

The omission of  $\text{INFO}(v)[5]$ , a boolean variable to be maintained at each node  $v$  of the connected interval tree, necessitates the changes below.

page	line	
7	+1	$\text{INFO}(v)[1:4] \rightarrow \text{INFO}(v)[1:5]$
7	+18	Insert: $\text{INFO}(v)[5]$ : a <i>boolean</i> which is <u>true</u> if each leaf interval of the subtree rooted at $v$ is covered by an alive line segment <i>which does not</i> cover all of the $v$ -interval, and is <u>false</u> otherwise.
7	+19	$\text{INFO}(v)[3] = \rightarrow \text{INFO}(v)[3] = \text{INFO}(v)[5] =$
7	+26	$\text{INFO}(\text{root})[1] \neq 0 \rightarrow \text{INFO}(\text{root})[1] \neq 0$ <u>or</u> $\text{INFO}(\text{root})[5]$
8-9		add a 5-th element "f" to each 4-tuple in the figures.
11	+8:+15	change to: Leftcond := $\text{INFO}(v_\ell)[1] \geq 1$ <u>or</u> $\text{INFO}(v_\ell)[5]$ ; rightcond := $\text{INFO}(v_r)[1] \geq 1$ <u>or</u> $\text{INFO}(v_r)[5]$ ; <u>if</u> leftcond <u>and</u> rightcond <u>then</u> $\text{INFO}(v)[4] := 0$ ; $\text{INFO}(v)[5] := \text{true}$ <u>else if</u> leftcond <u>then</u> $\text{INFO}(v)[4] := \text{INFO}(v_r)[4]$ <u>else if</u> rightcond <u>then</u> $\text{INFO}(v)[4] := \text{INFO}(v_\ell)[4]$
	-6	$\text{INFO}(v)[4] \rightarrow \text{INFO}(v)[4:5]$
12	+9:+13	change to: Leftcond := $\text{INFO}(v_\ell)[1] \geq 1$ <u>or</u> $\text{INFO}(v_\ell)[5]$ ; rightcond := $\text{INFO}(v_r)[1] \geq 1$ <u>or</u> $\text{INFO}(v_r)[5]$ ; <u>if</u> leftcond <u>and</u> rightcond <u>then</u> <u>else</u> $\text{INFO}(v)[5] := \text{false}$ ; <u>if</u> leftcond <u>then</u> $\text{INFO}(v)[4] := \text{INFO}(v_r)[4]$ <u>else if</u> rightcond <u>then</u> $\text{INFO}(v)[4] := \text{INFO}(v_\ell)[4]$
	-10	$\text{INFO}(v)[2:4] \rightarrow \text{INFO}(v)[2:5]$
	-5	append: <u>or</u> $\text{INFO}(\text{root})[5]$
15	-3	and $\text{INFO}(v)[3] \rightarrow ,\text{INFO}(v)[3]$ and $\text{INFO}(v)[5]$ .



**stichting  
mathematisch  
centrum**



---

AFDELING INFORMATICA  
(DEPARTMENT OF COMPUTER SCIENCE)

IW 121/79 NOVEMBER

P.M.B. VITÁNYI & D. WOOD

COMPUTING THE PERIMETER OF A SET  
OF RECTANGLES

Preprint

---

**2e boerhaavestraat 49 amsterdam**

*Printed at the Mathematical Centre, 49, 2e Boerhaavestraat, Amsterdam.*

*The Mathematical Centre, founded the 11-th of February 1946, is a non-profit institution aiming at the promotion of pure mathematics and its applications. It is sponsored by the Netherlands Government through the Netherlands Organization for the Advancement of Pure Research (Z.W.O).*

---

1980 Mathematics Subject Classification: 68B15, 68C25, 28.04,  
28A75, 26B15, 51M25

---

ACM-Computer Reviews-Categories: 5.25, 5.39, 4.34, 3.71

Computing the perimeter of a set of rectangles<sup>\*</sup>)

by

Paul M.B. Vitányi<sup>\*\*</sup>) & Derick Wood<sup>\*\*\*</sup>)

ABSTRACT

We describe an algorithm for computing the perimeter (sum of lengths of boundaries) of the area in the plane covered by a given set of  $n$  rectilinearly-oriented rectangles. With some modifications the algorithm also computes the measure (surface) of this area. For the latter task such an algorithm was available before. Our main thrust shall be a comparison of the worst-case performances of the algorithms under various computational assumptions. The results strengthen the conjecture that  $\Theta(n)$  space and  $\Theta(n \log n)$  time simultaneously is unattainable for the perimeter and measure problems when a realistic model of computation is assumed. The algorithms generalize easily to higher dimensions. Without substantially altering the time/storage requirements, the perimeter algorithm can be adapted to output the boundary of a set of intersecting rectangles (or intervals in  $d$ -space), which may be useful in computer graphics.

KEY WORDS & PHRASES: *Geometric complexity, algorithms and data structures, intersecting rectilinearly-oriented rectangles, intersecting intervals in  $d$ -space, perimeter, measure, boundary.*

---

<sup>\*</sup>) The work of the second author supported in part by a Natural Sciences and Engineering Research Council of Canada Grant No. A 7700 and in part by a grant of the Netherlands Organization for the Advancement of Pure Research (Z.W.O.)

<sup>\*\*</sup>) Mathematical Centre, 2e Boerhaavestraat 49,  
Amsterdam, The Netherlands

<sup>\*\*\*</sup>) Unit for Computer Science, McMaster University,  
Hamilton, Ontario L8S 4K1, Canada

This report will be submitted for publication elsewhere.



## 1. INTRODUCTION

Consider the problem of computing the perimeter of a set of  $n$  rectilinearly-oriented rectangles in the plane. That is, we want to determine the total length of the boundaries delineating the total area covered by the rectangles. This is the 2-dimensional particularization of a problem which we shall call the *perimeter problem*. More formally, the perimeter problem is stated as follows. A *closed interval* (generalized rectangular parallelepiped) in  $d$ -space consists of all points  $x = (x_1, x_2, \dots, x_d)$  such that  $l_i \leq x_i \leq u_i$  ( $1 \leq i \leq d$ ) for some fixed numbers  $l_1, l_2, \dots, l_d$  and  $u_1, u_2, \dots, u_d$  with  $l_i < u_i$  for all  $i$ . Given  $n$  closed intervals in  $d$ -space, say  $A_1, A_2, \dots, A_n$ , the perimeter problem asks for an efficient algorithm to compute the  $(d-1)$ -dimensional *measure* of the *boundary* of  $\bigcup_{i=1}^n A_i$ . The boundary of  $\bigcup_{i=1}^n A_i$  is defined as the intersection of  $\bigcup_{i=1}^n A_i$  and the closure of its complement in  $d$ -space. Therefore, the boundary will consist of pieces of  $(d-1)$ -dimensional hyperplanes, and the *perimeter* is the sum of the  $(d-1)$ -dimensional measures of these pieces. For  $d = 1$  the perimeter problem consists in computing the number of *end points* in  $\bigcup_{i=1}^n A_i$ , i.e., twice the number of disjoint closed intervals making up  $\bigcup_{i=1}^n A_i$ .

