

in order is done on the supercomputer, it can be handed off to a PC for finishing; at that point in time a second job can be fed to the supercomputer; when the second job is done on the supercomputer, it can proceed to a PC regardless of whether or not the first job is done (since the PCs work in parallel); and so on.

Let's say that a *schedule* is an ordering of the jobs for the supercomputer, and the *completion time* of the schedule is the earliest time at which all jobs will have finished processing on the PCs. This is an important quantity to minimize, since it determines how rapidly El Goog can generate a new index.

Give a polynomial-time algorithm that finds a schedule with as small a completion time as possible.

8. Suppose you are given a connected graph  $G$ , with edge costs that are all distinct. Prove that  $G$  has a unique minimum spanning tree.

9. One of the basic motivations behind the Minimum Spanning Tree Problem is the goal of designing a spanning network for a set of nodes with minimum *total* cost. Here we explore another type of objective: designing a spanning network for which the *most expensive* edge is as cheap as possible.

Specifically, let  $G = (V, E)$  be a connected graph with  $n$  vertices,  $m$  edges, and positive edge costs that you may assume are all distinct. Let  $T = (V, E')$  be a spanning tree of  $G$ ; we define the *bottleneck edge* of  $T$  to be the edge of  $T$  with the greatest cost.

A spanning tree  $T$  of  $G$  is a *minimum-bottleneck spanning tree* if there is no spanning tree  $T'$  of  $G$  with a cheaper bottleneck edge.

- (a) Is every minimum-bottleneck tree of  $G$  a minimum spanning tree of  $G$ ? Prove or give a counterexample.
- (b) Is every minimum spanning tree of  $G$  a minimum-bottleneck tree of  $G$ ? Prove or give a counterexample.
10. Let  $G = (V, E)$  be an (undirected) graph with costs  $c_e \geq 0$  on the edges  $e \in E$ . Assume you are given a minimum-cost spanning tree  $T$  in  $G$ . Now assume that a new edge is added to  $G$ , connecting two nodes  $v, w \in V$  with cost  $c$ .
- (a) Give an efficient algorithm to test if  $T$  remains the minimum-cost spanning tree with the new edge added to  $G$  (but not to the tree  $T$ ). Make your algorithm run in time  $O(|E|)$ . Can you do it in  $O(|V|)$  time? Please note any assumptions you make about what data structure is used to represent the tree  $T$  and the graph  $G$ .

- (b) Suppose  $T$  is no longer the minimum-cost spanning tree. Give a linear-time algorithm (time  $O(|E|)$ ) to update the tree  $T$  to the new minimum-cost spanning tree.

11. Suppose you are given a connected graph  $G = (V, E)$ , with a cost  $c_e$  on each edge  $e$ . In an earlier problem, we saw that when all edge costs are distinct,  $G$  has a unique minimum spanning tree. However,  $G$  may have many minimum spanning trees when the edge costs are not all distinct. Here we formulate the question: Can Kruskal's Algorithm be made to find all the minimum spanning trees of  $G$ ?

Recall that Kruskal's Algorithm sorted the edges in order of increasing cost, then greedily processed edges one by one, adding an edge  $e$  as long as it did not form a cycle. When some edges have the same cost, the phrase "in order of increasing cost" has to be specified a little more carefully: we'll say that an ordering of the edges is *valid* if the corresponding sequence of edge costs is nondecreasing. We'll say that a *valid execution* of Kruskal's Algorithm is one that begins with a valid ordering of the edges of  $G$ .

For any graph  $G$ , and any minimum spanning tree  $T$  of  $G$ , is there a valid execution of Kruskal's Algorithm on  $G$  that produces  $T$  as output? Give a proof or a counterexample.

