

TRIANGULATING A SIMPLE POLYGON  
IN LINEAR TIME

Bernard Chazelle

CS-TR-264-90

May 1990

---

# Triangulating a Simple Polygon in Linear Time

---

BERNARD CHAZELLE

*Department of Computer Science*

*Princeton University*

*Princeton, NJ 08544*

**Abstract:** We give a deterministic algorithm for triangulating a simple polygon in linear time. The basic strategy is to build a coarse approximation of a triangulation in a bottom-up phase and then use the information computed along the way to refine the triangulation in a top-down phase. The main tools used are (i) the polygon-cutting theorem, which provides us with a balancing scheme, and (ii) the planar separator theorem, whose role is essential in the discovery of new diagonals. Only elementary data structures are required by the algorithm. In particular, no dynamic search trees, finger trees, or point location structures are needed. We mention a few applications of our algorithm.

---

The author wishes to acknowledge the National Science Foundation for supporting this research in part under Grant CCR-8700917.

## 1. Introduction

Triangulating a simple polygon has been one of the most outstanding open problems in two-dimensional computational geometry. It is a basic primitive in computer graphics and, generally, seems the natural preprocessing step for most nontrivial operations on simple polygons [5,13]. Recall that to triangulate a polygon is to partition it into triangles without adding any new vertices. Despite its apparent simplicity, however, the triangulation problem has remained elusive. In 1978, Garey et al. [11] gave an  $O(n \log n)$ -time algorithm for triangulating a simple  $n$ -gon. While it was widely believed that triangulating should be easier than sorting, no proof was to be found until 1988, when Tarjan and Van Wyk [26] discovered an  $O(n \log \log n)$ -time algorithm. Following this breakthrough, Clarkson, Tarjan, and Van Wyk [7] discovered a Las Vegas algorithm with  $O(n \log^* n)$  expected time. Recently, Kirkpatrick, Klawe, and Tarjan [19] gave a new, conceptually simpler  $O(n \log \log n)$ -time algorithm, and they also derived an  $O(n \log^* n)$  bound for the case where vertices have polynomially-bounded integer coordinates [20]. Other results on the triangulation problem include linear or quasi-linear algorithms for restricted classes of polygons [6,10,16,28].

Our main result is a linear-time deterministic algorithm for triangulating a simple polygon. The algorithm is elementary in that it does not require the use of any complicated data structure; in particular, it does not need dynamic search trees, finger trees, or any fancy point location structures.

What makes fast polygon triangulation a difficult problem are the basic inadequacies of either a pure top-down or a pure bottom-up approach. To proceed top-down is to look at the whole polygon and infer global information right away. One can rely on the polygon-cutting theorem [4] which says that the polygon can be cut along a diagonal into two roughly equal-size pieces. The immediate dilemma is that to find such a diagonal appears just as difficult as triangulating the whole polygon to begin with. Moreover, one would actually need a sublinear method for finding such a diagonal to keep our hopes for an optimal triangulation algorithm alive. A bottom-up approach, on the other hand, involves computing, say, triangulations of subpieces of the polygon's boundary. This suffers from the obvious flaw that too much information gets to be computed. Indeed, diagonals for small pieces of the boundary are not guaranteed to be diagonals of the whole polygon and might therefore be wasted. Our solution is to mix bottom-up and top-down approaches together. The basic strategy is to build a coarse approximation of a triangulation in a bottom-up phase and then use the information computed along the way to refine the triangulation in a top-down phase. The main tools used are (i) the polygon-cutting theorem, which provides us with a balancing scheme, and (ii) the planar separator theorem [22], whose role is essential in the discovery of new diagonals.

Here is a more detailed overview of the algorithm. As was observed in [6,10] a triangulation of a polygon can be obtained directly from its horizontal visibility map: this is the partition of the polygon obtained by drawing horizontal chords from the vertices. We can extend this notion easily and speak, more generally, of the visibility map of any simple polygonal curve (Figure 2.1). Chazelle and Incerpi [6] showed how to build the visibility map of an  $n$ -vertex curve in  $O(n \log n)$  time, using divide-and-conquer. Their algorithm mimics mergesort: Assuming that  $n$  is a power of 2, at the  $k$ th stage ( $k = 1, 2, \dots, \log n$ ), the boundary of the polygon is decomposed into chains of size  $2^k$ , whose visibility maps are computed by piecing together the maps obtained at the previous stage. Each

stage can be accomplished in a linear number of operations, so computing the visibility map of the polygon takes  $O(n \log n)$  time.

The new algorithm follows exactly the same pattern: It goes through an *up-phase* of  $\log n$  stages, each involving the merging of maps obtained at the previous stage. The novelty we bring into this process is to use only coarse samples of the visibility maps during the merges. In this way, we can carry out an entire stage in sublinear time and beat  $n \log n$ . These samples are uniform submaps of the visibility maps; uniform in the sense that they approximate the visibility maps anywhere equally well. Of course, in the end, we also need an efficient way to *refine* the submap for the whole polygon into its full-fledged visibility map. After this is done, it takes only linear time to derive a triangulation. To refine a submap we go down through stages in reverse (a *down-phase*): each transition refines the submap incrementally, until we get back to the first stage, at which point the full visibility map emerges at last. Figure 1.1 illustrates the meaning of the up- and down-phases.

Perhaps it is the right time to wonder whether our approach is not inherently flawed from the start. How sound is it to mimic mergesort when our goal is to beat  $n \log n$ ? Any attempt to speed up mergesort by using “coarse samples” of the lists to be merged is trivially doomed. So, what’s so different about polygons? The difference is rooted in a notion which we call *conformality*. This is perhaps the single most important concept in our algorithm, for it is precisely where mergesort and triangulation part ways. Recall that the polygon-cutting theorem is a geometric analog of the centroid theorem for free trees: a visibility map has a tree structure and so can be written as a collection of “blobs” of roughly equal size, themselves interconnected in a tree pattern. These blobs are the constituents of a submap. Merging two submaps can thus be equated with “merging” two trees together. The mergesort equivalent of a submap would be a sublist (of one of the lists to be merged) obtained by picking keys at regular intervals. Notice that merging two such sublists might in the worst case produce a new sublist whose corresponding intervals are up to twice the size of the original intervals. This coarsening effect prevents us from speeding up mergesort, because repairing the damage might involve computing medians or things of that nature for which no shortcuts can be found. To be sure, as we shall see, equally bad things can happen with submaps; repairing the damage, however, can be done by simply adding new chords to submaps, which can be made to take only sublinear time. To make this possible we must keep the coarseness of submaps under control by requiring that the tree structure of a submap be of bounded degree: in our terminology “conformal” actually means degree at most 4. Enforcing conformality is the linchpin of the algorithm and, as one should expect, its most delicate and subtle part as well.

Although our algorithm is an outgrowth of the mergesort-like method of Chazelle and Incerpi, it goes far beyond it. For example, the algorithm *never* actually merges visibility maps but only submaps (which is done quite differently). We must also borrow ideas from a number of other sources: As we mentioned before, one of them is the polygon-cutting theorem, and more specifically, the hierarchical polygon decomposition of Chazelle and Guibas [5]. There exist optimal algorithms for computing such decompositions [3,13], but as it turns out, standard sub-optimal methods work just fine for our purposes. Merging submaps requires a primitive to detect new chords. In the Chazelle-Incerpi method, the detection is limited to constant-size domains, so it can be done naively. But here, the domains can be arbitrarily large so we need a sublinear ray-shooting method. This

cannot be done by using fast planar point location, like in the algorithm of Kirkpatrick, Klawe, and Tarjan [19]. The reason is that we need to support merging of two point location structures, and the known methods, even the dynamic ones, are inadequate in that regard. We turn the problem around by using a weak form of divide-and-conquer based on Lipton and Tarjan's planar separator theorem [22].

## 2. Merging Visibility Maps

Following the tradition of visibility algorithms, we begin by restricting their applicability to polygons where no two distinct vertices have the same  $y$ -coordinates. Of course, the standard excuse still works: We can easily get around this assumption by rotating the reference system appropriately or applying the perturbation techniques of [9,29]. Turning now to visibility maps, we observe that in the Chazelle-Incerpi method, chords are extended only toward the interior of the polygon, with a special rule to determine how far they should extend. Kirkpatrick et al. [19] use the more consistent (and simpler) scheme of extending chords on both sides. This amounts to thinking of a polygonal curve as a very thin polygon embedded in a cylindrical plane that lets chords wrap around infinity. For the reasons we give below, we will find it more convenient to embed our objects in a topological manifold, called the *spherical plane*, which is equivalent to a 2-sphere. Proofs of correctness in the polygonal merging business tend to be tricky and ridden with painful case-analyses. One reason for this, we found, is that these proofs often attempt to establish topological facts by geometric means, thus adding unnecessary complication. By sticking to topological considerations as much as possible, proofs become much simpler, provided, of course, that the ambient space has the "right" topology.

What is the right topology in this context? One problem with the cylindrical plane is that although it has a Jordan curve theorem, the two regions created by removing a simple closed curve may now have one of three types: (1) an open disk, (2) a cylinder  $S^1 \times (0, +\infty)$ , or (3) a perforated cylinder. We simplify all that by defining our ambient space to be the *spherical plane*: This is the product space  $[-\infty, \infty]^2$  with the following identification rules (i)  $(-\infty, y) = (+\infty, y)$ , (ii)  $(x, -\infty) = (-x, -\infty)$ , (iii)  $(x, +\infty) = (-x, +\infty)$ , for all  $x, y \in [-\infty, +\infty]$ . The spherical plane is homeomorphic to a 2-sphere, so we now have the nicest Jordan curve theorem of all: the removal of any simple closed curve creates two open disks.

The only remaining difficulty is the orientability of Jordan curves. If we fix an orientation of the spherical plane, the boundary of any region homeomorphic to a disk, being of codimension 1, will have an induced orientation, called *clockwise*; the reverse orientation is called *counterclockwise*. Conversely, given a Jordan curve, a *clockwise* tour of the curve will refer to the orientation induced by one of the two homeomorphic disks that it bounds. Unless one of these disks is already understood, we choose one unambiguously by fixing a reference point somewhere in the plane and looking at the unique disk that does not contain it. Of course, this means restricting ourselves to curves that avoid that point, but this is only a minor inconvenience. From now on, we will assume that the reference point is at  $(0, +\infty)$ . In this way, to speak of the clockwise traversal of a Jordan curve, without reference to an enclosed region, makes full sense.

**A. The visibility map.** Given a simple (nonclosed) polygonal curve  $C$  with vertices  $v_1, \dots, v_n$ , we define the *visibility map* of  $C$ , denoted  $V(C)$ , as the planar subdivision formed by extending two horizontal segments from each vertex  $v_i$ , one in each direction. Each segment, which we call a *chord*, is extended until it meets another point of  $C$ . If it goes to infinity in the Cartesian plane, then it will actually wrap around in the spherical plane until it hits  $C$  again. Adding chords to  $C$  subdivides the spherical plane into *regions*: triangles or trapezoids with two horizontal segments and two nonhorizontal segments; all of them possibly divided up into several collinear edges (Figure 2.1).