For  $d = 2$  the perimeter problem consists in finding the sum of the lengths of the line segments which form the boundary of  $\bigcup_{i=1}^n A_i$ , where the  $A$ 's are rectangles in the plane. We shall present an algorithm for solving the 2-dimensional perimeter problem and analyze its complexity. Under certain cost measures this algorithm is optimal. Virtually the same algorithm can be used to solve the case  $d = 1$ , and then yields running times of the same order of magnitude as does the case  $d = 2$ . Presumably, this is optimal for  $d = 1$ . Generalizations of the algorithm to the cases  $d \geq 3$  do not seem to be optimal.

Problems like the perimeter problem belong to the general area of *geometric complexity*. A closely related problem was proposed by KLEE [1977] who called it the *measure problem*. The measure problem asks us to design efficient algorithms to compute the  $d$ -dimensional measure of the union of  $n$  intervals  $A_1, A_2, \dots, A_n$  in  $d$ -space,  $d = 1, 2, \dots$ . It was soon found, that for the 1-dimensional case there was an algorithm with a worst case running time of  $\Theta(n \log n)$ , which was proved to be optimal.

BENTLEY [1977] designed a very efficient algorithm to solve the 2-dimensional case, viz. that of finding the total area in the plane covered by a set of rectangles. For this purpose he introduced a data-structure called the *segment tree*. The algorithm runs in  $\Theta(n \log n)$  time, i.e., the same as the optimal algorithm for the 1-dimensional case, and is therefore optimal too. The generalization of this algorithm to the  $d$ -dimensional case runs in time  $O(n^{d-1} \log n)$ ,  $d \geq 2$ . The algorithm is also described in VAN LEEUWEN and WOOD [1979]. In the latter paper BENTLEY's [1977] result also is improved for dimensions greater or equal to 3, by the exhibition of an algorithm which solves the measure problem in  $d$  dimensions in time  $O(n^{d-1})$ ,  $d \geq 3$ . (This is achieved by using as underlying data structure not the segment tree, but the *quad tree* due to FINKEL and BENTLEY [1974].) Another application of the segment tree has been in algorithms for reporting all pairs of intersecting rectangles from a given set of rectangles in the plane, see e.g. BENTLEY and WOOD [1979].

In the solution to the perimeter problem we give in section 2, the segment tree will once again prove its value as underlying data structure for algorithms solving problems connected with sets of intersecting rectangles. For completeness sake, we show in section 3 how to modify the given algorithm to obtain a variant of BENTLEY's [1977] algorithm for solving the measure problem in 2-space.

Our main thrust shall be to analyze the performances of the algorithms, both for the perimeter and the measure problem, under various models of computation for the representation and manipulation of numbers. It shall be shown, that under reasonable assumptions, such as corresponding to actual computer implementations of the algorithms,  $O(n \log n)$  running time implies  $\Omega(n \log n)$  space in the worst case for the measure problem, and seems impossible for the perimeter problem.

However, it will also appear that if we allow slightly more time, i.e.  $O(n \log^2 n)$ , then both problems can be solved in  $\Theta(n \log \log n)$  space in the worst case. On a RAM with uniform cost criterion both algorithms run in  $\Theta(n \log n)$  time and  $\Theta(n)$  space.

We conjecture that an algorithm to solve either problem in simultaneous  $O(n \log n)$  time and  $O(n)$  space does not exist under "reasonable" models of computation. What is reasonable here shall be argued in sections 4 and 5.



In section 5 we shall note that the generalization of the perimeter algorithm to  $d$  dimensions runs in time  $O(d n^{d-1} \log n)$  and that this may be improved to  $O(d n^{d-1})$  for  $d \geq 3$ . How optimal this is (for  $d \geq 3$ ) we do not know. The presented algorithms for the 2-dimensional case are optimal under the uniform (constant) cost criterion.

Note that from the results one gets the feeling that the perimeter problem is more difficult than the measure problem, but the evidence is not conclusive.

Problems connected with large sets of intersecting rectangles arise in many applications. For instance, the determination of all pairs of intersecting rectangles in the plane is a problem arising in maintaining architectural data bases and forms a crucial step in design rule checking for Very Large Scale Integrated circuitry (VLSI), see BENTLEY and WOOD [1979]. An application in computer graphics might be to determine the boundary of a set of rectangles. In section 5 we indicate how to adapt the perimeter algorithm to do so.

## 2. THE PERIMETER PROBLEM

Suppose we are given  $n$  rectilinearly-oriented rectangles  $A_1, A_2, \dots, A_n$  (in 2-space) represented as

$$A_i = [\ell_{i1} : u_{i1}; \ell_{i2} : u_{i2}], \quad 1 \leq i \leq n,$$

where  $\ell$  = "low" and  $u$  = "high" (or "upper"). For the area or *measure problem* we wish to compute the measure of the total area in the plane covered by the rectangles. For the *perimeter problem* we wish to compute the sum of the lengths of the boundaries delineating this area. An approach to a solution of such problems is usually based on the scanning technique. By scanning the set of rectangles from left-to-right in, say, the  $x$ -direction, and keeping track of appropriate information about the rectangles currently being scanned, the perimeter or measure is accumulated. Such an approach only requires the scanning of the  $2n$  endpoints of the rectangles, since these are the points in the scan at which changes occur. In Figure 1 seven rectangles are displayed together with the corresponding scan lines  $s_1, s_2, \dots, s_{14}$ . To compute the perimeter of the rectangles we first compute the contribution in the  $x$ -direction and secondly in the  $y$ -direction.