12. Suppose you have  $n$  video streams that need to be sent, one after another, over a communication link. Stream  $i$  consists of a total of  $b_i$  bits that need to be sent, at a constant rate, over a period of  $t_i$  seconds. You cannot send two streams at the same time, so you need to determine a *schedule* for the streams: an order in which to send them. Whichever order you choose, there cannot be any delays between the end of one stream and the start of the next. Suppose your schedule starts at time 0 (and therefore ends at time  $\sum_{i=1}^n t_i$ , whichever order you choose). We assume that all the values  $b_i$  and  $t_i$  are positive integers.

Now, because you're just one user, the link does not want you taking up too much bandwidth, so it imposes the following constraint, using a fixed parameter  $r$ :

(\*) For each natural number  $t > 0$ , the total number of bits you send over the time interval from 0 to  $t$  cannot exceed  $rt$ .

Note that this constraint is only imposed for time intervals that start at 0, not for time intervals that start at any other value.

We say that a schedule is *valid* if it satisfies the constraint (\*) imposed by the link.

**The Problem.** Given a set of  $n$  streams, each specified by its number of bits  $b_i$  and its time duration  $t_i$ , as well as the link parameter  $r$ , determine whether there exists a valid schedule.

**Example.** Suppose we have  $n = 3$  streams, with

$$(b_1, t_1) = (2000, 1), \quad (b_2, t_2) = (6000, 2), \quad (b_3, t_3) = (2000, 1),$$

and suppose the link's parameter is  $r = 5000$ . Then the schedule that runs the streams in the order 1, 2, 3, is valid, since the constraint (\*) is satisfied:

$t = 1$ : the whole first stream has been sent, and  $2000 < 5000 \cdot 1$

$t = 2$ : half of the second stream has also been sent,  
and  $2000 + 3000 < 5000 \cdot 2$

Similar calculations hold for  $t = 3$  and  $t = 4$ .

- (a) Consider the following claim:

Claim: There exists a valid schedule if and only if each stream  $i$  satisfies  $b_i \leq r t_i$ .

Decide whether you think the claim is true or false, and give a proof of either the claim or its negation.

- (b) Give an algorithm that takes a set of  $n$  streams, each specified by its number of bits  $b_i$  and its time duration  $t_i$ , as well as the link parameter  $r$ , and determines whether there exists a valid schedule. The running time of your algorithm should be polynomial in  $n$ .

13. A small business—say, a photocopying service with a single large machine—faces the following scheduling problem. Each morning they get a set of jobs from customers. They want to do the jobs on their single machine in an order that keeps their customers happiest. Customer  $i$ 's job will take  $t_i$  time to complete. Given a schedule (i.e., an ordering of the jobs), let  $C_i$  denote the finishing time of job  $i$ . For example, if job  $j$  is the first to be done, we would have  $C_j = t_j$ ; and if job  $j$  is done right after job  $i$ , we would have  $C_j = C_i + t_j$ . Each customer  $i$  also has a given weight  $w_i$  that represents his or her importance to the business. The happiness of customer  $i$  is expected to be dependent on the finishing time of  $i$ 's job. So the company decides that they want to order the jobs to minimize the weighted sum of the completion times,  $\sum_{i=1}^n w_i C_i$ .

Design an efficient algorithm to solve this problem. That is, you are given a set of  $n$  jobs with a processing time  $t_i$  and a weight  $w_i$  for each job. You want to order the jobs so as to minimize the weighted sum of the completion times,  $\sum_{i=1}^n w_i C_i$ .

**Example.** Suppose there are two jobs: the first takes time  $t_1 = 1$  and has weight  $w_1 = 10$ , while the second job takes time  $t_2 = 3$  and has weight

$w_2 = 2$ . Then doing job 1 first would yield a weighted completion time of  $10 \cdot 1 + 2 \cdot 4 = 18$ , while doing the second job first would yield the larger weighted completion time of  $10 \cdot 4 + 2 \cdot 3 = 46$ .

14. You're working with a group of security consultants who are helping to monitor a large computer system. There's particular interest in keeping track of processes that are labeled "sensitive." Each such process has a designated start time and finish time, and it runs continuously between these times; the consultants have a list of the planned start and finish times of all sensitive processes that will be run that day.