In order to distinguish between the two sides of an edge, we give each edge of  $C$  an infinitesimal width so as to make the curve  $C$  into a very thin simple polygon. (This is just a conceptual device, so in particular, no actual perturbation of the polygon needs to be performed on the computer.) The boundary of that polygon is called the *double boundary* of  $C$  and is denoted  $\partial C$ . By abuse of terminology we refer to the two *sides* of the double boundary (which is meaningful geometrically but not topologically). Each vertex of  $C$  that is not a local extremum in the  $y$ -direction is associated with two *companion* vertices of  $\partial C$ , one on each side of the curve (Figure 2.2.1). In this way, each vertex is incident upon exactly one chord (possibly the same chord for the two companions). But what about local extrema? For each such vertex of  $C$ , if it is not one of the two endpoints, we create a total of four vertices in  $\partial C$  (Figure 2.2.2): two companion pairs of duplicate vertices; one pair on each side of  $\partial C$ . The duplicate vertices in a given pair are next to each other along  $\partial C$ . By convention, we say that one of these pairs, the one on the “inside” of the turn, gives rise to a chord of null length. Finally, as shown in Figure 2.2.3, for each endpoint of  $C$  we create two duplicate vertices; these can also be called companions since they are next to each other as well as on both sides of  $\partial C$ . Figure 2.3 illustrates these definitions. Note that for simplicity we have not numbered vertices connected to null-length chords. We have now ensured that each vertex of  $\partial C$  is incident upon exactly one chord of the visibility map. Therefore, any vertex—and actually also any point—of  $\partial C$  has a unique horizontal “chord direction” (left or right) assigned to it. Roughly speaking, this direction points to the left of an observer walking clockwise around  $\partial C$ .

If a chord of  $V(C)$  has  $p$  and  $q$  as endpoints, we say that  $p$  and  $q$  *see* each other with respect to  $C$ , or simply, see each other, when  $C$  is understood. We also say that  $p$  and  $q$  are mutually *visible*. More generally, if  $p$  and  $q$  are two points (not necessarily vertices) of  $\partial C$  with the same  $y$ -coordinates, we say that  $p$  and  $q$  see each other with respect to  $C$  if one of the two (relatively open) segments joining  $p$  and  $q$  lies completely outside of  $C$  (regarded as a thin simple polygon). By extension the segment  $pq$  is also called a chord. For example, in Figure 2.3, points labeled 6 and 14 as well as points 3 and 4 and points 1 and 1' do *not* see each other. We close our discussion of visibility with a simple but useful fact.

**Lemma 2.1.** *If we remove a pair of mutually visible points from the double boundary of a simple polygonal curve, then no chord can connect the two resulting pieces.*

*Proof:* Let  $C_1$  and  $C_2$  be the two pieces of the double boundary resulting from the removal of a pair of visible points. Together with  $C_1$  and  $C_2$ , the chord  $c$  connecting the two points subdivides the spherical plane into three polygonal regions (i.e., regions bounded by simple polygonal curves), one

of which is the infinitesimally thin polygon  $C$  itself. Any chord connecting  $C_1$  and  $C_2$  lies outside the “polygon”  $C$ , so it must cross  $c$ . But this is impossible because chords are horizontal. ■

The circular sequence of chord endpoints in  $V(C)$  encountered during a clockwise traversal of  $\partial C$  is called the *canonical vertex enumeration* of  $V(C)$ : note that it contains other points besides the vertices of  $\partial C$ . Figure 2.3 provides the sequence  $1, 1', \dots, 17', 18$ , with the primes indicating duplicated vertices. We have omitted to number the endpoints of null-length chords. Speaking of which, note that null-length chords create empty regions. The traversal of the double boundary leads to a canonical enumeration of the regions (with repetitions). In the case of Figure 2.3, we have the list (including only the nonempty regions for simplicity)

$$(I, II, III, IV, V, VI, VII, VIII, VII, IX, VII, VI, V, IV, III, X, XI, X, III, II, I, XII, I, XIII).$$

It is easy to see that the dual graph of the subdivision is a tree, naturally called the *visibility tree* of  $C$ . This graph is defined by associating a distinct node with each region (empty or nonempty) of the visibility map and connecting any pair of nodes whose corresponding regions share a common chord. Note that this also includes null-length chords. Figure 2.1 shows the tree without the nodes associated with empty regions. Since any two consecutive regions in the canonical region enumeration have a common chord, the visibility tree is indeed connected. Why can't it have cycles? it suffices to show that removing any chord  $ab$  would disconnect the graph.

First, assume that  $a$  and  $b$  are not duplicates of each other (although they might be companions). Then, from Lemma 2.1, removing  $a$  and  $b$  splits the double boundary into two pieces which are closed under visibility. It follows that the boundary of any region contains segments from exactly one of the two pieces, and therefore can be classified by the piece to which these segments belong. Consequently,  $ab$  is the only chord separating a region associated with one piece from a region associated with the other piece. This concludes the first case. Assume now that  $a$  and  $b$  are duplicates of each other. If the chord  $ab$  has zero length then removing it isolates an empty region from the rest, and our claim holds. So, suppose that  $ab$  has nonzero length. Then  $a$  is either the highest or lowest point of  $\partial C$ , therefore removing  $ab$  would disconnect the upper or lower halfplane from the other regions. This completes the proof that the visibility tree is indeed a tree. (There is also a simple topological proof, which is given in Lemma 2.2 below.)

**B. Visibility submaps.** A useful operation is to remove chords from  $V(C)$ . Before we go any further, however, we wish to accommodate for the possibility that  $V(C)$  has been augmented with some additional chords (something that will often happen later). Obviously, these new chords cannot connect vertices of  $\partial C$  (since all have been used up) but rather arbitrary points on the curve. Although this clearly affects  $V(C)$  as well as the visibility tree, everything we said previously regarding canonical enumerations is trivially extended to this new situation. So, from now on, let us treat  $V(C)$  either in its original state or in some augmented form. We will specify which applies whenever the distinction needs to be made.

The operation we want to discuss now involves removing a given chord from  $V(C)$ . Since  $V(C)$  may have additional chords, we need to be careful about the meaning of a removal. A chord has two endpoints; none, one, or two of which are vertices of  $\partial C$ . So, the removal of a chord entails removing not only the chord itself but also those endpoints that are *not* vertices of  $\partial C$ , and glueing back  $\partial C$  at those points. This cleanup operation is meant to prevent the presence of vertices stranded in the middle of an edge of  $\partial C$ : in other words, any vertex that is not a vertex of  $\partial C$  *must* be the endpoint of a chord.

Removing one or several chords produces a polygonal subdivision of the spherical plane, called a *submap* of  $V(C)$  (Figure 2.4). The boundary of a nonempty region of the submap is an oriented circular sequence of horizontal segments, called *exit chords*, alternating with pieces of  $\partial C$ , referred to as the region's *arcs*. It is important to keep in mind that arcs are pieces of  $\partial C$  and *not* of  $C$ . For example, let us assume that all null-length chords have been removed in the submap of Figure 2.4. Then the boundary of the region labeled II consists of a two-edge arc (beginning at the big dot), followed by an exit chord, a one-edge arc, an exit chord, a three-edge arc (not a four-edge arc!), an exit chord, a one-edge arc, and one final exit chord. Note that some arcs may be of zero length, as is the case in the region labeled V. As illustrated in Figure 2.4, a clockwise traversal of  $\partial C$  induces a counterclockwise traversal of each region of the submap. More formally, we have the following fact.

**Lemma 2.2.** *Let  $A_1, \dots, A_k$  be the counterclockwise enumeration of the (oriented) arcs of a nonempty submap region (as induced from the region's orientation). Then each  $A_i$  is oriented clockwise with respect to  $\partial C$ . Moreover, the sequence  $A_1, \dots, A_k$  also occurs clockwise around  $\partial C$ .*

*Proof:* The curve  $\partial C$  is homeomorphic to a circle embedded in the spherical plane. Adding a chord is topologically equivalent to connecting two points on the circle by a simple curve lying on the same side of the circle. The requirement that all these curves should be mutually disjoint induces a parenthesis system which immediately reveals the tree structure of the dual graph. This is similar to the parenthesis systems in Jordan sorting [17]. As an example, Figure 2.5 depicts the topological equivalent of the submap of Figure 2.4. From this perspective, the lemma is now obvious. ■

As was the case with  $V(C)$ , a clockwise traversal of  $\partial C$  induces canonical vertex/region enumerations of the submap. Figure 2.4 gives the region enumeration: I,II,III,II,IV,II,V,II,I. (Recall that that particular submap is supposed to have had all its null-length chords removed and therefore has no empty regions.) Bold dots mark the points during the clockwise traversal of the double boundary where the canonically enumerated regions are first discovered. An important requirement is that a vertex enumeration of a submap should list only the endpoints of actual exit chords and thus might skip over many vertices of  $\partial C$ . In this way, canonical enumerations of any type will take time proportional to the number of regions and not to the number of vertices (which might be much higher). We define the *weight* of a region as 0 if the region is empty, or else as the maximum number of non-null length edges in any of its arcs. For example, regions I and II in Figure 2.4 have weights 4 and 3, respectively.



Although weights count only edges of nonzero length, chords of zero length (if any) are taken into account inasmuch as they separate arcs. In other words, an arc never contains any chord, whether that chord be of nonzero length or not. Also, note the important role played by the double boundary in the definition of a region's weight. Indeed, a region may have a very small weight because its arcs are all small; but this does not prevent any one of these arcs to have a huge number of vertices *on the other side* of the double boundary. Of course, because vertices of  $\partial C$  have companions those extra vertices have to be endpoints of exit chords.

Combinatorially, a region corresponds to a subtree of the visibility tree of  $C$ . The dual graph of a submap is obtained by contracting the edges of the visibility tree that correspond to the removed chords (Figure 2.6); Being derived from the visibility tree by graph-minor operations, the dual graph of a submap is itself a tree (as was clear in the proof of Lemma 2.2), which we simply call the "tree of the submap". Note that, conversely, contracting any edge of the visibility tree amounts to removing the corresponding chord from the visibility map. The *weight* of a node naturally refers to the weight of its corresponding region.

Since no two distinct vertices of  $C$  have the same  $y$ -coordinates, the degree of any node in a visibility tree cannot exceed 4. We should not expect this to be always true of submap trees, however, so we distinguish the trees of *conformal submaps* as those with node-degree at most 4. By analogy with the polygon-cutting theorem we can decompose a conformal submap in a hierarchical manner. The idea is to pick the centroid of the submap tree and observe that there exists at least one incident edge whose removal leaves two subtrees of size at most  $3/4$  the original size. Associating the removed edge with the root of a binary tree and recursing in this fashion with respect to the root's two children provides a *tree decomposition* of the submap. The tree has depth logarithmic in the number of regions (which is the number of chords plus one). The internal nodes (resp. leaves) are in bijection with the exit chords (resp. regions) of the submap. By using straightforward tree labeling techniques we can find the centroid node, and hence the first edge to be removed, in linear time. Proceeding recursively gives us a simple  $O(m \log m)$  time algorithm for computing the tree decomposition of a submap of  $m$  regions. We will use this result below because it is simple and practical. Optimal methods exist but they are all fairly complicated [3,13].

One final piece of classification concerns the respective sizes of a submap's regions. We say that a submap is  $\gamma$ -granular if (i) every node of its tree has weight at most  $\gamma$  and (ii) contracting any edge incident to at least one node of degree less than 3 produces a new node whose weight exceeds  $\gamma$ . Note that this weight might be less than the added weight of the two nodes of the contracted edge. This is because one or both endpoints of the removed chord might not be vertices of  $\partial C$  and might thus disappear. If only condition (i) holds, then the submap is  $\gamma$ -semigranular. Adding condition (ii) makes the semigranularity *maximal* in some sense. Finally, by default, if (i) holds but the submap has no exit chord, it is still said to be  $\gamma$ -granular.