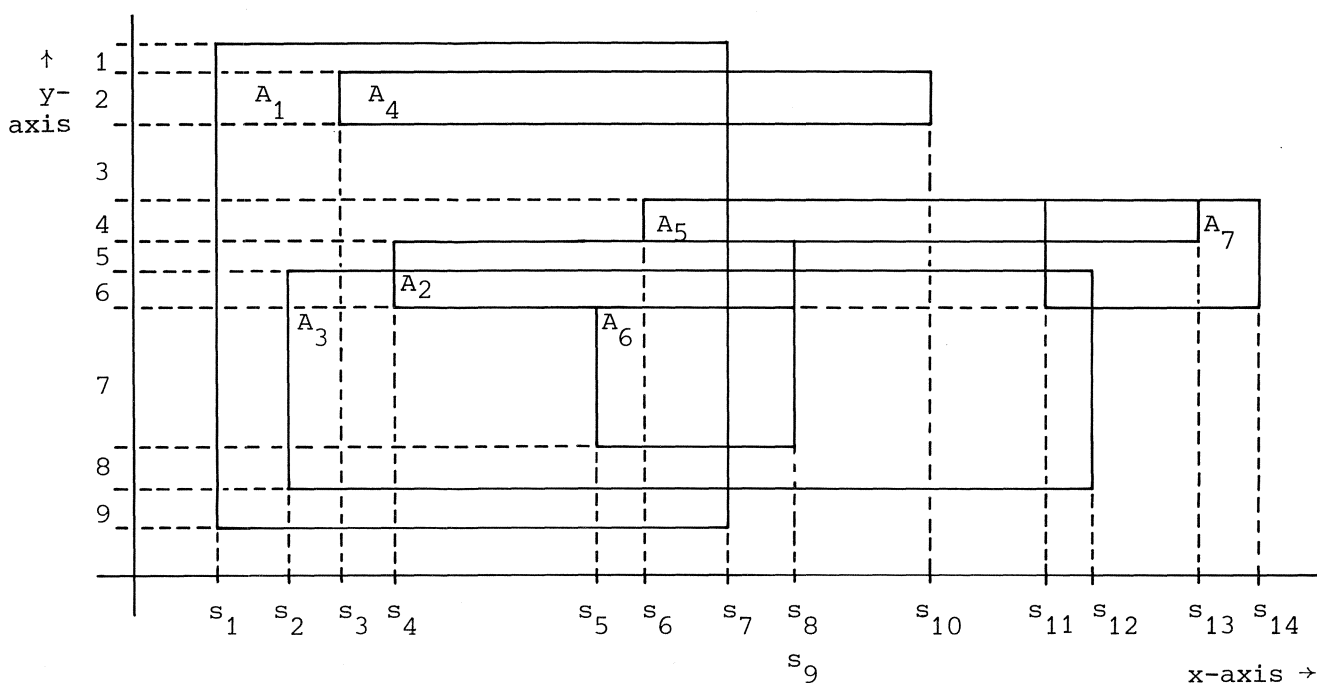


Figure 1.

The contribution  $P_x$  in the x-direction is the sum of the  $2n-1$  partial contributions in the x-direction defined by the  $2n$  scan lines. These, in turn, are simply twice the number of "maximal connected intervals" at the first one of a pair of consecutive scan lines multiplied by the distance in between these scan lines. More precisely:

DEFINITION 1:  $I$  is a *maximal connected interval*, or *m.c.i.*, at scan line  $s$  if:

- (i)  $I$  is a connected interval of  $s$ ;
- (ii)  $I$  is contained by  $\bigcup_{i=1}^n A_i$ ;
- (iii)  $I$  is maximal, i.e.,  $I$  is not contained by any superinterval  $I'$  for which (i) and (ii) hold;
- (iv) There is a number  $\epsilon > 0$  such that (i)-(iii) also hold with  $s$  substituted by  $s_\delta$ , a line parallel to  $s$  but  $\delta$  further away from the origin, for each  $\delta$ ,  $0 \leq \delta \leq \epsilon$ .

Condition (iv) is needed to ensure that upper boundaries of rectangles at  $s$  do not contribute to m.c.i.'s at  $s$ .

Let  $d_1^{(i)}, d_2^{(i)}, \dots, d_{n_i}^{(i)}$  be the m.c.i.'s at scan line  $s_i$  for  $1 \leq i \leq 2n$ . Then the perimeter of the  $n$  rectangles is given by  $P_x + P_y$  where  $P_x$  is the portion of the perimeter parallel to the  $x$ -axis and  $P_y$  is the portion of the perimeter parallel to the  $y$ -axis. We can compute  $P_x$  as

$$(1) \quad P_x = 2 \sum_{i=1}^{2n-1} n_i (s_{i+1} - s_i).$$

$P_y$  is computed similarly. The surface of the area (or the measure) defined by the  $n$  rectangles is given by

$$(2) \quad M = \sum_{i=1}^{2n-1} (s_{i+1} - s_i) \left( \sum_{j=1}^{n_i} \text{length}(d_j^{(i)}) \right),$$

where  $\text{length}(d)$  is the length (or 1-dimensional measure) of  $d$ . Note that when two scan lines  $s_i$  and  $s_{i+1}$  are coincident then  $s_{i+1} - s_i = 0$  and hence no contribution to the perimeter or measure is involved at step  $i$ .

Thus, to compute the perimeter of the set of rectangles we carry out two scans of the figure, while to compute the measure we scan the figure but once. (2) demonstrates that, to compute the two-dimensional measure, we need the one-dimensional measure in the  $y$ -direction at each scan line (perpendicular to the  $x$ -axis.) This fact forms also the basis for the algorithm to compute the measure as given by BENTLEY [1977] (see also VAN LEEUWEN and WOOD [1979]).

With each scan line  $s_i$  we need to associate

- (i) its  $x$ -coordinate (tacitly assumed to be identical with  $s_i$ );
- (ii) whether it corresponds to a left or right (lower or higher) end of a rectangle; and
- (iii) the interval or line segment in the  $y$ -direction the rectangle concerned defines.

To be able to scan the figure in a left-to-right manner we first need to sort the scan lines by their  $x$ -coordinates. To compute  $P_x$  as given by (1) we then move through the scan lines in the sorted order while computing at step  $i$  of this procedure the number of maximal connected intervals at scan line  $s_i$ . This we are able to do by using a suitable modification of the segment tree (BENTLEY [1977], BENTLEY and WOOD [1979]) which we shall call

the *connected interval tree*. The major distinction between the segment tree and the connected interval tree is that identification of the inserted and deleted line segments is not required in the latter case while it is required in the former case.

Let us now describe the connected interval tree, its construction and manipulation. In Figure 1 7 rectangles are displayed which give rise to 9 line fragments in the y-direction, according to the division of the covered part of the y-axis by the y-coordinates of the 7 rectangles. Thus, the projection of each rectangle on the y-axis is made up of a contiguous subset of these fragments. For example, the projection of  $A_3$  on the y-axis is made up of the line fragments numbered 6, 7 and 8. We now construct a minimal height binary tree with as many leaves as there are fragments in the y-direction induced by the rectangles. In our example we construct a minimal binary tree with 9 leaves representing the fragments 1-9 in left-to-right order as displayed in Figure 2: the *skeletal* connected interval tree.

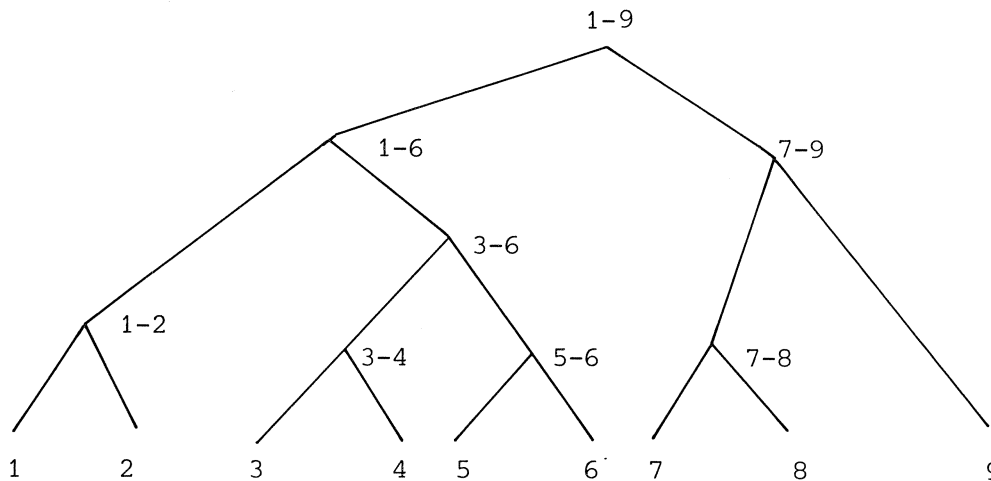


Figure 2. The skeletal connected interval tree

The internal nodes of the tree represent the total interval spanned by their sons; hence for a node  $v$  we write  $v$ -interval to mean the interval represented by  $v$ . A line segment  $I$  covers a  $v$ -interval if the interval represented by  $v$  is contained in  $I$ . At each node  $v$  in the tree we place four

additional pieces of information given by  $\text{INFO}(v)$  [1:4]. The different fields of  $\text{INFO}(v)$  are updated at each scan line  $s_i$  so as to contain the following information at the end of the update.