As a simple first step, they've written a program called `status_check` that, when invoked, runs for a few seconds and records various pieces of logging information about all the sensitive processes running on the system at that moment. (We'll model each invocation of `status_check` as lasting for only this single point in time.) What they'd like to do is to run `status_check` as few times as possible during the day, but enough that for each sensitive process  $P$ , `status_check` is invoked at least once during the execution of process  $P$ .

- (a) Give an efficient algorithm that, given the start and finish times of all the sensitive processes, finds as small a set of times as possible at which to invoke `status_check`, subject to the requirement that `status_check` is invoked at least once during each sensitive process  $P$ .
- (b) While you were designing your algorithm, the security consultants were engaging in a little back-of-the-envelope reasoning. "Suppose we can find a set of  $k$  sensitive processes with the property that no two are ever running at the same time. Then clearly your algorithm will need to invoke `status_check` at least  $k$  times: no one invocation of `status_check` can handle more than one of these processes."

This is true, of course, and after some further discussion, you all begin wondering whether something stronger is true as well, a kind of converse to the above argument. Suppose that  $k^*$  is the largest value of  $k$  such that one can find a set of  $k$  sensitive processes with no two ever running at the same time. Is it the case that there must be a set of  $k^*$  times at which you can run `status_check` so that some invocation occurs during the execution of each sensitive process? (In other words, the kind of argument in the previous paragraph is really the only thing forcing you to need a lot of invocations of `status_check`.) Decide whether you think this claim is true or false, and give a proof or a counterexample.

15. The manager of a large student union on campus comes to you with the following problem. She's in charge of a group of  $n$  students, each of whom is scheduled to work one *shift* during the week. There are different jobs associated with these shifts (tending the main desk, helping with package delivery, rebooting cranky information kiosks, etc.), but we can view each shift as a single contiguous interval of time. There can be multiple shifts going on at once.

She's trying to choose a subset of these  $n$  students to form a *supervising committee* that she can meet with once a week. She considers such a committee to be *complete* if, for every student not on the committee, that student's shift overlaps (at least partially) the shift of some student who is on the committee. In this way, each student's performance can be observed by at least one person who's serving on the committee.

Give an efficient algorithm that takes the schedule of  $n$  shifts and produces a complete supervising committee containing as few students as possible.

**Example.** Suppose  $n = 3$ , and the shifts are

Monday 4 P.M.–Monday 8 P.M.,  
Monday 6 P.M.–Monday 10 P.M.,  
Monday 9 P.M.–Monday 11 P.M.

Then the smallest complete supervising committee would consist of just the second student, since the second shift overlaps both the first and the third.

16. Some security consultants working in the financial domain are currently advising a client who is investigating a potential money-laundering scheme. The investigation thus far has indicated that  $n$  suspicious transactions took place in recent days, each involving money transferred into a single account. Unfortunately, the sketchy nature of the evidence to date means that they don't know the identity of the account, the amounts of the transactions, or the exact times at which the transactions took place. What they do have is an *approximate time-stamp* for each transaction; the evidence indicates that transaction  $i$  took place at time  $t_i \pm e_i$ , for some "margin of error"  $e_i$ . (In other words, it took place sometime between  $t_i - e_i$  and  $t_i + e_i$ .) Note that different transactions may have different margins of error.

In the last day or so, they've come across a bank account that (for other reasons we don't need to go into here) they suspect might be the one involved in the crime. There are  $n$  recent *events* involving the account, which took place at times  $x_1, x_2, \dots, x_n$ . To see whether it's plausible that this really is the account they're looking for, they're wondering

whether it's possible to associate each of the account's  $n$  events with a distinct one of the  $n$  suspicious transactions in such a way that, if the account event at time  $x_i$  is associated with the suspicious transaction that occurred approximately at time  $t_j$ , then  $|t_j - x_i| \leq e_j$ . (In other words, they want to know if the activity on the account lines up with the suspicious transactions to within the margin of error; the tricky part here is that they don't know which account event to associate with which suspicious transaction.)