**Lemma 2.3.** *If  $C$  is a polygonal curve with  $n$  vertices, any  $\gamma$ -granular conformal submap of the (possibly augmented) visibility map of  $C$  has  $O(n/\gamma + 1)$  regions and each region is bounded by  $O(\gamma)$  edges.*

*Proof:* We can assume that the tree of the submap has at least one edge, otherwise the lemma is trivial. Among the edges of that tree, let  $E$  be the set of those incident to at least one node of degree less than three. It is trivial to show by induction on the size of the tree that  $E$  accounts for at least a fixed fraction of all the edges. Now, contracting any edge in  $E$  produces a merged region of weight greater than  $\gamma$ , meaning that it has an arc with more than  $\gamma$  edges of nonzero length. Since a vertex of  $C$  can give rise to at most four vertices of  $\partial C$ , and removed chords do not leave extra vertices behind except those of  $\partial C$ , such an arc must involve at least on the order of  $\gamma$  distinct vertices of  $C$ . If contracting any edge of  $E$  were always to produce a *distinct* merged region then it would follow from the pigeon-hole principle that  $E$ , and hence the whole tree, has  $O(n/\gamma + 1)$  nodes. Unfortunately, two edges of  $E$  might produce overlapping merged regions. From the conformality of the submap, however, we know that a given vertex of  $C$  can still be used at most a constant number of times in this counting argument, therefore  $E$  has indeed  $O(n/\gamma + 1)$  nodes and the first part of the lemma is established. The second part derives from the conformality of the submap, which ensures that there is a bounded number of arcs per region and hence that the size of any region, i.e., its number of bounding edges, is at most proportional to its weight. ■

We close this discussion of visibility submaps by establishing a useful result on the topology of regions and chords. Let  $D$  be a closed disk on the spherical plane, let  $\bar{D}$  be its boundary, and let  $ab$  denote a diametrical chord of  $D$ . Pick two distinct points  $c, d$  on  $\bar{D}$  such that  $a, c, b$  occur in clockwise order (with respect to  $D$ ), and let  $A$  be a simple curve lying inside  $D$  and running from  $c$  to  $d$ . Consider the circular arc that runs clockwise from  $d$  to  $c$  and let  $B_i$  ( $i = 1, 2$ ) be the closures of its intersections with the two circular arcs of  $\bar{D} \setminus \{a, b\}$  (Figure 2.7). Note that  $B_i$  consists of 0, 1, or 2 circular arcs. We say that a subset  $\alpha$  of  $A$  is *shielded from*  $B_j$  if either  $B_j$  is empty or else no point of  $\alpha$  can be connected to any point of  $B_j$  by a curve (understood here as a closed set) that lies entirely inside  $D$  and does not intersect either  $ab$  or  $A$ . In Figure 2.7.1, for example, the piece of  $A$  running from  $c$  to  $a'$  is shielded from  $B_2$  since none of its points can be connected to  $B_2$  without crossing  $ab$ . Similarly, the piece from  $a'$  to  $b'$  is shielded from  $B_1$  since a connection to it would have to cross  $ab$  or  $A$ .

**Lemma 2.4.** *If  $a \in B_1 \cup B_2$  (resp.  $b \in B_1 \cup B_2$ ), let  $a'$  (resp.  $b'$ ) be the first (resp. last) point of  $ab \cap A$  encountered when traversing the diametrical chord  $ab$  from  $a$  to  $b$ . The points  $a'$  and  $b'$  (which might not exist) subdivide  $A$  into a total of one, two, or three connected curves, each of which is shielded from some  $B_j$  (not necessarily the same  $j$  for all curves). Furthermore, an appropriate  $B_j$  can be identified simply on the basis of  $a, b, c, d, a', b'$ .*

*Proof:* We can assume that both  $B_1$  and  $B_2$  are nonempty and that  $A$  intersects  $ab$  (else the lemma is trivially correct). By attaching  $B_1 \cup B_2$  to  $A$ , we obtain a simple closed curve within  $D$ , which is, therefore, the boundary of a region  $R$  of  $D$  homeomorphic to a disk. If  $a$  (resp.  $b$ ) belongs to  $B_1 \cup B_2$  then the segment  $aa'$  (resp.  $bb'$ ) lies within  $R$  and thus, acting like a chord, subdivides  $R$  into two regions. Since  $aa'$  and  $bb'$  cannot intersect, together they subdivide  $R$  into two or three disk-like regions. The boundary of each such region intersects the boundary of  $D$  in a single connected arc and therefore avoids one  $B_i$  (outside of  $a$  and  $b$ ) entirely. Figure 2.7 illustrates the two possible cases; note that the third case, where the counterclockwise traversal from  $c$  to  $d$  avoids both  $a$  and

$b$ , was eliminated earlier, since it corresponds to a situation where one of the  $B_i$ 's is empty. None of the curves obtained by removing  $a'$  and  $b'$  (if they exist) from  $A$  can belong to more than one of the subdividing regions, so each of them is shielded from some  $B_i$ . Which one can be determined by simple examination of the relative order of the points  $a, b, c, d, a', b'$  around the boundary of  $R$ . ■

**C. Representation issues.** How do we represent visibility maps and submaps as data structures? We first describe our mode of representation, then we point out some of its idiosyncracies and explain why they are needed. Let  $P$  be the input polygonal curve (the one whose visibility map is sought) and let  $C$  be the piece of its boundary whose visibility map (or submap) we wish to represent. We shall assume that  $P$  is nonclosed; this is not restrictive since a little hole can always be punctured if it is closed to begin with. We assume that the edges of  $P$  are stored in a table (the *input table*) in the order in which they occur along the boundary of  $P$ . (A doubly-linked list would also do.) Note that the notion of double boundary need not be encoded explicitly, i.e., no edges are duplicated in the table. The input table is read-only: it is never to be modified or even copied. A visibility submap of  $V(C)$  is represented by its own data structure: arcs are encoded by pointing directly into the input table. More precisely, each arc is represented by a separate *arc-structure*. Null-length arcs can be represented explicitly so let us assume that the arc has nonzero length. Let  $e_1, \dots, e_t$  be the edges of an arc in clockwise order along the double boundary, where  $e_1$  and  $e_t$  are the edges adjacent to the two chords connected by the arc. If  $t = 1$  then the arc-structure consists of a single pointer into the input table to the edge  $e$  of  $P$  that contains  $e_1$ . Since  $e_1$  is an edge of the double boundary, we also need to indicate by a flag which side of  $e$  is to be understood. We don't need to record the endpoints of the arc because chords take care of that. If  $t > 1$ , we store the same information as above but now with respect to both  $e_1$  and  $e_t$  in that order. We say that a map or submap is given in *normal form* if the following information is provided.

- (i) The tree of the submap (or map) is represented in standard edge/node adjacency fashion.
- (ii) Each edge of the tree contains a record describing the corresponding chord as well as pointers to the arc-structures of the two, three, or four arcs adjacent to it. Conversely, each arc-structure has a pointer to the node of the tree whose corresponding region is incident upon the arc in question.
- (iii) The arc-structures are stored in a table (the *arc-sequence table*) in the order corresponding to a canonical traversal of the double boundary  $\partial C$ .
- (iv) If the submap is conformal then its tree decomposition should be available.

We choose what may seem to be a contrived representation of a submap in order to use storage proportional not to the number of edges in the submap but rather to the number of regions (which is roughly the same as the number of chords and arcs). It is essential to avoid duplication of information because we will need to encode a collection of submaps whose number of *distinct* features is  $O(n)$  but whose combined size is  $\Theta(n \log n)$ . Note that our representation is powerful enough to let us perform canonical vertex/region enumerations in optimal time. If we wish to, we can also enumerate all the vertices of  $\partial C$  in clockwise order directly from a canonical vertex enumeration of the submap, since any arc can be reconstructed explicitly from the succinct information given by the arc-structure: it

suffices to explore the input table between the locations indicated by the two pointers of the arc-structure. Note that caution must be used since the arc might wrap around both sides of  $\partial C$ . This can be detected when we traverse the arc as soon as we reach an edge of  $P$  incident to an endpoint of  $C$ .

Perhaps a less obvious task is to retrieve the arc structure corresponding to an arc, given one of its edges. More specifically, suppose that we are given an edge  $e$  of  $C$  and a point  $q$  on it. The question is to find the arc-structures of all those arcs in the submap that pass through the point  $q$ . By passing through, we don't care whether the arc is on the right or wrong side of the double boundary, so for example, if  $q$  is not an endpoint of any chord in the submap then there are at most two distinct arcs to be found; otherwise there are at most four of them. Once we have located the two endpoints of  $C$  in the arc-sequence table (i.e., which arcs pass through them) we can conceptually break up the circular arc sequence into two linear sequences and perform in each of them a binary search, using the name of the containing edge  $e$  as a query. Either search might take us to a unique arc-structure, in which case we are done, or to a contiguous interval of arc-structures: this might happen if  $e$  contributes several arcs. We can disambiguate by pursuing the binary search, now using, say, the  $y$ -coordinate of  $q$  as a query. The total running time is logarithmic in the number of arcs, or more conservatively,  $O(\log |C|)$ . This operation will be very useful later: we call it the *double identification* of a point of  $C$ .

We have said repeatedly that a submap has a tree structure. But let us change our perspective a little and look at a submap as a standard planar subdivision, without distinguishing between chords and arc edges. There are many standard representations of planar graphs [2,14,23] which allow us to navigate through such a subdivision in constant time per step taken. Normal-form representations are not quite that powerful. The only problem is if we attempt to cross from one side of an arc to the other along, say, a straight line. In order to find which region we are about to enter we must perform a double identification. The problem here is that unlike with standard graph representations we do not keep adjacency information between regions and edges (except for chords). Also, we do not provide explicit correspondence between the features on the two sides of an edge of  $C$ . For reasons which will become clear later, it would be a very bad idea to try to do so.

**D. Merging two submaps.** We consider the problem of merging two conformal submaps. Let  $C_1$  and  $C_2$  be two polygonal curves of  $n_1$  and  $n_2$  vertices respectively, whose union  $C$  forms a connected vertex-to-vertex piece of the input (simple and nonclosed) polygonal curve  $P$ . Let  $S_i$  ( $i = 1, 2$ ) be a  $\gamma_i$ -granular conformal submap of  $V(C_i)$ , with  $\gamma_1 \leq \gamma_2$ . From now on, we shall use the notation  $\bar{\alpha}$  to refer to the portion of  $C$  to which an arc  $\alpha$  of  $\partial C$  corresponds. We assume that each  $S_i$  is given in normal form and that the following set of primitives is available. For each region arc  $\alpha$  of  $S_i$  ( $i = 1, 2$ ) specified by a pointer to its arc-structure:

- (i) There exists a ray-shooting oracle which, given any point  $p$  with a horizontal direction (left or right) and any subarc  $\alpha'$  of  $\alpha$  specified by its two endpoints (along with two pointers to the input table to indicate the names of the edges of  $P$  that contain these two endpoints as well as two flags indicating which side of  $\partial P$  is to be understood), reports the single point of  $\alpha'$  (if any) that a ray of light shot from  $p$  in the given direction would hit in the absence of any obstacle

except  $\alpha'$ . In addition to the point hit, the report should also include the name of the edge of  $P$  that contains it. The report should take  $O(f(\gamma_i))$  time, where  $f$  is a nondecreasing function.

- (ii) There exists an arc-cutting oracle which, in  $O(g(\gamma_i))$  time, subdivides the subarc  $\alpha'$  into at most  $g(\gamma_i)$  subarcs  $\alpha_1, \alpha_2, \dots$ , such that (1) each  $\alpha_j$  is specified by its two endpoints and a pair of pointers into the input table to indicate which edges of  $P$  contain these endpoints; the pair should be ordered to reflect a clockwise traversal along  $\partial P$  and two flags should be included to indicate on which sides of  $\partial P$  these endpoints fall; (2) the relative interior of no  $\alpha_j$  contains a point of  $\partial C_i$  that corresponds to an endpoint of  $C_i$ , that is, each subarc must lie entirely on one side of  $\partial C$ ; and (3) a normal-form  $h(\gamma_i)$ -granular conformal submap of  $V(\bar{\alpha}_j)$  has already been computed and is available in normal form, along with its tree decomposition. Again,  $g$  and  $h$  are assumed to be nondecreasing functions.