$\text{INFO}(v)[1]$ : an *integer* which equals the number of times the  $v$ -interval is completely covered by an *alive* line segment, that is, the left (lower) border of a rectangle which has been inserted in the tree at the present or a previous scan line, and not yet been deleted from the tree as the result of meeting a right (upper) border at the present or an earlier scan line. (Cf. also VAN LEEUWEN and WOOD [1979]).

$\text{INFO}(v)[2]$ : a *boolean* which is *true* if the leftmost leaf-interval of the subtree rooted at  $v$  is covered by an alive line segment *which does not* cover all of the  $v$ -interval, and *false* otherwise.

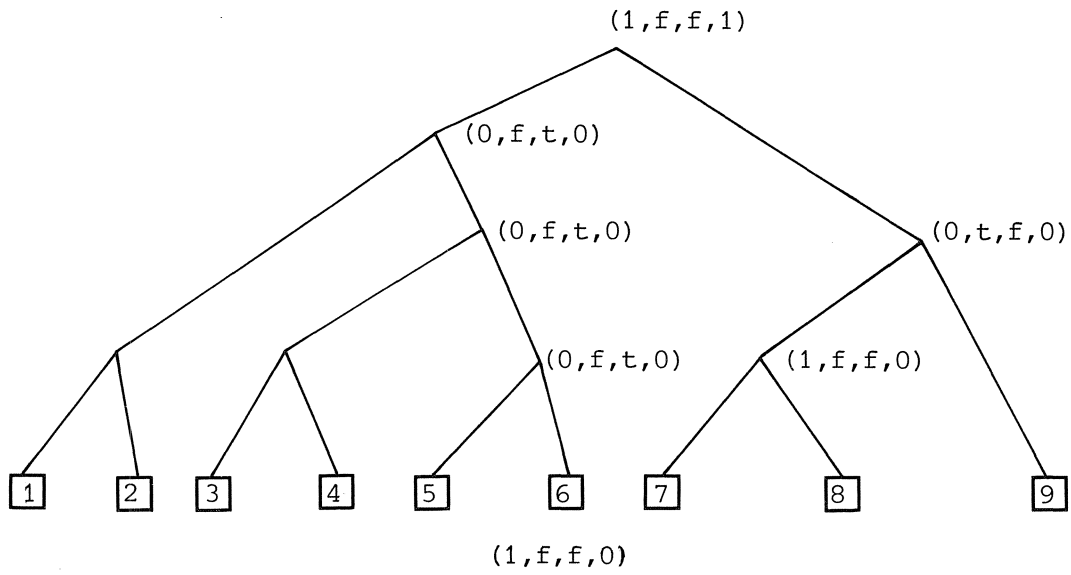
$\text{INFO}(v)[3]$ : a *boolean* as  $\text{INFO}(v)[2]$  but about rightmost leaf-intervals.

$\text{INFO}(v)[4]$ : an *integer* which equals the number of maximal connected intervals, contained in the  $v$ -interval, which do not contain the leftmost leaf-interval nor the rightmost leaf-interval of the subtree rooted at  $v$ .

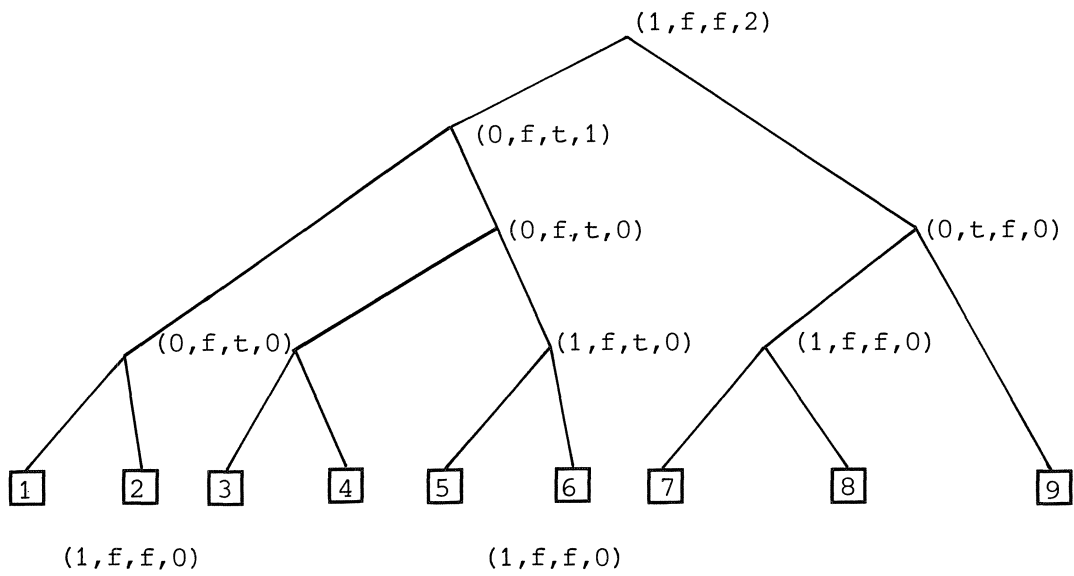
Initially,  $\text{INFO}(v)[1] = \text{INFO}(v)[4] = 0$  and  $\text{INFO}(v)[2] = \text{INFO}(v)[3] = \text{false}$  for all nodes  $v$  of the connected interval tree  $T$ . At each scan line  $s_i$ ,  $1 \leq i \leq 2n$ , we either insert a line segment (viz. the lower border of a rectangle) or delete a line segment (viz. the upper border of a rectangle) from the connected interval tree. Subsequent to the updating which is involved,  $\text{INFO}(\text{root})$  contains the information needed to determine the number of maximal connected intervals at  $s_i$ . Specifically, the number of m.c.i.'s equals 1 if  $\text{INFO}(\text{root})[1] \neq 0$  and equals  $\text{INFO}(\text{root})[2] + \text{INFO}(\text{root})[3] + \text{INFO}(\text{root})[4]$  otherwise (where we assume *true*  $\equiv 1$  and *false*  $\equiv 0$ ). For example, at scan lines  $s_2$ ,  $s_4$  and  $s_6$  of Figure 1 the connected interval tree will be as shown in Figure 3 a, b and c, were only the significant INFO values are displayed. Since  $\text{INFO}(\text{root})[1] \neq 0$  in all three cases there is only 1 m.c.i. as expected. However, if  $A_1$  is deleted from the set of rectangles and from the trees in Figures 1 and 3 respectively, then at  $s_2$  there is one m.c.i. since  $\text{INFO}(\text{root})[2] + \text{INFO}(\text{root})[3] + \text{INFO}(\text{root})[4] = 1$ ; at  $s_4$  there are 2 m.c.i.'s and at  $s_6$  again 2 m.c.i.'s.