Give an efficient algorithm that takes the given data and decides whether such an association exists. If possible, you should make the running time be at most  $O(n^2)$ .

17. Consider the following variation on the Interval Scheduling Problem. You have a processor that can operate 24 hours a day, every day. People submit requests to run *daily jobs* on the processor. Each such job comes with a *start time* and an *end time*; if the job is accepted to run on the processor, it must run continuously, every day, for the period between its start and end times. (Note that certain jobs can begin before midnight and end after midnight; this makes for a type of situation different from what we saw in the Interval Scheduling Problem.)

Given a list of  $n$  such jobs, your goal is to accept as many jobs as possible (regardless of their length), subject to the constraint that the processor can run at most one job at any given point in time. Provide an algorithm to do this with a running time that is polynomial in  $n$ . You may assume for simplicity that no two jobs have the same start or end times.

**Example.** Consider the following four jobs, specified by (*start-time, end-time*) pairs.

(6 P.M., 6 A.M.), (9 P.M., 4 A.M.), (3 A.M., 2 P.M.), (1 P.M., 7 P.M.).

The optimal solution would be to pick the two jobs (9 P.M., 4 A.M.) and (1 P.M., 7 P.M.), which can be scheduled without overlapping.

18. Your friends are planning an expedition to a small town deep in the Canadian north next winter break. They've researched all the travel options and have drawn up a directed graph whose nodes represent intermediate destinations and edges represent the roads between them.

In the course of this, they've also learned that extreme weather causes roads in this part of the world to become quite slow in the winter and may cause large travel delays. They've found an excellent travel Web site that can accurately predict how fast they'll be able to travel along the roads; however, the speed of travel depends on the time of year. More precisely, the Web site answers queries of the following form: given an



edge  $e = (v, w)$  connecting two sites  $v$  and  $w$ , and given a proposed starting time  $t$  from location  $v$ , the site will return a value  $f_e(t)$ , the predicted arrival time at  $w$ . The Web site guarantees that  $f_e(t) \geq t$  for all edges  $e$  and all times  $t$  (you can't travel backward in time), and that  $f_e(t)$  is a monotone increasing function of  $t$  (that is, you do not arrive earlier by starting later). Other than that, the functions  $f_e(t)$  may be arbitrary. For example, in areas where the travel time does not vary with the season, we would have  $f_e(t) = t + \ell_e$ , where  $\ell_e$  is the time needed to travel from the beginning to the end of edge  $e$ .

Your friends want to use the Web site to determine the fastest way to travel through the directed graph from their starting point to their intended destination. (You should assume that they start at time 0, and that all predictions made by the Web site are completely correct.) Give a polynomial-time algorithm to do this, where we treat a single query to the Web site (based on a specific edge  $e$  and a time  $t$ ) as taking a single computational step.

19. A group of network designers at the communications company CluNet find themselves facing the following problem. They have a connected graph  $G = (V, E)$ , in which the nodes represent sites that want to communicate. Each edge  $e$  is a communication link, with a given available bandwidth  $b_e$ .

For each pair of nodes  $u, v \in V$ , they want to select a single  $u$ - $v$  path  $P$  on which this pair will communicate. The *bottleneck rate*  $b(P)$  of this path  $P$  is the minimum bandwidth of any edge it contains; that is,  $b(P) = \min_{e \in P} b_e$ . The *best achievable bottleneck rate* for the pair  $u, v$  in  $G$  is simply the maximum, over all  $u$ - $v$  paths  $P$  in  $G$ , of the value  $b(P)$ .