Let  $\gamma$  be any integer at least as large as  $\gamma_2 \geq \gamma_1$ . Our goal in this section is to compute a normal-form  $\gamma$ -granular conformal submap of  $V(C)$  in  $O((n_1/\gamma_1 + n_2/\gamma_2 + 1)f(\gamma_2)g(\gamma_2)(h(\gamma_2) + \log(n_1 + n_2)))$  time. We proceed in three stages. First, we find which points of  $\partial C$  can be seen by the endpoints of the exit chords of  $S_i$  ( $i = 1, 2$ ) and by the companion vertices resulting from the duplication of  $C_1 \cap C_2$  and we use this information to create a submap  $S$  of  $V(C)$ , called the *fusion* of  $S_1$  and  $S_2$ . In a second stage we ensure that the submap is conformal. We bring it to the desired granularity in the final stage.

*Fusion:*

By symmetry, we may limit our discussion to the problem of determining the points of  $\partial C$  that are seen by the endpoints of the exit chords of  $S_1$  and by the companion vertices resulting from the duplication of  $C_1 \cap C_2$ . The idea is to repeat the work described below with respect to  $S_2$ . As it turns out, it is possible to unify the two passes into one, which might be the preferred solution in practice, but in order to make the proof of correctness simpler we choose to keep the two passes separate. Let  $a_0$  and  $a_{m+1}$  be the companion vertices, as they appear next to each other clockwise around  $\partial C_1$ , resulting from the duplication of the vertex  $C_1 \cap C_2$  in  $\partial C_1$ . Let  $a_1, a_2, \dots, a_m$  be the canonical vertex enumeration of  $S_1$ . Recall that this enumerates the exit chord endpoints in  $S_1$  as we encounter them going clockwise around  $\partial C_1$ . Since the sequence is circular we can assume that  $a_0$  precedes  $a_1$  and  $a_{m+1}$  follows  $a_m$ . Note that it could happen that  $a_0$  and  $a_{m+1}$  are already part of the sequence, but for the sake of generality, we assume that they are not and therefore add them on. A clockwise tour around  $\partial C_1$  that begins at  $a_0$  thus ends at  $a_{m+1}$ . We shall compute the points of  $\partial C$  seen by  $a_0, a_1, \dots, a_{m+1}$  in that order.

We begin with a simple observation. Given any point  $p$  of  $\partial C_i$  and the arc to which it belongs (specified by its arc-structure), we can determine which point of  $\partial C_i$  it sees (with respect to  $C_i$ ) in  $O(f(\gamma_i))$  time. This operation is called *local shooting*. If  $p$  is an endpoint of an exit chord we can easily do that (even in constant time). If not, then  $p$  belongs to a unique region of  $S_i$ , which we can determine in constant time (i.e., via its node in the submap tree), and the point of  $\partial C_i$  that it sees lies on one of the region's arcs. Using the appropriate ray-shooting oracles, we can find that point by checking each arc in turn. The claim on the time follows from the conformality of  $S_i$  which ensures that at most four arcs need to be checked. Note that local shooting is still possible even if

$p$  does not lie on  $\partial C_i$ : it can lie anywhere in the spherical plane as long as a horizontal direction (left or right) has been specified and we know in which region of  $S_i$  the point lies. We still call this operation local shooting.

We need a start-up phase to launch our fusion. Using local shooting, we find the point of  $\partial C_1$  that  $a_0$  sees with respect to  $C_1$ . Although  $a_0$  does not lie on  $\partial C_2$  we obviously know in which region of  $S_2$  it lies therefore, as we just pointed out, we can also use local shooting to find the point of  $\partial C_2$  that  $a_0$  sees with respect to  $C_2$ . These two pieces of information combine to give us the unique point  $c_0$  of  $\partial C$  that  $a_0$  sees with respect to  $C$ . We distinguish between two cases. If  $c_0$  belongs to  $\partial C_2$  then we set  $p = a_0$  and we call the region of  $S_2$  crossed by  $a_0 c_0$  *current*: the start-up phase is over (Figure 2.8.1). Otherwise, from Lemma 2.1, the chord  $a_0 c_0$  splits  $\partial C$  into two curves, each closed under visibility. One of these curves, the one running from  $a_0$  to  $c_0$  clockwise, is a piece of  $\partial C_1$  (Figure 2.8.2), so the points of  $\partial C$  that its exit chord endpoints see all belong to  $\partial C_1$ , and thus are available directly from  $S_1$ . We can therefore skip all the way to  $c_0$ . Now, however,  $c_0$  sees a point of  $\partial C_2$ , namely  $a_0$ , so we set  $p = c_0$  and call the region of  $S_2$  containing  $a_0$  *current*.

Technically, it is not quite true that  $a_0$  is always a point of  $\partial C_2$ . It coincides with one most often, but when it sits at a local extremum (in the  $y$ -direction) it is not one. What is true, however, is that when  $c_0$  cannot see a point of  $\partial C_2$ , an infinitesimal deformation of  $\partial C_2$  locally around  $a_0$  can make  $c_0$  see one. This is a minor technicality which will not affect the remainder of the fusion algorithm, so for simplicity we will still go on saying that  $c_0$  sees a point of  $\partial C_2$  with respect to  $C$ . Another minor problem is that  $a_0 c_0$  might lie on an exit chord of  $S_2$  and thus there might be more than one candidate for the status of current region. We break ties by electing the region that we locally enter as we leave  $p$  in a clockwise traversal of  $\partial C_1$ . This concludes the start-up phase. At this point, we have the following situation (all visibility being understood with respect to  $C$ ):

- A. The points of  $\partial C$  that are seen by the exit chord endpoints of  $S_1$  on the portion of  $\partial C_1$  running clockwise from  $a_0$  to the *current* point  $p$  have all been determined already; we should mention here that later on  $p$  might not always be one of the  $a_i$ 's.
- B. The point  $q$  of  $\partial C$  that is seen by  $p$  belongs to  $\partial C_2$  and the chord  $pq$  lies in the region of  $S_2$  called *current*. If  $p$  lies on the boundary of the current region then leaving  $p$  clockwise around  $\partial C_1$  locally enters the current region.

These two conditions will form our loop invariant, that is, they will hold prior to every iteration of the process which we now describe. Let  $A_i$  denote the region arc of  $S_1$  running from  $a_{i-1}$  to  $a_i$  (in clockwise order around  $\partial C_1$ ); by extension  $A_1$  (resp.  $A_{m+1}$ ) stands for the subarc extending from  $a_0$  to  $a_1$  (resp.  $a_m$  to  $a_{m+1}$ ). Let  $A_k$  be the arc containing  $p$ . In the likely event that  $p$  is the endpoint of a chord of  $S_1$  and thus belongs to two arcs, we must choose the one starting (and not ending) at  $p$ , i.e., we set the condition  $p \neq a_k$ . If  $p = a_{m+1}$  however, the algorithm simply terminates and no  $A_k$  need be defined. Let  $R$  denote the *current* region prior to entering the following iteration: Loop through  $j = k, k + 1, \dots$ , until

- (i)  $a_j$  lies in  $R$  and the point of  $\partial C$  that  $a_j$  sees belongs to  $\partial C_2$  (Figure 2.9.1), or
- (ii) the previous condition (i) does not hold, but  $R$  has an exit chord such that the point of  $\partial C$  seen by one of its endpoints belongs to  $A_j$  but is distinct from  $p$  (Figure 2.9.2), or

(iii)  $j = m + 2$ ; in that case, we terminate.

If case (i) occurs, find which point of  $\partial C$  is seen by  $a_j$ , declare that all  $a_i$ 's ( $k \leq i < j$ ) see points of  $\partial C_1$  (with respect to  $C$ ), set  $p = a_j$ , let the current region still be  $R$ , and iterate. If case (ii) occurs, then of all the candidate endpoints, i.e., those endpoints satisfying (ii), determine the one that sees the point  $p'$  that is the first encountered as we traverse  $\partial C_1$  clockwise starting from  $p$ . In Figure 2.9.2, for example,  $p'$  is the point labeled  $p_1$  and the chosen endpoint is labeled  $q_1$ . Next, declare that all  $a_i$ 's ( $k \leq i < j$ ) see points of  $\partial C_1$ , set  $p = p'$ , make current the region of  $S_2$  which we enter as we cross the exit chord at  $p'$  along  $\partial C_1$ , and iterate. In case (iii) we stop and declare that all  $a_i$ 's ( $k \leq i \leq m + 1$ ) see points of  $\partial C_1$ . We have made several claims and skipped over important implementation issues in order to get the main idea of the algorithm across. Next, we fill in all the missing parts and substantiate our claims.

**Lemma 2.5.** *It is possible to compute the fusion  $S$  of  $S_1$  and  $S_2$  in  $O((n_1/\gamma_1 + n_2/\gamma_2 + 1)(f(\gamma_2) + \log(n_1 + n_2)))$  time.*

*Proof:* Since  $\gamma_1 \leq \gamma_2$  and  $f$  is nondecreasing, it suffices to consider the case  $i = 1$ . We have already shown that the start-up phase leads to a situation which satisfies the loop invariant, so it suffices to establish the correctness of the inner loop past  $a_0$ . In case (i) we know that  $a_j$  lies in  $R$  (actually in its interior) and sees a point of  $\partial C_2$ , so Invariant (B) is satisfied. How about (A)? We made the claim that  $a_k, \dots, a_{j-1}$  all see points of  $\partial C_1$ . But actually the negation of (i) for  $a_k, \dots, a_{j-1}$  is not strong enough to reach the necessary conclusion about what  $a_k, \dots, a_{j-1}$  must see. Any of these points (if they exist) either indeed sees  $\partial C_1$  or perhaps lies outside of  $R$ . Why should lying outside  $R$  imply seeing  $\partial C_1$ ? Suppose that, for some  $\ell$  ( $k \leq \ell < j$ ),  $a_\ell$  lies in region  $R'$  distinct from  $R$  (like  $a_{k+1}$  in Figure 2.9.1) but also sees  $\partial C_2$ . We will derive a contradiction. Let  $\mathcal{A}$  denote the directed portion of  $\partial C_1$  as we go from  $p$  to  $a_\ell$  clockwise, and let  $q$  (resp.  $q'$ ) be the point of  $\partial C$  seen by  $p$  (resp.  $a_\ell$ ). The union of  $\mathcal{A}$ , the chords  $pq$  and  $a_\ell q'$ , and the portion  $\mathcal{B}$  of  $\partial C_2$  running clockwise (with respect to  $C_2$ ) from  $q'$  to  $q$  forms the boundary of a simple polygon (Figure 2.10). Since the dual graph of a submap is a tree, there is a unique exit chord  $ab$  of  $R$  that leads to  $R'$  (note that  $ab$  need not be an exit chord of  $R'$ , since there might be one or even several regions separating  $R$  from  $R'$ ). The chord  $ab$  is cut by  $\mathcal{A}$ , therefore one of its endpoints, say,  $a$ , must lie in  $\mathcal{B}$ . Let  $a'$  be the point of  $ab \cap \mathcal{A}$  first encountered as we go from  $a$  to  $b$  along the chord. The points  $a$  and  $a'$  see each other with respect to  $\partial C$ , and  $a'$  lies in  $A_h$ , for some  $h$  between  $k$  and  $\ell$ . Because, in clockwise order around  $\partial C_1$ , the point  $a'$  is leaving  $R$  locally, it cannot be equal to  $p$ . Therefore, the inequality  $\ell < j$  implies that case (ii) must have already occurred when  $j$  was equal to some integer between  $k$  and  $\ell$ , which is impossible.

Having shown that the loop invariant remains satisfied through case (i), we must do the same with case (ii). Let  $\mathcal{A}$  now denote the directed portion of  $\partial C_1$  as we go from  $p$  to  $a_j$  clockwise. The new assignment of  $p$  is the first point of  $\mathcal{A}$  that sees an endpoint of an exit chord of  $R$ . Certainly, the new assignment of the current region satisfies Invariant (B). (If one is bothered by the fact that  $p$  is now on the boundary between two regions, one might conceptually want to move  $p$  an infinitesimal amount along  $\mathcal{A}$  so as to be squarely in the new region.) Figure 2.9.2 shows three candidate endpoints, with  $q_1$  winning the contest. Turning now to Invariant (A), we must prove our

claim that the points of  $\partial C$  seen by  $a_k, \dots, a_{j-1}$  all belong to  $\partial C_1$ . We omit the proof, since it is identical to the previous one, with the role of  $a_j$  now played by  $p'$ . To summarize, the proof given in the previous paragraph shows that if  $a_\ell$  sees  $\partial C_2$  then either  $a_\ell$  lies in  $R$  (case (i)) or it does not, but then, we must fall in case (ii) upon or prior to getting to  $a_\ell$ . We will use this result again below.