8

after  $s_2$ :



after  $s_4$ :



after  $s_6$ :

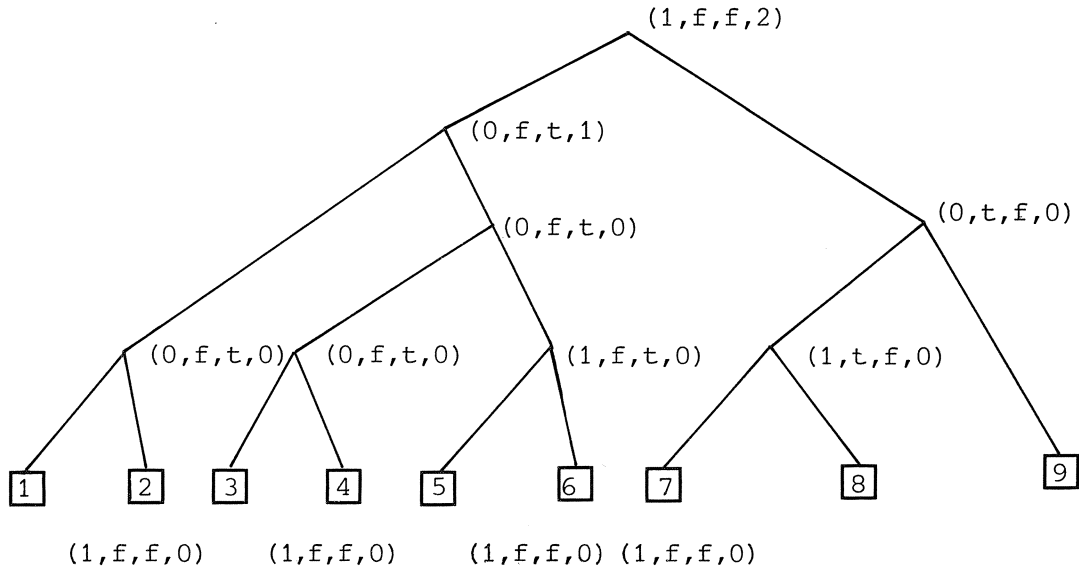


Figure 3.

Having introduced the connected interval tree, we now describe the insertion and deletion procedures for the line segments. Note that the structure of the tree remains unchanged during these operations, only the information contained in the tree is modified.

For convenience sake let  $I$  denote the line segment to be inserted or deleted, and let  $I(v)$  denote the  $v$ -interval for each node  $v$  in the tree  $T$ . Furthermore, let  $v_l$ ,  $v_r$  and  $v_f$  denote the leftson, rightson and father of a node  $v$ ; and let  $\text{LEFT}(I(v))$  and  $\text{RIGHT}(I(v))$  denote the leftmost and rightmost leaf-intervals, respectively, of the subtree rooted at  $v$ .

Basically, the insertion and deletion in the connected interval tree follows the strategy used in the segment tree of BENTLEY [1977] and BENTLEY and WOOD [1979].

On insertion of  $I$  into the connected interval tree  $T$ , we first visit the root. At each node  $v$  visited during insertion one of the following four conditions hold.

- (i)  $I(v) \subseteq I$ .

- (ii)  $I(v) \not\subseteq I$  and  $\text{LEFT}(I(v)) \subseteq I$ .
- (iii)  $I(v) \not\subseteq I$  and  $\text{RIGHT}(I(v)) \subseteq I$ .
- (iv)  $I(v) \not\subseteq I$  and  $\text{LEFT}(I(v)) \not\subseteq I$  and  $\text{RIGHT}(I(v)) \not\subseteq I$ .

Note that conditions (i)-(iv) exclude each other. Condition (i) will increment  $\text{INFO}(v)[1]$  by 1. Condition (ii) causes  $\text{INFO}(v)[2]$  to be set to *true*. Condition (iii) causes  $\text{INFO}(v)[3]$  to be set to *true*. Condition (iv) may cause a change in  $\text{INFO}(v)[4]$  to be discussed below.

On visiting a node  $v$ , the decision to visit any of the sons of  $v$  depends on whether

- (a)  $I(v) \subseteq I$ .
- (b)  $I(v) \not\subseteq I$  and  $I(v_\ell) \cap I \neq \emptyset$ .
- (c)  $I(v) \not\subseteq I$  and  $I(v_r) \cap I \neq \emptyset$ .

In case (a) neither of  $v$ 's sons is visited. In case (b)  $v_\ell$  is visited and in case (c)  $v_r$  is visited. Note that both (b) and (c) may hold, causing a visit to both sons of  $v$ . However, it is not difficult to see (and shown in BENTLEY [1977] and BENTLEY and WOOD [1979]) that at each level of the tree at most 4 nodes may be visited. Since there are at most  $2n-1$  leaves in the tree, and the tree is of minimal height by construction, on each insertion of a line segment into  $T$   $O(\log n)$  nodes are visited. During deletion of  $I$  from  $T$  exactly the same strategy is followed as during insertion. The only differences lie in the way in which the  $\text{INFO}(v)$  associated with the visited nodes  $v$  is updated. Hence also during the deletion of a line segment from  $T$   $O(\log n)$  nodes are visited. Both insertion and deletion in  $T$  have a downward and an upward phase. We now present the insertion and deletion algorithms for connected interval trees.

Insert( $I,v$ );  $I$  is the interval to be inserted,  $v$  is the root node of the connected interval tree;

begin

⋄ Downward phase ⋄

if  $I(v) \subseteq I$  then  $\text{INFO}(v)[1] := \text{INFO}(v)[1] + 1$

else

if  $I(v_\ell) \cap I \neq \emptyset$



```

then if LEFT(I(v))  $\subseteq$  I then INFO(v)[2] := true fi;
      Insert (I, vℓ) fi;
if I(vr)  $\cap$  I  $\neq$   $\emptyset$ 
then if RIGHT(I(v))  $\subseteq$  I then INFO(v)[3] := true fi;
      Insert (I, vr) fi
fi;

```

‡ Upward phase ‡

```

if INFO(vℓ)[1]  $\geq$  1 and INFO(vr)[1]  $\geq$  1
then INFO(v)[4] := 0
else
  if INFO(vℓ)[1]  $\geq$  1
  then INFO(v)[4] := INFO(vr)[4]
  else
    if INFO(vr)[1]  $\geq$  1
    then INFO(v)[4] := INFO(vℓ)[4]
    else
      if INFO(vℓ)[3] or INFO(vr)[2]
      then INFO(v)[4] := INFO(vℓ)[4] + INFO(vr)[4] + 1
      else INFO(v)[4] := INFO(vℓ)[4] + INFO(vr)[4]
      fi
    fi
  fi
fi

```

end of insertion.

Note that during the upward phase of the insertion algorithm  $\text{INFO}(v)[4]$  is re-computed for all nodes of  $T$  on the insertion paths. During the downward phase  $\text{INFO}(v)[1:3]$  is updated for all nodes of  $T$  on the insertion paths.

