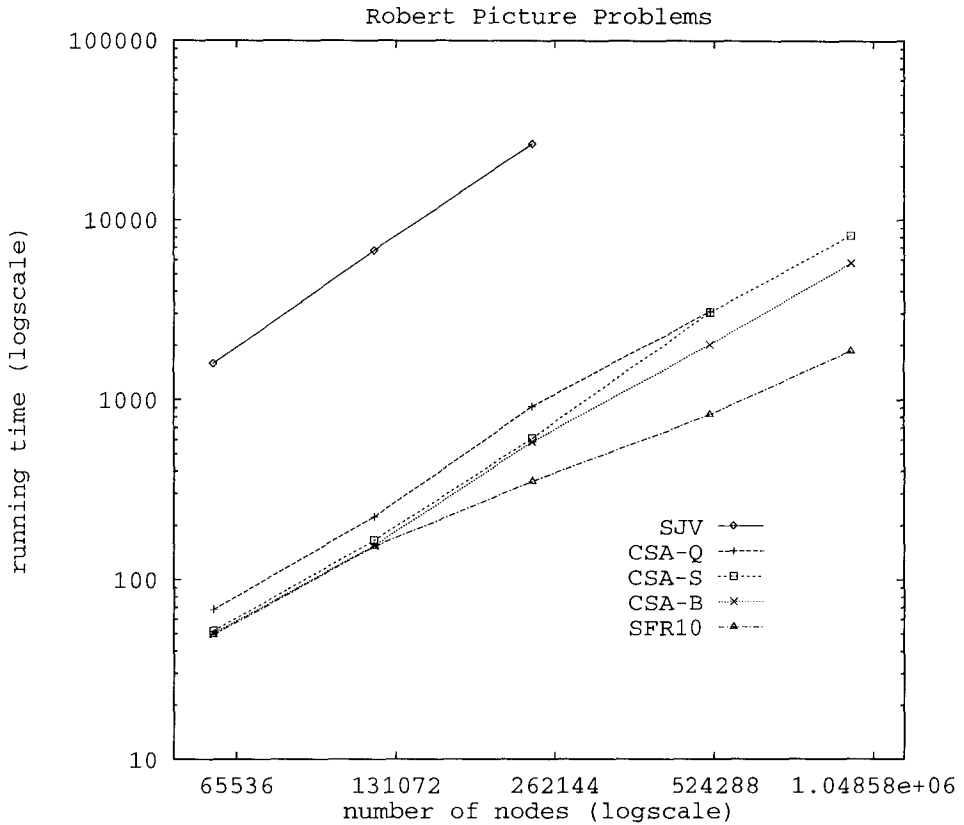


Fig. 12. Running times for problems from Robert's picture.

6.7. Picture problems

Although the pictures used had very similar characteristics, the tentative conclusions we draw here about the relative performance of the codes seem to apply to a broader class of images. We performed trials on a variety of images generated and transformed by various techniques, and found no substantial differences in relative performance, although some pictures seem to yield more difficult assignment problems than others. See Figs. 11 and 12 and Tables 10 and 11. On the picture problems we tried, SFR10 performs better than any of the CSA implementations; we believe that the “reverse-auction” phases performed by SFR10 [6] are critical to this performance difference. We were unable to obtain times for SJV and CSA-Q on the largest problem instance from each picture, nor from SFR10 on the largest problem instance from one of the pictures because the codes required too much memory. On the second-largest instance from each picture, our experiments suggested that SJV would require more than a day of CPU-time, so we did not collect data for these cases. On picture problems CSA-Q

performs significantly worse than either of the other two CSA implementations. This situation is no surprise because CSA-Q performs an additional pointer dereference each time it examines an arc. In such a sparse graph, the four arcs stored at each node exhaust the list of arcs incident to that node, so no benefit is to be had from the k th-best heuristic.

7. Concluding remarks

Castañon [6] gives running times for an auction code called SF5 in addition to performance data for SFR10; SF5 and SFR10 are the fastest among the robust codes discussed. The data in [6] show that on several classes of problems, SF5 outperforms SFR10 by a noticeable margin. Comparing Castañon's reported running times for SFR10 with the data we obtained for the same code allows us to estimate roughly how SF5 performs relative to our codes. The data indicate that CSA-S and CSA-Q should perform at least as well as SF5 on all classes for which data are available, and that CSA-Q should outperform SF5 by a wide margin on some classes. A possible source of error in this technique of estimation is that Castañon reports times for test runs on cost-minimization problems, whereas all the codes we test here (including SFR10) are configured to maximize cost. The difference in every case is but a single line of code, but while on some classes minimization and maximization problems are similar, on other classes we observed that minimization problems were significantly easier for all the codes. This difference is unlikely to be a large error source, however, since the relative performance of the codes we tested was very similar for minimization problems and maximization problems.