It's getting to be very complicated to keep track of a path for each pair of nodes, and so one of the network designers makes a bold suggestion: Maybe one can find a spanning tree  $T$  of  $G$  so that for every pair of nodes  $u, v$ , the unique  $u$ - $v$  path in the tree actually attains the best achievable bottleneck rate for  $u, v$  in  $G$ . (In other words, even if you could choose any  $u$ - $v$  path in the whole graph, you couldn't do better than the  $u$ - $v$  path in  $T$ .)

This idea is roundly heckled in the offices of CluNet for a few days, and there's a natural reason for the skepticism: each pair of nodes might want a very different-looking path to maximize its bottleneck rate; why should there be a single tree that simultaneously makes everybody happy? But after some failed attempts to rule out the idea, people begin to suspect it could be possible.

Show that such a tree exists, and give an efficient algorithm to find one. That is, give an algorithm constructing a spanning tree  $T$  in which, for each  $u, v \in V$ , the bottleneck rate of the  $u$ - $v$  path in  $T$  is equal to the best achievable bottleneck rate for the pair  $u, v$  in  $G$ .

20. Every September, somewhere in a far-away mountainous part of the world, the county highway crews get together and decide which roads to keep clear through the coming winter. There are  $n$  towns in this county, and the road system can be viewed as a (connected) graph  $G = (V, E)$  on this set of towns, each edge representing a road joining two of them. In the winter, people are high enough up in the mountains that they stop worrying about the *length* of roads and start worrying about their *altitude*—this is really what determines how difficult the trip will be.

So each road—each edge  $e$  in the graph—is annotated with a number  $a_e$  that gives the altitude of the highest point on the road. We'll assume that no two edges have exactly the same altitude value  $a_e$ . The *height* of a path  $P$  in the graph is then the maximum of  $a_e$  over all edges  $e$  on  $P$ . Finally, a path between towns  $i$  and  $j$  is declared to be *winter-optimal* if it achieves the minimum possible height over all paths from  $i$  to  $j$ .

The highway crews are going to select a set  $E' \subseteq E$  of the roads to keep clear through the winter; the rest will be left unmaintained and kept off limits to travelers. They all agree that whichever subset of roads  $E'$  they decide to keep clear, it should have the property that  $(V, E')$  is a connected subgraph; and more strongly, for every pair of towns  $i$  and  $j$ , the height of the winter-optimal path in  $(V, E')$  should be no greater than it is in the full graph  $G = (V, E)$ . We'll say that  $(V, E')$  is a *minimum-altitude connected subgraph* if it has this property.

Given that they're going to maintain this key property, however, they otherwise want to keep as few roads clear as possible. One year, they hit upon the following conjecture:

*The minimum spanning tree of  $G$ , with respect to the edge weights  $a_e$ , is a minimum-altitude connected subgraph.*

(In an earlier problem, we claimed that there is a unique minimum spanning tree when the edge weights are distinct. Thus, thanks to the assumption that all  $a_e$  are distinct, it is okay for us to speak of *the* minimum spanning tree.)

Initially, this conjecture is somewhat counterintuitive, since the minimum spanning tree is trying to minimize the *sum* of the values  $a_e$ , while the goal of minimizing altitude seems to be asking for a fairly different thing. But lacking an argument to the contrary, they begin considering an even bolder second conjecture:

A subgraph  $(V, E')$  is a minimum-altitude connected subgraph if and only if it contains the edges of the minimum spanning tree.

Note that this second conjecture would immediately imply the first one, since a minimum spanning tree contains its own edges.

So here's the question.