For deletion we have a similar algorithm, which we now give:

Delete (I, v); I is the interval to be deleted; v is the root of the connected interval tree;

```

begin † Downward phase †
    if  $I(v) \subseteq I$  then  $INFO(v)[1] := INFO(v)[1] - 1$ 
    else if  $I(v_\ell) \cap I \neq \emptyset$  then Delete  $(I, v_\ell)$  fi;
        if  $I(v_r) \cap I \neq \emptyset$  then Delete  $(I, v_r)$  fi
    fi;
    † Upward phase †
     $INFO(v)[2] := INFO(v_\ell)[2]$  or  $(INFO(v_\ell)[1] \geq 1)$ ;
     $INFO(v)[3] := INFO(v_r)[3]$  or  $(INFO(v_r)[1] \geq 1)$ ;
    if  $INFO(v_\ell)[1] \geq 1$  and  $INFO(v_r)[1] \geq 1$ 
    then  $INFO(v)[4] := 0$ 
    else if  $INFO(v_\ell)[1] \geq 1$ 
        then  $INFO(v)[4] := INFO(v_r)[4]$ 
        else if  $INFO(v_r)[1] \geq 1$ 
            then  $INFO(v)[4] := INFO(v_\ell)[4]$ 
            else if  $INFO(v_\ell)[3]$  or  $INFO(v_r)[2]$ 
                then  $INFO(v)[4] := INFO(v_\ell)[4] + INFO(v_r)[4] + 1$ 
                else  $INFO(v)[4] := INFO(v_\ell)[4] + INFO(v_r)[4]$ 
            fi
        fi
    fi
    fi
end delete.

```

Here during the downward phase  $INFO(v)[1]$  is updated and during the upward phase  $INFO(v)[2:4]$ .

For both insertion and deletion the number of nodes visited is  $O(\log n)$ .

The number of maximal connected intervals at a scan line  $s_i$  is contained in  $INFO(\text{root})$  subsequent to the insertion or deletion at this scan line.

To obtain this number we do a query:

```

Query(T) := if  $INFO(\text{root})[1] \neq 0$ 
    then 1
    else  $INFO(\text{root})[4] + INFO(\text{root})[3] +$ 
         $INFO(\text{root})[2]$ 
    fi † we take true  $\equiv 1$  and false  $\equiv 0$  †

```

The algorithm is correct if subsequent to each insertion or deletion the  $\text{INFO}(v)$  for each node  $v$  in  $T$  is correct. This is easily ascertained to be so.

### 3. THE MEASURE PROBLEM

To determine the measure is easier than to determine the perimeter. We use the same connected interval tree skeleton with different additional information at each node  $v$ :  $\text{INFO}(v)[1:3]$ .

$\text{INFO}(v)[1]$ : as before.

$\text{INFO}(v)[2]$ : a *real* equal to the length of the  $v$ -interval.

$\text{INFO}(v)[3]$ : a *real* equal to the sum of the lengths of the maximal connected intervals, covered by inserted and not yet deleted line segments, contained in the  $v$ -interval.

$\text{INFO}(v)[2]$  is determined at the time of setting up the connected interval tree.  $\text{INFO}(v)[1]$  and  $\text{INFO}(v)[3]$  are computed during insertion and deletion, quite similar to  $\text{INFO}(v)[1]$  and  $\text{INFO}(v)[4]$  in the algorithm for the perimeter problem.  $\text{INFO}(v)[3]$  can be determined by the following piece of program (in an upward phase):

```

if  $\text{INFO}(v_\ell)[1] \geq 1$  and  $\text{INFO}(v_r)[1] \geq 1$ 
then  $\text{INFO}(v)[3] := \text{INFO}(v_\ell)[2] + \text{INFO}(v_r)[2]$ 
else if  $\text{INFO}(v_\ell)[1] \geq 1$ 
  then  $\text{INFO}(v)[3] := \text{INFO}(v_\ell)[2] + \text{INFO}(v_r)[3]$ 
  else if  $\text{INFO}(v_r)[1] \geq 1$ 
    then  $\text{INFO}(v)[3] := \text{INFO}(v_\ell)[3] + \text{INFO}(v_r)[2]$ 
    else  $\text{INFO}(v)[3] := \text{INFO}(v_\ell)[3] + \text{INFO}(v_r)[3]$ 
  fi
fi
fi

```

Note that  $\text{INFO}(\text{leaf})[3] = 0$ . Finally, for the measure we have

```

QUERY( $T$ ) := if  $\text{INFO}(\text{root})[1] > 0$  then  $\text{INFO}(\text{root})[2]$ 
  else  $\text{INFO}(\text{root})[3]$ 
fi

```

#### 4. TIME AND SPACE REQUIREMENTS

Using the connected interval tree and its associated updating algorithms we can compute the perimeter of a set of rectangles based on formula (1) and the measure of a set of rectangles based on formula (2). Thus, in both cases the time taken for  $n$  rectangles can be expressed as

$$(3) \quad \text{TIME}(n) = O(n \log n + n + \text{UPDATE}(n) + \text{QUERY}(n)).$$

The endpoints of the rectangles need to be sorted (in both the  $x$ - and the  $y$ -directions for the perimeter problem), taking  $O(n \log n)$  time; the skeletal connected interval tree can be constructed in  $O(n)$  time; and at each scan line there is a deletion or insertion taking a total of  $\text{UPDATE}(n)$  time and a query taking a total of  $\text{QUERY}(n)$  time. Similarly, the space needed for  $n$  rectangles can be expressed as:

$$(4) \quad \text{SPACE}(n) = O(n + \text{INF}(n)),$$

since the  $n$  rectangles and the skeletal connected interval tree, that is without  $\text{INFO}(v)$  at each node  $v$ , require  $O(n)$  space.  $\text{INF}(n)$  denotes the space required to store the  $\text{INFO}(v)$ 's for  $n$  rectangles.

Thus we see that the actual space/time costs depend to a large extent on the costs required to update and store  $\text{INFO}(v)$  or, more precisely, on how much we charge for reals, boolean and integers and the manipulations performed on them as far as relevant in this setting.

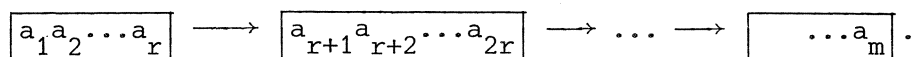
Booleans: we may clearly charge a constant amount in space for storage and in time for manipulation.

Reals : Obviously here everything depends on the precision we require. However, it seems reasonable to charge a constant amount in space for storage and in time for manipulation by keeping reals in  $O(1)$  words.

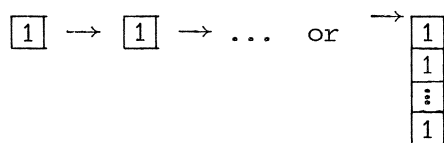
Integers: we will consider three measures:

- (I) Constant cost for storage and time for manipulation.
- (II) Logarithmic cost for storage and time for manipulation;

we can for example keep an  $m$ -bit integer  $a_1 a_2 \dots a_m$  as a linked list of length  $m/r$ :



(III) The cost resulting of storing integers in unary notation, for example as a linked list or counter:



In the following we shall analyze the costs, resulting from the assumed constant costs for booleans and reals together with each of the costs (I), (II) and (III) for integers, in time and space for both the perimeter - and the measure algorithm. Unless stated otherwise, all estimates for running time and storage cost are tacitly assumed to be for the *worst case*.