What about termination? Obviously, the three cases exhaust all possible situations. Every time we fall in either of the two cases (i,ii) we determine more visibility information, so that all visibility relations are known from  $a_0$  all the way to the current position of  $p$ . How about the last iteration, the one leading to case (iii). We claimed that all  $a_i$ 's ( $k \leq i \leq m+1$ ) see points of  $\partial C_1$ , which follows directly from our last observation in the previous paragraph. The proof of correctness is now complete.

Let us now analyze the complexity of the algorithm. To test whether  $a_j$  lies in  $R$  can be done in  $O(f(\gamma_2))$  time by using the ray-shooting oracles on each arc: first, we find which point of an arc is hit by a ray of light shot from  $a_j$  in its assigned chord direction. If there is no hit,  $a_j$  is not in  $R$ . Else, let  $s$  be the first point hit by the ray over all the arcs of  $R$ . Whether  $p$  lies in  $R$  or not can be directly inferred from the local orientation of the hit at  $s$  and which side of the double boundary is hit. If  $a_j$  lies in  $R$ , then  $s$  is the point of  $\partial C_2$  seen by  $p$  with respect to  $C_2$ . Next we use local shooting within  $S_1$  to determine the point  $t$  of  $\partial C_1$  hit by a ray of light shot from  $a_j$  in its assigned direction. Note that most often (i.e., when  $0 < j \leq m$ )  $a_j$  is the endpoint of a chord of  $S_1$  so  $t$  is just the other endpoint of the chord. We are in case (i) if and only if  $sa_j$  lies in  $ta_j$ . The test takes  $O(f(\gamma_1) + f(\gamma_2)) = O(f(\gamma_2))$  time.

Regarding (ii), we must assess how fast we can find the point of  $\partial C$  that is seen by an endpoint  $a$  of a given exit chord  $ab$  of  $R$ . Actually, we must find that point only if it belongs to  $A_j$ . So, we can shoot a ray of light from  $a$  toward  $A_j$  in the appropriate direction and see what happens, which takes  $O(f(\gamma_1))$  time. If we don't get a hit or if the hit does not lie on  $ab$ , the endpoint can be disqualified. Otherwise, we must find whether the point  $s$  of  $A_j$  hit by the ray of light can see  $a$  with respect to  $C$ : the problem here is that other arcs  $A_i$  ( $i \neq j$ ) might get in the way. Using local shooting in  $S_1$ , however, we can shoot a ray of light from  $s$  toward  $a$ . If this is not the natural shooting direction from  $s$  (which we can decide in constant time by checking the local orientation of the edge containing  $s$ ) then we know that the "companion" point of  $s$  prevents it from seeing  $a$ . Otherwise, the point  $t$  hit by the ray is found in  $O(f(\gamma_1))$  time. If  $st$  contains  $sa$  then  $s$  and  $a$  see each other with respect to  $C$  and we fall in case (ii), else we know that case (ii) cannot occur with respect to the endpoint  $a$  of the chord  $ab$  (although it might occur with respect to other endpoints of exit chords in  $R$ ). This shows that testing case (ii) takes  $O(f(\gamma_1))$  time.

We thus have shown that every elementary test (i,ii) can be performed in time  $O(f(\gamma_1) + f(\gamma_2)) = O(f(\gamma_2))$ . At each such test we advance through the list of  $A_i$ 's or we report a pair of visible points in  $\partial C$ , one of which is an endpoint of an exit chord of  $S_2$ . These reports are never duplicated because we always move forward among the  $A_i$ 's. Therefore, to discover all the chords of the fusion  $S$  of  $S_1$  and  $S_2$  takes time  $O(mf(\gamma_2))$  time, where  $m$  is the total number of arcs and exit chords in  $S_1$  and  $S_2$ , which from Lemma 2.3, is  $O(n_1/\gamma_1 + n_2/\gamma_2 + 1)$ . Note that among the chords to be included in the fusion  $S$ , we not only have the newly discovered chords that connect  $\partial C_1$  and  $\partial C_2$  as well as the old chords of  $S_1$  and  $S_2$  that still form visible pairs of points with respect to  $C$ ,



but we also have all the null-length chords of  $S_1$  and  $S_2$  as well as the chords incident upon the vertices of  $\partial C$  resulting from the duplication of the vertex  $C_1 \cap C_2$ .

To set up the submap  $S$  in normal form can be done quite easily, now that we have all the chords of  $S$  handy, along with the names of the edges to which they are adjacent (in the way of pointers into the input table). In order to allow canonical vertex enumerations, let us sort the endpoints of these chords along  $\partial C$ , which is done by sorting the names of the edges of  $P$  on which these arcs abut, and then sorting the endpoints falling within the same edges by considering  $y$ -coordinates. This allows us to set up the required arc-sequence table. Note that merging can also be used instead of sorting but this step is not the dominant cost, anyway. With this information it is now straightforward to set up the tree of  $S$  with all the necessary arc-structures and their relevant pointers. Now, we must compute the tree decomposition of  $S$ , which as we mentioned earlier can be done very simply in time  $O((n_1/\gamma_1 + n_2/\gamma_2 + 1) \log(n_1 + n_2))$ . Putting everything together, we derive the upper bound of  $O((n_1/\gamma_1 + n_2/\gamma_2 + 1)(f(\gamma_2) + \log(n_1 + n_2)))$  on the time needed to complete the fusion of  $S_1$  and  $S_2$ . ■

*Conformality:*

Of course, there is no reason to believe that the fusion  $S$  should be conformal. Things can never be too bad, however. Indeed, let  $A_1, A_2, \dots$  be the arcs of a region  $R$  of  $S$  in counterclockwise order and let  $B_1, B_2, \dots$  be a partition of the sequence of arcs into runs, meaning that  $B_j = A_{i_j}, A_{i_j+1}, \dots, A_{i_{j+1}-1}$  is a maximal subsequence of arcs from either  $\partial C_1$  or  $\partial C_2$  (but not both). In the definition of maximal, we regard  $A_1, A_2, \dots$  as a circular sequence. Because any exit chord endpoint of  $S_i$  is still an endpoint of a chord in  $S$  and, with the possible exception of the chords incident to  $a_0$  or  $a_{m+1}$ , every chord of  $S$  that connects two points of  $\partial C_i$  is also a chord of  $S_i$ , it follows by conformality that a run associated with  $\partial C_i$  cannot have more than five exit chords in its midst. Therefore, a run cannot have more than six arcs. On the other hand, it follows from Lemma 2.2 that there are at most two runs. Why is that so? The lemma says that if we walk along  $\partial C$  clockwise we will in effect traverse (among other things) the boundary of  $R$  counterclockwise (minus the exit chords). If we begin our walk at one of the two points of  $\partial C$  corresponding to the vertex  $C_1 \cap C_2$ , we will first exhaust, say,  $\partial C_1$  and then  $\partial C_2$ . Therefore, the counterclockwise traversal of the boundary of  $R$  must exhaust first the runs  $B_i$  contributed by  $S_1$  and then the runs contributed by  $S_2$ . Obviously, this leaves only the possibility of having at most one run of each type, and hence a total of at most two runs. The conclusion to draw is that, although not necessarily conformal, the submap  $S$  has no region with more than a bounded number of arcs. For this reason,  $S$  can be immediately put in normal form. Our goal now is to reduce the number of arcs per region to 4 or less by adding chords into  $S$ , if necessary. We begin with two technical results.

**Lemma 2.6.** *Given two arcs  $A_1$  and  $A_2$  of the same region of  $S$ , we can check whether they have a pair of mutually visible points, one of which is a vertex of  $\partial C$  (meaning that, say,  $A_1$  contains a vertex  $v$  such that the point of  $\partial C$  seen by  $v$  lies in  $A_2$ ) and, if that is the case, find such a pair, all of it in time  $O(f(\gamma_2)g(\gamma_2)(h(\gamma_2) + \log \gamma_2))$ .*

*Proof:* We will concentrate on finding whether, say,  $A_1$  contains a vertex  $v$  such that the point of  $\partial C$  seen by  $v$  lies in  $A_2$ . To begin with, observe that  $A_1$  and  $A_2$  are arcs or subarcs of either  $S_1$  or  $S_2$  but cannot overlap in both  $\partial C_1$  and  $\partial C_2$ . The reason is that we added chords incident upon the vertices of  $\partial C$  resulting from the duplication of the vertex  $C_1 \cap C_2$ . Therefore, their number of edges (for each of them) is at most  $\gamma_2$  (recall that  $\gamma_1 \leq \gamma_2$ ). Using local shooting, we can efficiently check, in time  $O(f(\gamma_2))$ , whether a given vertex of  $A_1$  qualifies as the vertex  $v$  sought. (Again, one must be careful that local shooting reports edges of  $P$  and does not tell us if the point hit is on the desired arc or is the companion of a point of the arc. We already discussed how local checking can decide which way it is in constant time, so we will not make further mention of that minor difficulty.) Of course, we should not check all the vertices of  $A_1$  because there might just be too many of them. Instead, we need to do some kind of binary search among the vertices of  $A_1$ . For that purpose we invoke the appropriate arc-cutting oracle and subdivide  $A_1$  into at most  $g(\gamma_\ell)$  subarcs, for  $\ell = 1$  or  $\ell = 2$ , such that for each subarc  $\alpha$  we have a normal-form  $h(\gamma_\ell)$ -granular conformal submap  $S_\alpha$  of  $V(\bar{\alpha})$  as well as the tree decomposition of  $S_\alpha$ . We will search each subarc in turn, stopping as soon as we find a good vertex, if ever. Since each subarc  $\alpha$  comes with the tree decomposition  $T$  of the submap  $S_\alpha$ , we will be able to check the candidacy of  $\alpha$  in its entirety in  $O(f(\gamma_2)(h(\gamma_2) + \log \gamma_2))$  time, provided that the following test is performed in  $O(f(\gamma_2))$  time: Given a chord  $ab$  of  $S_\alpha$ , either determine that  $a$  or  $b$  is a vertex of  $\alpha$  and sees  $A_2$  (with respect to  $C$ ), or find some  $i$  such that  $\alpha \cap \alpha_i$  is empty or has no point that sees  $A_2$ , where  $\alpha_1$  and  $\alpha_2$  denote the two pieces of  $\partial \bar{\alpha}$  between  $a$  and  $b$ . (Note that  $\partial \bar{\alpha}$  is a superset of  $\alpha$ , with twice the number of vertices.)

Indeed, we begin by applying this test with respect to the chord corresponding to the root of the tree  $T$ . Either we terminate with a positive answer or else we identify the child of the root that corresponds to  $\alpha_j$  ( $j \neq i$ ) and iterate on this process from that node. This will lead us to termination at some internal node of  $T$  or perhaps will take us to the bottom of the tree. Note that determining which node to branch to at each step is trivial once we have identified the  $\alpha_i$  to be rejected. (So, we can perform the test just as stated above without having to “resize”  $\alpha$  to reflect the current status of the ever-shrinking set of candidates.) If we reach a leaf, we examine each vertex of the region associated with it and, among those belonging to  $\alpha$ , we check whether any of them can see  $A_2$ . Since there are only  $O(h(\gamma_2))$  edges in the region and the depth of the tree is  $O(\log \gamma_2)$  the running time of the algorithm is  $O(f(\gamma_2)(h(\gamma_2) + \log \gamma_2))$ , as claimed.