It is interesting that SJV is asymptotically worse than all its competitors on every sparse class, and that SJV and DJV are asymptotically very similar to their competitors on the dense classes. DJV performs very well on the uniform dense problem class, but we feel SJV provides a more genuine reference point, since the other combinatorial codes could be sped up on dense problems by replacing their central data structures with an adjacency matrix representation similar to that in DJV.

From our tests and data from [25] and [6], we conclude that CSA-Q is a robust, competitive implementation that should be considered for use by those who wish to solve assignment problems in practice.

Acknowledgements

The authors would like to thank David Castañon for supplying and assisting with the SFR10 code, Anil Kamath and K.G. Ramakrishnan for their assistance in interpreting results reported in [25], Jianxiu Hao for supplying and assisting with the SJV and DJV implementations, and Serge Plotkin for his help producing the digital pictures.

The second author would like to thank Scott Burson for providing computing facilities during the development of the codes.

Appendix A. Generator inputs

The assignment instances on which we ran our tests were generated as follows. Problems in the high-cost, low-cost, fixed-cost and dense classes were generated using the DIMACS generator `assign.c`. Problems in the two-cost class were generated using `assign.c` with output post-processed by the DIMACS awk script `twocost.a`. Problems in the geometric class were generated using the DIMACS generator `dcube.c` with output post-processed by the DIMACS awk script `geomasn.a`. Picture problems were generated from images in the Portable Grey Map format using our program `p5pgmtoasn`. To obtain the DIMACS generators, use anonymous ftp to `dimacs.rutgers.edu`, or obtain the `csa` package (which includes the generators) as described below.

In each class except the picture class, we generated instances of various numbers of nodes N and used various seeds K for the random number generator. For each problem type and each N , either three or fifteen values of K were used; the values were integers 270 001 through 270 003 or through 270 015. For picture problems, we tested the codes on a single instance of each size.

A.1. The high-cost class

We generated high-cost problems using `assign.c` from the DIMACS distribution. The input parameters given to the generator are as follows, with the appropriate values substituted for N and K :

```
nodes  $N$ 
sources  $\frac{1}{2}N$ 
degree  $2 \log_2 N$ 
maxcost 100 000 000
seed  $K$ 
```

A.2. The low-cost class

Like high-cost problems, low-cost problems are generated using the DIMACS generator `assign.c`. The parameters to the generator are identical to those for high-cost problems, except for the maximum edge cost:

```
nodes  $N$ 
sources  $\frac{1}{2}N$ 
degree  $2 \log_2 N$ 
maxcost 100
seed  $K$ 
```

A.3. *The two-cost class*

Two-cost instances are derived from low-cost instances using the UNIX awk program and the DIMACS awk script `twocost.a`. The instance with N nodes and seed K was generated using the following UNIX command line, with input parameters identical to those for the low-cost problem class:

```
assign | awk -f twocost.a
```

A.4. *The fixed-cost class*

We generated fixed-cost instances using `assign.c`, with input parameters as follows:

```
nodes  $N$ 
sources  $\frac{1}{2}N$ 
degree  $\frac{1}{16}N$ 
maxcost 100
multiple
seed  $K$ 
```

A.5. *The geometric class*

We generated geometric problems using the DIMACS generator `dcube.c` and the DIMACS awk script `geomasn.a`. We gave input parameters to `dcube` as shown below, and used the following UNIX command line:

```
dcube | awk -f geomasn.a
nodes  $N$ 
dimension 2
maxloc 1000 000
seed  $K$ 
```

A.6. *The dense class*

We generated dense problems using `assign.c`, with input parameters as follows:

```
nodes  $N$ 
sources  $\frac{1}{2}N$ 
complete
maxcost 1000 000
seed  $K$ 
```

Appendix B. Obtaining the CSA codes

To obtain a copy of the CSA codes, DIMACS generators referred to in this paper, and documentation files, send mail to `ftp-request@theory.stanford.edu` and use send

csas.tar as the subject line; you will automatically be mailed a uuencoded copy of a tar file.