### I. Constant cost criterion

Each update requires  $\Theta(\log n)$ , viz. the number of visited nodes, and each query requires  $\Theta(1)$ . Hence we immediately have  $\text{UPDATE}(n) = \Theta(n \log n)$ ,  $\text{QUERY}(n) = \Theta(n)$ ; and furthermore  $\text{INF}(n) = \Theta(n)$ ; for both the perimeter and measure algorithms. Hence for both algorithms we have under the constant cost criterion that

$$(5) \quad \text{TIME}_{\text{const.}}(n) = \Theta(n \log n)$$

and

$$(6) \quad \text{SPACE}_{\text{const.}}(n) = \Theta(n).$$

### II. Logarithmic cost criterion

First note that the contribution of  $\text{INFO}(v)[2]$  and  $\text{INFO}(v)[3]$  is  $\Theta(1)$  for each  $v$  in  $T$  for both the perimeter and the measure connected interval trees.

There are at most  $n$  insertions in the connected interval tree. Therefore, for each node  $v$  in  $T$ ,  $\text{INFO}(v)[1]$  needs to count up to  $n$ , which costs  $O(\log n)$  space under the logarithmic cost. Since there are  $\theta(n)$  nodes in the tree, all in all it seems to take  $O(n \log n)$  space to maintain  $\text{INFO}(\cdot)[1]$  for the total tree. However, this estimate is much too crude as we now show.

Since the connected interval tree  $T$  has at most  $2n-1$  leaves, each insertion visits at most  $2 \log 2n$  nodes. (I.e., twice the height of the tree.) At each of these visited nodes  $v$ ,  $\text{INFO}(v)[1]$  might be incremented with 1 during an insertion. Hence the total count of the  $4n$  nodes in the tree, with respect to  $\text{INFO}(\cdot)[1]$  satisfies

$$(7) \quad \sum_{v \in T} \text{INFO}(v)[1] \leq 2n \log 2n$$

since there are  $n$  insertions. If all insertions visit exactly the same nodes in  $T$ , then  $T$  must consist of only the root and the space used by  $\text{INFO}(\text{root})[1]$  is  $O(\log n)$ . We obtain an upper bound on the space used by the  $\text{INFO}(\cdot)[1]$ 's as follows. The amount of space used by the combined  $\text{INFO}(v)[1]$ 's is maximized if they all count to the same maximum, i.e.  $(\log 2n)/2$ . This takes  $O(\log \log n)$  space for each  $\text{INFO}(v)[1]$ . Hence

$$(8) \quad \sum_{v \in T} \text{length}(\text{INFO}(v)[1]) \leq \theta(n \log \log n).$$

Secondly, we need to bound the space used by the  $\text{INFO}(\cdot)[4]$ 's. Clearly this is maximized when the number of connected intervals is maximized. This occurs when the leaf-intervals  $1, 3, 5, \dots, 2n-1$  are covered and  $2, 4, \dots, 2n-2$  are not. Consider a node  $v$  at level  $i$  in the tree  $T$ . Then  $\text{INFO}(v)[4]$  is approximately equal to  $2^{\log n - i}$ . (Hence, e.g.  $\text{INFO}(\text{root})[4] = n$ ,  $\text{INFO}(\text{leaf})[4] = 0$  and  $\text{INFO}(\text{father of leaf})[4] = 1$ ).

Now since there are  $2^i$  nodes  $v$  at level  $i$  in  $T$ , and there are  $\log n$  levels in  $T$ , we obtain:

$$(9) \quad \sum_{v \in T} \text{length}(\text{INFO}(v)[4]) \\ = \sum_{i=0}^{\log n} \sum_{\substack{v \text{ in level} \\ i \text{ of } T}} \text{length}(\text{INFO}(v)[4])$$

$$\begin{aligned}
&= \theta\left(\sum_{i=0}^{\log n} 2^i (\log n - i)\right) \\
&= \theta\left(2^{\log n} \sum_{j=0}^{\log n} j \cdot 2^{-j}\right), \quad \text{by } j = \log n - i, \\
&= \theta\left(n \sum_{j=0}^{\log n} j \cdot 2^{-j}\right).
\end{aligned}$$

Since for  $2^j > j^3$ , e.g. for  $j \geq 10$ , we have that  $j \cdot 2^{-j} < j^{-2}$ , and since  $\sum_{j=1}^{\infty} j^{-2}$  converges, we obtain

$$\begin{aligned}
(10) \quad \sum_{j=0}^{\log n} j \cdot 2^{-j} &< \sum_{j=0}^{\infty} j \cdot 2^{-j} \\
&< \sum_{j=0}^9 j \cdot 2^{-j} + \sum_{j=10}^{\infty} j^{-2} \\
&= c,
\end{aligned}$$

for some constant  $c$ . Hence from (9) and (10) we find that

$$(11) \quad \sum_{v \in T} \text{length}(\text{INFO}(v)[4]) = \theta(n).$$

From (8) and (11) and the previous remarks we find that the worst case space requirements under the logarithmic cost for both the perimeter and measure algorithm are

$$(12) \quad \text{SPACE}_{\log} (n) = \theta(n \log \log n).$$

### *Time requirements*

We now analyze  $\text{TIME}(n)$ , that is,  $\text{UPDATE}(n)$  and  $\text{QUERY}(n)$ . First consider the contribution from  $\text{INFO}(\cdot)[1]$ . We saw above, that the space used is maximized when each node contains a count of  $\theta(\log n)$ . However, time taken is maximized when  $\theta(\log n)$  nodes contain a count of  $\theta(n)$ , that is, when  $\theta(n)$  identical line segments  $I$  are inserted which occasion the visit of  $\theta(\log n)$  nodes. In this case a further insertion or deletion of  $I$  takes  $\theta(\log^2 n)$  time, yielding a total contribution of  $\theta(n \log^2 n)$  for the manipulation time of  $\text{INFO}(\cdot)[1]$ , as opposed to  $\theta(n \log \log n)$  time in the former case. By

merging  $\frac{1}{2}n$  rectangles consuming maximal space and  $\frac{1}{2}n$  rectangles consuming maximal time, as indicated above, we reach simultaneously the worst-case space and time requirements. The time bound given is indeed worst-case since there are at most  $\Theta(n)$  insertions/deletions and each insertion/deletion necessitates the visiting of  $\Theta(\log n)$  nodes, each of which requires at most  $\Theta(\log n)$  time for the updating of  $\text{INFO}(\cdot)[1]$ . Hence, to maintain  $\text{INFO}(\cdot)[1]$  during  $n$  insertions/deletions might take a worst-case time of order

$$(13) \quad \Theta(n \log^2 n).$$

Now consider the contribution of  $\text{INFO}(\cdot)[4]$ . Based on the space analysis of  $\text{INFO}(\cdot)[4]$  we observe that the worst situation which can occur is when we need to add two values of  $\text{INFO}(\cdot)[4]$  one of which is maximal. Since there are at most four nodes visited on each level during an insertion or deletion, and length  $(\text{INFO}(v)[4]) \leq \log n - i$  for a node  $v$  at level  $i$  of the tree, we obtain a contribution of

$$(14) \quad \begin{aligned} & \mathcal{O}\left(\sum_{i=0}^{\log n} (\log n - i)\right) \\ &= \mathcal{O}\left(\sum_{j=0}^{\log n} j\right) = \mathcal{O}(\log^2 n). \end{aligned}$$