Whenever we discover a successful candidate vertex, the search can obviously be stopped right there. What remains to be seen is why upon reaching a leaf the corresponding region is the only one which can still provide the desired answer. Let us assume that the search ends up at a leaf. At the very beginning, each vertex of  $\partial \bar{\alpha}$  is a candidate, but every time we branch down the tree we eliminate the candidacy of a whole interval of  $\partial \bar{\alpha}$ . At the end, the remaining candidates constitute a set of vertices in  $\partial \bar{\alpha}$  which no node of  $T$ , i.e., no chord of  $S_\alpha$ , can further cut apart. This means that the set consists of vertices in a single region of  $S_\alpha$ , namely, the region associated with the leaf. This proves the correctness of the procedure. So, to summarize, if the basic test stated earlier can be performed in  $O(f(\gamma_2))$  time, then we can solve the entire problem in  $O(f(\gamma_2)g(\gamma_2)(h(\gamma_2) + \log \gamma_2))$  time, which proves the lemma.

Removing  $\partial \bar{\alpha}$  from the spherical plane leaves two open regions, each polygonal and homeomor-

phic to a disk. One of them is infinitesimally small; let  $D$  be the closure of the other one. It is important that  $D$  should be homeomorphic to a closed disk and not to a 2-sphere, so the interior of  $\bar{\alpha}$ , and more generally, of  $C$ , should be understood as being very small but nonempty. Let  $c$  and  $d$  be the endpoints of  $\alpha$  on  $\partial\bar{\alpha}$ . Removing  $c$  and  $d$  from  $\partial C$  leaves  $\alpha$  and a curve  $A$ , both lying in  $D$ , so we have set the stage for Lemma 2.4. Figure 2.11 illustrates the correspondence: the disk  $D$  represents the outside of the snake on the left picture, while it is the inside of the circle on the right. The subarc  $\alpha$  runs between  $c$  and  $d$  and corresponds to  $B_1 \cup B_2$ . To compute  $a'$  and  $b'$  (if defined) can be done by local shooting in  $O(f(\gamma_2))$  time. If ever  $a$  (resp.  $b$ ) is a vertex of  $\alpha$  and  $a'$  (resp.  $b'$ ) belongs to  $A_2$  then obviously we are done and successful in our search, so we can assume that neither does. But, in that case,  $A_2$  lies entirely within one of the connected components of the curve  $A$  after it has been cut up by removing  $a'$  and  $b'$  (if they exist). Therefore, by Lemma 2.4,  $A_2$  must be shielded from some  $B_j$ , which means that it cannot be connected to  $B_j$  without crossing  $ab$  or  $A$ . Furthermore we know that  $B_j$  can be identified in constant time. The key observation now is that  $B_j$  coincides precisely with one of  $\alpha \cap \alpha_1$  or  $\alpha \cap \alpha_2$ , say, the first one. It follows that no point of  $\alpha \cap \alpha_1$  can see  $A_2$ , and the test is completed. ■

**Lemma 2.7.** *Let  $A_1, \dots, A_k$  be the clockwise circular sequence of arcs around a region of  $S$ . If  $k > 4$  then there exist  $i, j$ , such that (i)  $i - j \not\equiv -1, 0, 1 \pmod{k}$  and (ii)  $A_i$  has a vertex that sees  $A_j$  (with respect to  $C$ ).*

*Proof:* Recall that the region is associated with a subtree of the visibility tree of  $C$ . Let us add to this subtree the edges that connect it to the rest of the visibility tree. With respect to this augmented subtree, each exit chord of the region is associated with a distinct node of degree 1 (but the converse may not be true). Note that some of these chords might be of zero length. Consider the Steiner minimal tree of these particular degree-1 nodes (i.e., the smallest connected set of edges that join these nodes together), and form a tree homomorphic to it by ignoring nodes of degree 2. Now embed this new tree in the plane and enclose it by a simple closed curve that connects together all its degree-1 nodes (Figure 2.12). The area inside that curve is partitioned into  $k$  sub-areas, each corresponding to a distinct arc  $A_i$ . It is straightforward to prove that because there are at least five degree-1 nodes, the maximum node-degree is 4, and no node has degree 2, then at least one edge of the tree must be incident upon two sub-areas associated with  $A_i$  and  $A_j$ , respectively, where  $i - j \not\equiv -1, 0, 1 \pmod{k}$ . Of all the chords in the region only the exit chords might fail to be incident upon at least one vertex of  $\partial C$  (indeed, recall that the fusion algorithm may introduce chords incident upon no vertices of  $\partial C$ ), therefore the edge in question corresponds to one or several chords of  $V(C)$  (in its original, non-augmented state, i.e., without the extra chords found during fusion) that connect  $A_i$  and  $A_j$ . ■

Equipped with the two previous lemmas, making  $S$  conformal is now quite easy. Recall that although  $S$  might not be conformal, none of its regions has more than a bounded number of arcs. For any region with more than four arcs, let us apply Lemma 2.6 to every pair of nonconsecutive arcs until we find a chord which we can add to  $S$ , all the while keeping  $S$  in normal form. We iterate on this process until no region has more than four arcs. Note that although  $S$  keeps changing, Lemma 2.6 always holds since region arcs can only become smaller. Lemma 2.7 tells us that this

chord addition process will not stop until  $S$  becomes conformal. Since the total number of arcs in  $S$  is  $O(n_1/\gamma_1 + n_2/\gamma_2 + 1)$  arcs, we conclude:

**Lemma 2.8.** *The submap  $S$  can be made conformal in time  $O((n_1/\gamma_1 + n_2/\gamma_2 + 1)f(\gamma_2)g(\gamma_2)(h(\gamma_2) + \log \gamma_2))$ .*

*Granularity:*

Since by making  $S$  conformal we did not remove any exit chord, it is still the case that, as observed in the proof of Lemma 2.6, no arc has more than  $\gamma_2$  edges. Therefore,  $S$  is conformal and  $\gamma_2$ -semigranular. We must now check whether the second criterion for  $\gamma_2$ -granularity holds. This criterion says that contracting any edge of the submap tree that is adjacent to at least one node of degree less than 3 produces a new node whose weight exceeds  $\gamma_2$ . This is very easy to enforce: If an edge does not pass the test, we just contract it by removing its corresponding exit chord (and those endpoints that are not vertices of  $\partial C$ ). Note that this will not cause a violation of the first criterion, since the size of all the arcs will always remain within  $\gamma_2$ . We process each exit chord in turn and check whether it is removable. Chords need be processed only once since the removals cannot make any chord removable if it was not already so before. Therefore,  $\gamma_2$ -granularity, and more generally  $\gamma$ -granularity, for any  $\gamma \geq \gamma_2$ , can be enforced in this nondeterministic fashion in time linear in the size of the submap tree, that is,  $O(n_1/\gamma_1 + n_2/\gamma_2 + 1)$ . From Lemmas 2.5 and 2.8 we then derive the following result.

**Lemma 2.9.** *Let  $C_1$  and  $C_2$  be two polygonal curves of  $n_1$  and  $n_2$  vertices respectively, whose union forms a connected vertex-to-vertex piece of the input (simple and nonclosed) polygonal curve  $P$ . Suppose that we are given a normal-form  $\gamma_i$ -granular conformal submap of each  $V(C_i)$ , where  $\gamma_1 \leq \gamma_2$ , along with their ray-shooting oracles. Then, for any  $\gamma \geq \gamma_2$ , it is possible to compute a normal-form  $\gamma$ -granular conformal submap of  $V(C)$  in time  $O((n_1/\gamma_1 + n_2/\gamma_2 + 1)f(\gamma_2)g(\gamma_2)(h(\gamma_2) + \log(n_1 + n_2)))$ .*

Before we close this section, we wish to mention a slight extension of our merging algorithm. Suppose that  $R$  is a region of some conformal submap. In order to make its boundary into a simple nonclosed curve, we break it open by removing one vertex. Let  $C$  be the resulting curve. If  $C_1$  and  $C_2$  form a partition of  $C$  and we have  $\gamma_i$ -granular conformal submaps of each  $V(C_i)$  as well as the relevant oracles, then we want to argue that Lemma 2.9 can be applied verbatim. The only difficulty here is that  $C$  may now have up to four horizontal edges, which were previously banished. By tilting the edges ever so slightly if need be or just adapting our algorithms directly, however, it should be clear that the techniques leading to Lemma 2.9 are still applicable verbatim. This will be useful in Section 4 when we describe the down-phase of the visibility algorithm.

### 3. Toward Efficient Oracles

Of the two oracles defined earlier, the ray-shooter is the more challenging to implement; the reason being that it addresses the key issue in the triangulation business, which is the discovery of new chords. Fast planar point location could be the answer. But traditional methods, e.g., [8,18], are inadequate for several reasons, the most crucial of which is their inability to support merging in sublinear time. We turn this problem around by exploiting further the approximation scheme provided by the concept of granularity.

Let  $C$  be a connected vertex-to-vertex piece of the input polygonal curve  $P$  and let  $m$  be its number of vertices. Let  $S$  be a normal-form  $\gamma$ -granular conformal submap of  $V(C)$ . So far, we have focused mostly on the tree structure of  $S$ . But let us now regard  $S$  as a planar graph. For this purpose, we must temporarily forget the fact that  $C$  has been given a double boundary. We define  $S^*$  to be the planar subdivision obtained by taking  $S$  and making every vertex (vertices of  $\partial C$  and chord endpoints) a vertex on *both* sides of the double boundary. As a result, the edges of  $S^*$  might be smaller than those of  $S$  but, unlike in  $S$ , no edge of  $S^*$  is of zero length (zero-length edges are now “contracted” into vertices). More important, each face of  $S^*$  coincides exactly with a distinct region of  $S$ , except for the fact that it might have many more vertices incident upon it. Indeed, a region’s only vertices are the endpoints of its own exit chords along with some vertices of  $\partial C$ , whereas the vertices of a face include all of the above plus all the chord endpoints that abut on the corresponding region from the outside. Note that the correspondence face/region is not bijective because empty regions have no associated faces. Besides being planar, the graph  $S^*$  has two remarkable properties: (i) from Lemma 2.3 we know that it has  $O(m/\gamma + 1)$  faces, which is much smaller than the number of vertices (when  $\gamma$  is large), and (2) although a given face might be very complex (i.e., incident upon many edges) its number of noncollinear edges is small, i.e.,  $O(\gamma)$ . These two features will allow us to implement an efficient ray-shooting oracle.

Let  $G$  be the dual graph of  $S^*$ , that is, the graph obtained by associating a distinct node with each face of  $S^*$  and connecting two nodes if and only if they are distinct and their corresponding faces share a common edge. It is a classical result of graph theory that  $G$  is planar. How hard is it to compute  $G$ , say, in the form of adjacency lists? Two faces are adjacent if and only if either they share a chord or one of them has a chord endpoint that abuts on a non-null length arc of the region associated with the other face. The first type of adjacencies can be detected immediately from  $S$ . The latter can be done by double identification, as discussed in Section 2.C. It is then clear that  $G$  can be computed in  $O(\mu \log m)$  time, where  $\mu$  is the number of nodes in  $G$ .