- (a) Is the first conjecture true, for all choices of  $G$  and distinct altitudes  $a_e$ ? Give a proof or a counterexample with explanation.
  - (b) Is the second conjecture true, for all choices of  $G$  and distinct altitudes  $a_e$ ? Give a proof or a counterexample with explanation.
21. Let us say that a graph  $G = (V, E)$  is a *near-tree* if it is connected and has at most  $n + 8$  edges, where  $n = |V|$ . Give an algorithm with running time  $O(n)$  that takes a near-tree  $G$  with costs on its edges, and returns a minimum spanning tree of  $G$ . You may assume that all the edge costs are distinct.
  22. Consider the Minimum Spanning Tree Problem on an undirected graph  $G = (V, E)$ , with a cost  $c_e \geq 0$  on each edge, where the costs may not all be different. If the costs are not all distinct, there can in general be many distinct minimum-cost solutions. Suppose we are given a spanning tree  $T \subseteq E$  with the guarantee that for every  $e \in T$ ,  $e$  belongs to *some* minimum-cost spanning tree in  $G$ . Can we conclude that  $T$  itself must be a minimum-cost spanning tree in  $G$ ? Give a proof or a counterexample with explanation.
  23. Recall the problem of computing a minimum-cost arborescence in a directed graph  $G = (V, E)$ , with a cost  $c_e \geq 0$  on each edge. Here we will consider the case in which  $G$  is a directed acyclic graph—that is, it contains no directed cycles.
 

As in general directed graphs, there can be many distinct minimum-cost solutions. Suppose we are given a directed acyclic graph  $G = (V, E)$ , and an arborescence  $A \subseteq E$  with the guarantee that for every  $e \in A$ ,  $e$  belongs to *some* minimum-cost arborescence in  $G$ . Can we conclude that  $A$  itself must be a minimum-cost arborescence in  $G$ ? Give a proof or a counterexample with explanation.
  24. Timing circuits are a crucial component of VLSI chips. Here's a simple model of such a timing circuit. Consider a complete balanced binary tree with  $n$  leaves, where  $n$  is a power of two. Each edge  $e$  of the tree has an associated length  $\ell_e$ , which is a positive number. The *distance* from the root to a given leaf is the sum of the lengths of all the edges on the path from the root to the leaf.

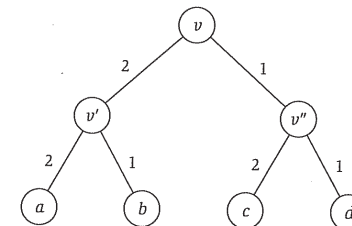


Figure 4.20 An instance of the zero-skew problem, described in Exercise 23.

The root generates a *clock signal* which is propagated along the edges to the leaves. We'll assume that the time it takes for the signal to reach a given leaf is proportional to the distance from the root to the leaf.

Now, if all leaves do not have the same distance from the root, then the signal will not reach the leaves at the same time, and this is a big problem. We want the leaves to be completely synchronized, and all to receive the signal at the same time. To make this happen, we will have to *increase* the lengths of certain edges, so that all root-to-leaf paths have the same length (we're not able to shrink edge lengths). If we achieve this, then the tree (with its new edge lengths) will be said to have *zero skew*. Our goal is to achieve zero skew in a way that keeps the sum of all the edge lengths as small as possible.

Give an algorithm that increases the lengths of certain edges so that the resulting tree has zero skew and the total edge length is as small as possible.

**Example.** Consider the tree in Figure 4.20, in which letters name the nodes and numbers indicate the edge lengths.

The unique optimal solution for this instance would be to take the three length-1 edges and increase each of their lengths to 2. The resulting tree has zero skew, and the total edge length is 12, the smallest possible.

25. Suppose we are given a set of points  $P = \{p_1, p_2, \dots, p_n\}$ , together with a distance function  $d$  on the set  $P$ ;  $d$  is simply a function on pairs of points in  $P$  with the properties that  $d(p_i, p_j) = d(p_j, p_i) > 0$  if  $i \neq j$ , and that  $d(p_i, p_i) = 0$  for each  $i$ .

We define a *hierarchical metric* on  $P$  to be any distance function  $\tau$  that can be constructed as follows. We build a rooted tree  $T$  with  $n$  leaves, and we associate with each node  $v$  of  $T$  (both leaves and internal nodes) a *height*  $h_v$ . These heights must satisfy the properties that  $h(v) = 0$  for each