Therefore, to maintain  $\text{INFO}(\cdot)[4]$  during  $n$  insertions or deletions might take time of order

$$(15) \quad \Theta(n \log^2 n).$$

Hence, by (13) and (15), we obtain that for both the perimeter and measure algorithm we have a worst-case  $\text{UPDATE}(n) = \Theta(n \log^2 n)$ . Since each query for the perimeter algorithm takes  $\Theta(\log n)$  time ( $\Theta(\log n)$  time in the worst-case) we have  $\text{QUERY}(n) = \Theta(n \log n)$  for the perimeter algorithm. Because of the constant cost for reals we have  $\text{QUERY}(n) = \Theta(n)$  for the measure algorithm. Therefore, under the logarithmic cost criterion, we have for both the perimeter and the measure algorithms a worst-case running time of:

$$(16) \quad \text{TIME}_{\log} (n) = \Theta(n \log^2 n).$$



### III. Unary cost criterion

#### Space requirements

By (7) we need all in all  $\Theta(n \log n)$  space to store all the  $\text{INFO}(\cdot)[1]$ 's. From (9) we compute the space needed for the  $\text{INFO}(\cdot)[4]$ 's, namely:

$$\begin{aligned}
 (17) \quad & \sum_{v \in T} \text{INFO}(v)[4] \\
 &= \Theta\left(\sum_{i=0}^{\log n} 2^i \cdot 2^{\log n - i}\right) \\
 &= \Theta(n \log n).
 \end{aligned}$$

Hence  $\text{INF}(n) = \Theta(n \log n)$  for both perimeter and measure algorithms and therefore, for both,

$$(18) \quad \text{SPACE}_{\text{unary}}(n) = \Theta(n \log n).$$

#### Time requirements

Incrementing and decrementing an integer in unary notation is a constant cost operation, and so is checking for 0. Hence  $\text{INFO}(\cdot)[1]$ 's contribution to  $\text{UPDATE}(n)$  is  $\Theta(n \log n)$  for  $n$  insertion/deletions which visit  $\Theta(\log n)$  nodes apiece. The contribution of the  $\text{INFO}(\cdot)[4]$ 's is simply the cost in time of catenating two integers in unary notation. Clearly this can be carried out in constant time, say in time equal  $c$ . Hence the total of the contribution of the  $\text{INFO}(\cdot)[4]$ 's is given by

$$\Theta\left(n \sum_{i=0}^{\log n} c\right) = \Theta(n \log n)$$

for  $n$  insertions/deletions. Thus, for both the perimeter and the measure algorithm we find under the unary cost criterion:

$$\text{UPDATE}(n) = \Theta(n \log n).$$

A difference between the measure and the perimeter algorithm comes in

with the query. Whereas a query in the measure algorithm extracts a real from  $\text{INFO}(\text{root})$ , at constant cost  $\Theta(1)$ , a query in the perimeter algorithm extracts an integer written in unary from  $\text{INFO}(\text{root})$  at cost  $\Theta(n)$ . Therefore we have for the measure algorithm that

$$\text{QUERY}(n) = \Theta(n)$$

while for the perimeter algorithm

$$\text{QUERY}(n) = \Theta(n^2),$$

since there may be  $\Theta(n)$  connected intervals in the worst case.

Combining the above we obtain:

$$(19) \quad \text{TIME}_{\text{unary}}(n) = \Theta(n \log n) \text{ for the measure problem}$$

and

$$(20) \quad \text{TIME}_{\text{unary}}(n) = \Theta(n^2) \text{ for the perimeter problem.}$$

## 5. CONCLUDING REMARKS

We have presented an algorithm for the perimeter problem and analyzed its time-space requirements (for the worst-case) under three different cost measures for storing and manipulating integers. These are summarized in the following table together with a similar analysis of the related algorithm for the measure problem.

	PERIMETER		MEASURE	
	space	time	space	time
Constant cost	$\Theta(n)$	$\Theta(n \log n)$	$\Theta(n)$	$\Theta(n \log n)$
Logarithmic cost	$\Theta(n \log \log n)$	$\Theta(n \log^2 n)$	$\Theta(n \log \log n)$	$\Theta(n \log^2 n)$
Unary cost	$\Theta(n \log n)$	$\Theta(n^2)$	$\Theta(n \log n)$	$\Theta(n \log n)$

Figure 4. Table

An analysis of the solution for the measure problem as in BENTLEY [1977] or VAN LEEUWEN and WOOD [1979] yields the same results.

The only distinction between the performances of the algorithms for the perimeter and for the measure problem occurs when using unary notation for the integers. This reflects the fact that the one-dimensional measure of the connected intervals is stored as a real, while their number is stored as an integer.

The results strengthen the conjecture that  $\Theta(n)$  space and  $\Theta(n \log n)$  time are unattainable as the worst-case performance of algorithms, for both the perimeter and measure problems, under realistic cost assumptions.

Note that it is realistic to assume a constant cost for storage and manipulation of reals, but that it is unrealistic to assume a constant cost for the storage and manipulation of integers in the discussed algorithms. This is so, because in many applications (like the ones sketched in the introduction) we may deal with a very great number of rectangles in a very limited area. A computation of the perimeter or the measure according to our algorithms *requires* exact bookkeeping of integers: it is easy to construct examples which give completely different answers otherwise. The algorithms in question *approximate* the values of the perimeter and measure insofar as the real values of the rectangle coordinates in the plane are precise.

The *d-dimensional* perimeter problems can be solved by a divide-and-conquer technique in  $O(dn^{d-1} \log n)$  time,  $d \geq 2$ , under the constant cost criterion. This can probably be improved to an  $O(dn^{d-1})$  time algorithm for  $d \geq 3$  by techniques similar to those in VAN LEEUWEN and WOOD [1979].

It is easy to see that the algorithm as presented in section 2 can be changed so as to output the *boundary* of  $\bigcup_{i=1}^n A_i$  too. This algorithm for computing the boundary of a set of intersecting rectangles may be of use to computer graphics and would have the same time/storage requirements as the perimeter algorithm under the unary cost criterion. Hint: maintain the real end points of the m.c.i.'s concerned in INFO(v)[2:4] in the perimeter algorithm, instead of their presence and/or number.

## REFERENCES

- BENTLEY, J.L., [1977], *Algorithms for Klee's rectangle problems*, unpublished notes, Carnegie-Mellon Univ., 1977.
- BENTLEY, J.L., & D. WOOD, [1979], *An optimal worst-case algorithm for reporting intersections of rectangles*, Techn. Rept. 79-CS-13, McMaster University, Unit for Computer Science, 1979.
- FINKEL, R.A. & J.L. BENTLEY, [1974], *Quad trees, A data structure for retrieval on composite keys*, Acta Informatica 4 (1974) 1-9.
- KLEE, V., [1977], *Can the measure of  $\cup_{i=1}^n [a_i, b_i]$  be computed in less than  $O(n \log n)$  steps?*, American Mathematical Monthly 84 (1977) 284-285.
- VAN LEEUWEN, J. & D. WOOD, [1979], *The measure problem for intersecting ranges in d-space*, manuscript in preparation, Univ. of Utrecht, Vakgroep Informatica.