If  $\mu = 1$  then ray-shooting can be done trivially in  $O(m)$  time, so let us assume that  $\mu > 1$ . The planarity of  $G$  will work wonders for us. The first payoff is that the number of edges is at most  $3\mu - 6$ . The second reward is that we can apply the linear-time algorithm of Lipton and Tarjan [22] to find a good separator. This partitions the nodes of  $G$  into three subsets  $A, B, C$ , such that (i) no edge joins a node of  $A$  with a node of  $B$ , (ii) neither  $A$  nor  $B$  contains more than  $2\mu/3$  nodes, and (iii)  $C$  contains at most  $\sqrt{8\mu}$  nodes. Let  $G_A$  (resp.  $G_B$ ) be the graph obtained by keeping only the nodes of  $A$  (resp.  $B$ ) and the edges of  $G$  that join only nodes of  $A$  (resp.  $B$ ). We repeat the procedure over with respect to each of  $G_A$  and  $G_B$  and iterate in this fashion until none of the

graphs have more than  $\mu^\delta$  nodes, for some fixed  $\delta$  ( $0 < \delta < 1$ ). Let  $C^*$  be the set of all separators, i.e., the union of all the  $C$ 's. We easily verify that  $C^* = O(\mu^\delta)$ , provided that  $\delta$  is chosen large enough, e.g.,  $\delta \geq 2/3$  [21]. In  $O(\mu \log \mu)$  time we can compute  $C^*$  and partition the remaining nodes into subsets  $D_1, D_2$ , etc., each of size at most  $\mu^\delta$ , such that no path of  $G$  can join two nodes in distinct subsets without passing through a node of  $C^*$ .

What is the utility of  $C^*$  for ray-shooting? Take a vertical line passing to the right of all the vertices of  $P$ , and intersect it with the chords of the regions in  $S$ . This breaks up the line into segments, every one of which falls entirely within some region; to split up the line and identify the regions cut by each segment takes  $O(\mu)$  time. We now claim that ray-shooting from any point can be done in  $O(\gamma\mu^\delta)$  time. Our first task is to shoot within each region that is dual to a node of  $C^*$ , using a naive algorithm which involves checking all the  $O(\gamma)$  edges of the region (and not the edges of the face, which might be much more numerous!). Assume that the ray of light hits a point  $p$ , meaning that  $p$  is the first point hit by the ray among all the edges of the regions dual to the nodes of  $C^*$ . Let  $R$  be the last region of  $S$  traversed before the hit. To identify  $R$  can be done by double identification, followed by checking on the local orientation of the hit. If  $R$  is a region dual to a node  $v$  of  $C^*$  then the starting point of the ray lies in  $R$  (otherwise an earlier hit would have been detected) and we are trivially done. Otherwise, let  $R'$  be the region incident upon the (region) edge containing the point of  $\partial C$  hit by the ray-shooting. (This is the region that we are looking for.) If  $R$  and  $R'$  are not the same then the two regions can be connected by a horizontal line segment that avoids all the regions dual to  $C^*$ . It follows that the node  $w$  associated with  $R'$  can be reached by a path from  $v$  that avoids  $C^*$ . Consequently,  $v$  and  $w$  both lie in the same  $D_i$ . We can find  $w$ , and from there, answer the ray-shooting query, by first finding  $D_i$ , which takes constant time since we know  $R$ , and then naively checking all the regions dual to nodes in  $D_i$ , which takes  $O(\gamma\mu^\delta)$  time. Returning to our earlier case-analysis, assume now that the ray of light hits no region dual to a node in  $C^*$ . Then the ray-shooting takes place entirely within the regions dual to the nodes of a single  $D_i$ . To find out which one, we shoot toward the vertical line and find which segment of the line is hit. This takes  $O(\log \mu)$  time by binary search. We can now identify the region  $R$  immediately. The remainder of the algorithm is unchanged. We conclude that ray-shooting can be done in  $O(m)$  time if  $\mu = 1$ , and  $O(\gamma\mu^\delta)$  time if  $\mu > 1$ , after  $O(\mu \log m)$  preprocessing.

**Lemma 3.1.** *Let  $C$  be a connected vertex-to-vertex piece of the input polygonal curve  $P$  and let  $m$  be its number of vertices. Let  $S$  be a normal-form  $\gamma$ -granular conformal submap of  $V(C)$ . Then it is possible to preprocess  $S$  in  $O(m(\log m)/\gamma + 1)$  time so that ray-shooting within  $S$  can be done in time  $O(\gamma^{1-\delta}m^\delta)$ , for some fixed  $\delta < 1$ .*

#### 4. The Visibility Algorithm

Let  $P$  be a simple nonclosed polygonal curve with  $n$  vertices  $v_1, \dots, v_n$ . By padding the curve with additional vertices, if necessary, we can assume that  $n = 2^p + 1$ . Any subcurve of  $P$  of the form  $v_a, \dots, v_b$ , where  $a - 1$  is a multiple of  $2^\lambda$  and  $b - a = 2^\lambda$  is called a *chain in grade*  $\lambda$ . Obviously, (i)

a grade- $\lambda$  chain has  $2^\lambda + 1$  vertices, (ii) there are  $2^{p-\lambda}$  chains in grade  $\lambda$ , and (iii) there are  $p + 1$  grades:  $0, 1, \dots, p$ . Put

$$\beta = \frac{1 - \delta}{4 - 2\delta}.$$

Note that  $\beta$  lies strictly between 0 and 1. We begin our work bottom-up, computing submaps of granularity roughly  $m^\beta$ , where  $m$  is the size of the underlying curve. We pursue this task until the submap for the whole polygon has been obtained. Then we reverse the process and work top-down until the submap has been completely refined into its full-fledged visibility map.

**A. The up-phase.** We begin with a piece of terminology: Given a curve  $C$  consisting of  $m$  contiguous edges of  $P$ , we say that a submap of  $V(C)$  is *canonical* if it is  $2^{\lceil \beta \lambda \rceil}$ -granular, conformal, and represented in normal form, where  $\lambda = \lceil \log m \rceil$  (logarithm to the base 2). Note that a canonical submap for a chain in grade  $\lambda$  is  $2^{\lceil \beta \lambda \rceil}$ -granular. For each grade  $\lambda = 0, 1, \dots, p$ , in that order, we compute a canonical submap of  $V(C)$  for each chain  $C$  in that grade. We preprocess each submap for ray-shooting along the lines of Lemma 3.1, setting  $\gamma$  to the value  $2^{\lceil \beta \lambda \rceil}$ , and we finally compute its tree decomposition. This work can be done naively for the early grades, so let us pick up the action at a grade  $\lambda$  larger than some appropriate constant, assuming that all grades less than  $\lambda$  have been processed already. We need the following result.

**Lemma 4.1.** *Given any portion  $D$  of  $P$  of the form  $v_a, \dots, v_b$ , where  $2^{\lambda-1} < b - a \leq 2^\lambda$ , we can compute a canonical submap of  $V(D)$  in time proportional to  $\lambda^2 (\log \lambda) 2^{\lambda(1 + \beta(\delta-1) + \beta^2(2-\delta))}$ .*

*Proof:* In  $O(\lambda)$  time we can partition  $D$  into  $j \leq 2\lambda$  chains,  $D_1, \dots, D_j$ , in grades less than  $\lambda$ . This means that, for each  $i = 1, \dots, j$ , a canonical submap  $S_i$  of  $V(D_i)$  is available. Let  $\gamma = 2^{\lceil \beta \lambda \rceil}$  be the granularity of a canonical submap of  $V(D)$ . Since the granularity of canonical submaps grows monotonically with the size of the underlying polygonal curve, we can trivially reset the granularity of each  $S_i$  to  $\gamma$  (Section 2). The time to do that is proportional to the total number of chords in all the  $S_i$ 's which, from Lemma 2.3, is on the order of  $\sum_{0 \leq k < \lambda} 2^{k - \lceil \beta k \rceil}$ , that is,  $O(2^{\lambda(1-\beta)})$ .

Let us now merge these submaps two-by-two ( $D_1$  with  $D_2$ ,  $D_3$  with  $D_4$ , etc.). More generally, we consider a complete binary tree whose leaves are in bijection with the  $D_i$ 's and we merge submaps bottom-up by following the tree pattern. Application of Lemma 2.9 results in a canonical submap of  $V(D)$  provided, of course, that the required oracles are available. But are they? Notice that during any merge any arc  $\alpha$  in either of the two input submaps consists of at most  $\gamma$  edges. Therefore, any subarc  $\alpha' \subseteq \alpha$  can be subdivided into a constant number of contiguous pieces whose corresponding portions of  $P$  consist of single line segments and vertex-to-vertex pieces of  $P$ , each with at most  $2^{\lceil \beta \lambda \rceil}$  edges. Each of these pieces can be partitioned into a collection of  $O(\lambda)$  chains in grades at most  $\lceil \beta \lambda \rceil$ . Our work at previous grades ensures that we have ray-shooting structures for each of these chains. Therefore, we can implement the ray-shooting oracles with the following specification:

$$f(\gamma) = \lambda 2^{\lceil \beta \lambda \rceil \delta + \lceil \beta \lceil \beta \lambda \rceil \rceil (1-\delta)}.$$

For the same reasons, the arc-cutting oracles can be implemented so that,

$$g(\gamma) = O(\lambda),$$

and

$$h(\gamma) \leq 2^{\lceil \beta \lceil \beta \lambda \rceil \rceil}.$$

There are  $O(\log \lambda)$  levels of merging to be performed, each incurring a cost at most proportional to

$$\lambda^2 2^{\lambda(1+\beta(\delta-1)+\beta^2(2-\delta))}$$

by virtue of Lemma 2.9. Since the initial cost of resetting the granularity is only  $O(2^{\lambda(1-\beta)})$ , the lemma follows readily. ■

Let us turn now to the processing of grade  $\lambda$ . Lemma 4.1 can be called upon to compute a canonical submap of the visibility map of each chain in grade  $\lambda$ . Preprocessing each chain for ray-shooting can be done by using Lemma 3.1, while as we pointed out earlier, computing the tree decomposition is easily done in time  $O(m \log m)$ , where  $m$  is the size of the tree. Putting all these facts together, we conclude that processing grade  $\lambda$  requires time at most proportional to

$$n\lambda^2(\log \lambda)2^{\beta\lambda(\beta(2-\delta)+\delta-1)} + n\lambda 2^{-\beta\lambda}.$$

From our choice of  $\beta$ , we find that the time is at most proportional to

$$n\lambda^2(\log \lambda)2^{-\beta^2\lambda(2-\delta)} + n\lambda 2^{-\beta\lambda},$$

therefore processing all  $p + 1$  grades, and thereby completing the up-phase, takes linear time.

**B. The down-phase.** Now that we have canonical submaps for each chain in each grade, along with the oracle structures and the tree decompositions, we are ready to refine the canonical submap of  $V(P)$ . This will be done incrementally by going down the tree. The following lemma provides the key to the algorithm:

**Lemma 4.2.** *Let  $C$  be an arbitrary chain in any grade  $\ell \geq \lambda > 0$  and suppose that a  $2^{\lceil \beta \lambda \rceil}$ -granular conformal submap of  $V(C)$  is available in normal form. Then it is possible to compute  $V(C)$  in time at most  $(c - 1/\lambda)2^\ell$ , where  $c$  is some constant large enough.*

*Proof:* We proceed by induction on  $\lambda$ . Let  $S$  be the  $2^{\lceil \beta \lambda \rceil}$ -granular conformal submap of  $V(C)$ . The case where  $\lambda$  is a constant is trivial since the regions of  $S$  have bounded size, and therefore the missing chords can be provided in constant time per region. So, let us switch directly to the inductive case, assuming that  $\lambda$  is large enough. Let  $R$  be a region of  $S$ . Because of conformality, the union of all the arcs of  $R$  can be partitioned into a constant number of single edges and vertex-to-vertex pieces of  $\partial C$  with at most  $2^{\lceil \beta \lambda \rceil}$  edges. Applying Lemma 4.1, we can compute a canonical submap for each connected polygonal piece in the partition in time at most proportional to

$$\lambda^2(\log \lambda)2^{\beta\lambda(1+\beta(\delta-1)+\beta^2(2-\delta))}.$$

Each of these submaps has granularity at most  $2^{\lceil \beta \lceil \beta \lambda \rceil \rceil}$ , so we can pursue the merging by putting together all these submaps and thus create a single normal-form  $2^{\lceil \beta \lceil \beta \lambda \rceil \rceil}$ -granular conformal submap of  $V(R^*)$ , where  $R^*$  is the boundary of  $R$  minus a vertex (to ensure that it is nonclosed). For the consistency of our discussion, we should regard  $R^*$  as a standard polygonal curve and *not* as part of a double boundary. The operation requires a constant number of merges, so we can carry it out effectively by replicating the argument of Lemma 4.1.