References

- [1] R.K. Ahuja, J.B. Orlin, C. Stein and R.E. Tarjan, “Improved algorithms for bipartite network flow,” *SIAM Journal on Computing* 23 (1994) 906–933.
- [2] R.J. Anderson and J.C. Setubal, “Goldberg’s algorithm for the maximum flow in perspective: a computational study,” in: D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993) pp. 1–18.
- [3] D.P. Bertsekas, “The auction algorithm: a distributed relaxation method for the assignment problem,” *Annals of Operations Research* 14 (1988) 105–123.
- [4] D.P. Bertsekas, *Linear Network Optimization: Algorithms and Codes* (MIT Press, Cambridge, MA, 1991).
- [5] R.G. Bland, J. Cheriyan, D.L. Jensen and L. Ladañyi, “An empirical study of min cost flow algorithms,” in: D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993) pp. 119–156.
- [6] D.A. Castañón, “Reverse auction algorithms for the assignment problems,” in: D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993) pp. 407–430.
- [7] U. Derigs, “The shortest augmenting path method for solving assignment problems—motivation and computational experience,” *Annals of Operations Research* 4 (1985–1986) 57–102.
- [8] U. Derigs and W. Meier, “Implementing Goldberg’s max-flow algorithm—a computational investigation,” *Zeitschrift für Operations Research* 33 (1989) 383–403.
- [9] S. Fujishige, K. Iwano, J. Nakano and S. Tezuka, “A speculative contraction method for the minimum cost flows: toward a practical algorithm,” in: D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993) pp. 219–246.
- [10] H.N. Gabow and R.E. Tarjan, “Faster scaling algorithms for network problems,” *SIAM Journal on Computing* 18 (1989) 1013–1036.
- [11] A.V. Goldberg, “Efficient graph algorithms for sequential and parallel computers,” Ph.D. Thesis, Massachusetts Institute of Technology, Cambridge, MA (1987); also: Technical Report TR-374, Laboratory for Computer Science, Massachusetts Institute of Technology, Cambridge, MA (1987).
- [12] A.V. Goldberg, “An efficient implementation of a scaling minimum-cost flow algorithm,” in: E. Balas and J. Clausen, eds., *Proceedings of the Third Integer Programming and Combinatorial Optimization Conference* (Springer, Berlin, 1993) pp. 251–266.
- [13] A.V. Goldberg and R. Kennedy, “Global price updates help,” Technical Report STAN-CS-94-1509, Department of Computer Science, Stanford University, CA (1994).
- [14] A.V. Goldberg and M. Kharitonov, “On implementing scaling push-relabel algorithms for the minimum-cost flow problem,” in: D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993) pp. 157–198.
- [15] A.V. Goldberg, S.A. Plotkin and P.M. Vaidya, “Sublinear-time parallel algorithms for matching and related problems,” *Journal of Algorithms* 14 (1993) 180–213.
- [16] A.V. Goldberg and R.E. Tarjan, “A new approach to the maximum flow problem,” *Journal of the Association for Computing Machinery* 35 (1988) 921–940.
- [17] A.V. Goldberg and R.E. Tarjan, “Finding minimum-cost circulations by successive approximation,” *Mathematics of Operations Research* 15 (1990) 430–466.
- [18] D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993).
- [19] R. Jonker and A. Volgenant, “A shortest augmenting path algorithm for dense and sparse linear assignment problems,” *Computing* 38 (1987) 325–340.
- [20] D. Knuth, Personal communication (1993).

- [21] H.W. Kuhn, “The Hungarian method for the assignment problem,” *Naval Research Logistics Quarterly* 2 (1955) 83–97.
- [22] E.L. Lawler, *Combinatorial Optimization: Networks and Matroids* (Holt, Rinehart & Winston, New York, 1976).
- [23] Q.C. Nguyen and V. Venkateswaran, “Implementations of Goldberg–Tarjan maximum flow algorithm,” in: D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993) pp. 19–42.
- [24] J.B. Orlin and R.K. Ahuja, “New scaling algorithms for the assignment and minimum cycle mean problems,” Sloan Working Paper 2019-88, Sloan School of Management, Massachusetts Institute of Technology, Cambridge, MA (1988).
- [25] K.G. Ramakrishnan, N.K. Karmarkar and A.P. Kamath, “An approximate dual projective algorithm for solving assignment problems,” in: D.S. Johnson and C.C. McGeoch, eds., *Network Flows and Matching: First DIMACS Implementation Challenge* (American Mathematical Society, Providence, RI, 1993) pp. 431–452.