There is a small subtlety in this last round of merges, so we will say a few more words. To take a simple example, suppose that  $R^*$  has two arcs and two exit chords:  $\alpha_1, a_1b_1, \alpha_2, a_2b_2$ , in cyclic order. Recall that when we merge two submaps we first compute the chords incident upon the vertices emanating from the intersection of the two input curves. Here,  $\alpha_1$  and  $\alpha_2$  need not have any point in common, so we use  $a_1b_1$  instead as one of these starting chords. The only difference on which we insist is that in the last round of merges the granularity enforcement procedure be skipped. (We will take care of this later.) The reason for this shortcut is that the exit chords along  $R^*$  should *not* be removed (at least not now), so we don't want to take chances. In this way,  $R^*$  corresponds to a cycle in the resulting conformal submap. Note that because the exit chords along  $R^*$  remain valid chords after the merges, conformality is guaranteed. Applying the entire procedure over each region of  $S$  requires time at most proportional to

$$2^\ell \lambda^2 (\log \lambda) 2^{\beta^2 \lambda (\beta(2-\delta) + \delta - 1)}.$$

We can now extract the relevant information, i.e., the exit chords within each region  $R$ . This involves checking the exit chords of the computed submap of  $V(R^*)$  and keeping only those whose endpoints lie on the arcs of the region  $R$ . This leads to a new map  $S^*$  of  $V(C)$  which is a refinement of  $S$ : all its arcs originate from the previous merges, therefore  $S^*$  is a  $2^{\lceil \beta \lceil \beta \lambda \rceil}$ -semigranular conformal submap of  $V(C)$ . It is now time to ensure  $2^{\lceil \beta \lceil \beta \lambda \rceil}$ -granularity. This is done by removing exit chords if necessary, which, as we saw in Section 2, takes time linear in the number of exit chords in  $S^*$ .

For  $\lambda$  large enough, we have  $\lceil \beta \lambda \rceil \leq \lambda - 1$ , so we can apply the induction hypothesis and derive  $V(C)$  from  $S^*$  in time at most  $(c - 1/(\lambda - 1))2^\ell$ . Putting everything together, the total running time for the construction of  $V(C)$  is at most

$$a2^\ell \lambda^2 (\log \lambda) 2^{\beta^2 \lambda (\beta(2-\delta) + \delta - 1)} + \left(c - \frac{1}{\lambda - 1}\right) 2^\ell,$$

for some constant  $a > 0$ , which is no more than

$$a2^\ell \lambda^2 (\log \lambda) 2^{-\beta^3 \lambda (2-\delta)} + \left(c - \frac{1}{\lambda - 1}\right) 2^\ell \leq \left(c - \frac{1}{\lambda}\right) 2^\ell,$$

for  $\lambda$  large enough. ■

From the up-phase we have a  $2^{\lceil \beta p \rceil}$ -granular conformal submap of  $V(P)$ , therefore Lemma 4.2 tells us that  $V(P)$  can be obtained in linear time. As demonstrated in [6,10], a triangulation can be derived from the visibility map in linear time, so our main goal has been reached.

**Theorem 4.3.** *It is possible to compute the visibility map of a simple polygonal curve, and hence, a triangulation of a simple polygon, in linear time.*

## 5. Some Applications and a Few Parting Words

There is a plethora of applications for a fast triangulation algorithm. We will mention only two of them and refer the reader to [1,12,13,15,19,24,25,27] for some pointers to other applications. Kirkpatrick [18] and Edelsbrunner, Guibas, Stolfi [8] have given optimal planar point location algorithms which require linear preprocessing time provided that the planar subdivision has triangulated or monotone faces. With our algorithm the preprocessing can be made linear as long as the graph is connected.

The second application follows from the observation that the triangulation algorithm can be adapted to check whether a polygonal curve is simple. The idea is to compute its visibility map and check that it is, indeed, a valid map. To do so, we verify that each face is a trapezoid or a triangle incident upon edges of the polygon and, conversely, that each vertex and each edge are incident upon the right number of faces, that these incident faces form a region homeomorphic to a disk, and that the obvious local rules (e.g. nonhorizontal edges lie on the polygon's boundary) are satisfied. In particular, we can verify that each point of the spherical plane belongs to at least one region and that, with respect to each vertex and each edge, incident faces are non-overlapping. If the map passes all those tests then we claim that it is, indeed, a valid visibility map and therefore the polygon is simple. Why is that so? To make things easier, we can begin by taking some point  $p$  and checking whether more than one face overlap at that point. If yes, the polygon is not simple. Otherwise, assume that the map passes all the tests but that two or more faces overlap in a point  $q$ . Then the segment  $pq$  must contain at least one point such that exactly one region resides locally on one side along  $pq$  and at least two of them overlap on the other side. This implies that  $q$  lies on an edge of the subdivision that is incident upon either only one face or at least two overlapping faces. Either way, this contradicts having passed all the local consistency checks. We conclude that testing simplicity requires only linear time.

From this we easily derive a new result: testing whether two simple polygons intersect in linear time. To do that, take the highest points of both polygons and shoot horizontally from the lower point  $p$  toward the other polygon in both directions. If there is no hit, then no intersection can occur. Otherwise, we can immediately infer from the local orientation of the hits whether  $p$  lies inside or outside the polygon. In the first case we conclude with a positive answer, else we connect the two polygonal boundaries into a single one by adding one of the connecting segment and duplicating it, thus reducing the problem to that of testing simplicity.

It is interesting to notice that our algorithm can be used to perform Jordan sorting in linear time. This provides a completely different alternative to Hoffman et al.'s algorithm [17]. Recall that the problem is to sort a sequence of numbers  $a_1, \dots, a_n$ , which correspond to the intersection points of a Jordan curve with the  $x$ -axis. Assuming that we do not know anything else about the curve, we construct the polygon  $P$  with vertices  $a_1, b_1, a_2, b_2, \dots, a_n, b_n$ , where

$$b_i = \left( \frac{a_i + a_{i+1}}{2}, \frac{a_{i+1} - a_i}{2} \right).$$

It is immediate that  $P$  is simple and that its visibility map immediately gives us the  $a_i$ 's in sorted order. It is instructive to compare the two methods: Hoffman et al.'s algorithm is on-line and

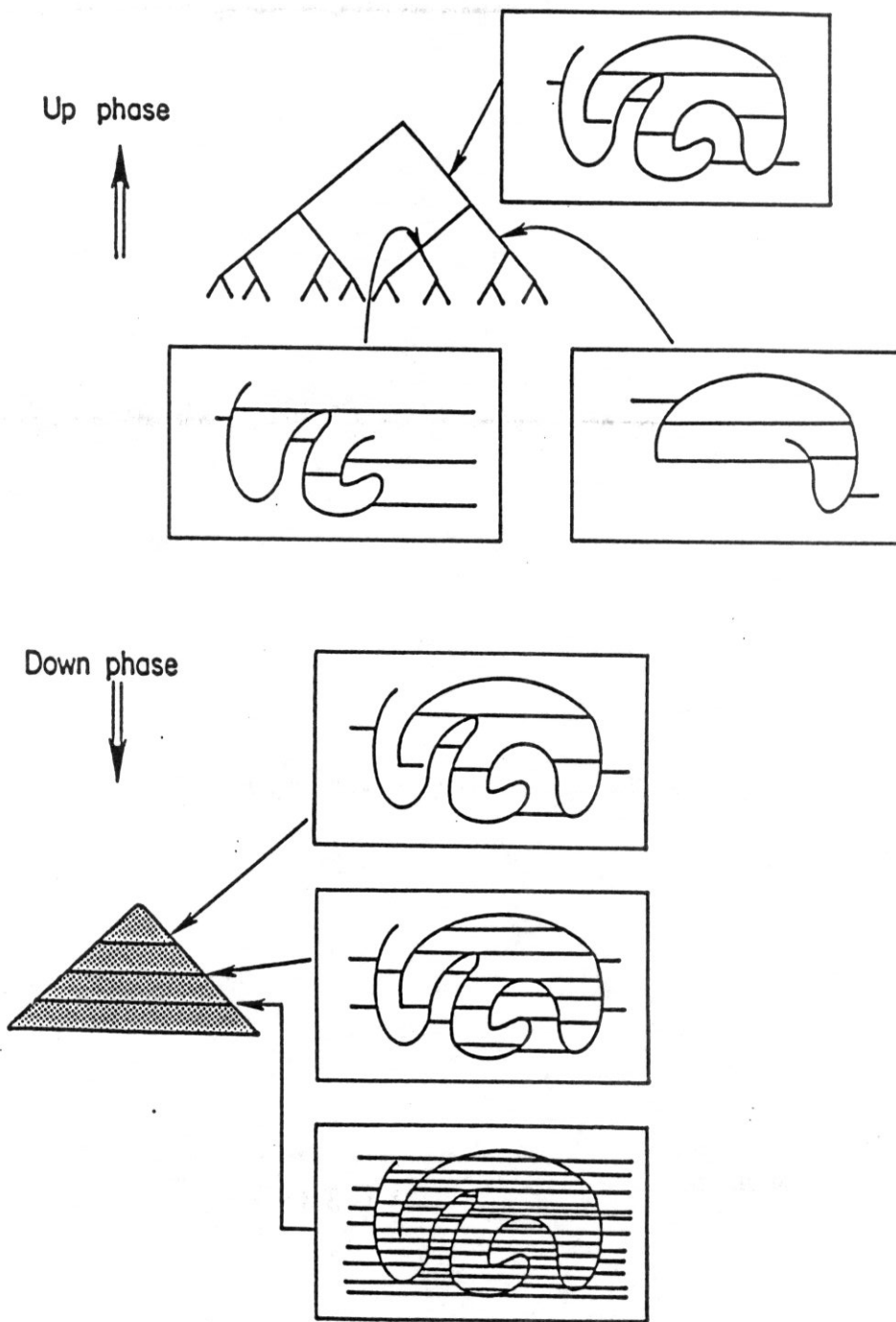
corresponds to an optimal enumeration scheme for Jordan permutations. Our method works off-line and uses divide-and-conquer. It is an intriguing open question whether triangulation can be done on-line in a manner similar to [17]. Tarjan and Van Wyk's method (almost) falls in that category but it is not optimal. Can their algorithm be made optimal? More generally, is it possible to maintain the visibility map optimally under on-line insertion of new edges? Obviously not in an explicit fashion, since a new edge can cut through a linear number of edges, hence creating a quadratic blowup. Even implicitly, however, it can be trivially shown that identifying the order types of the visibility maps of all the prefixes of a simple  $n$ -vertex polygonal curve requires  $\Theta(n \log n)$  bits, which is bad news. Thus, something with less information contents should be maintained by an optimal on-line algorithm. But what? This question might seem rhetorical in light of our linear-time algorithm, but the underlying issue is whether a fundamentally different optimal algorithm exists, in particular, one that works on-line. Also, is there a simple optimal probabilistic triangulation algorithm, say, as simple as Clarkson et al.'s [7]? Our algorithm can be modified by replacing the planar separator step by a randomized construction, but whether the resulting algorithm is really simpler is debatable. Thus the existence of a truly simple linear triangulation algorithm remains an open question. We believe that such an algorithm is unlikely.

## REFERENCES

---

1. Aggarwal, A., Wein, J. *Computational Geometry*, MIT Technical Report, MIT/LCS/RSS 3, August 1988.
2. Baumgart, B.G. *A polyhedron representation for computer vision*, 1975 National Computer Conference, AFIPS Conf. Proceedings, 44, AFIPS Press (1975), 589-596.
3. Booth, H., Tarjan, R.E. *Fast algorithms for some problems on trees and graphs*, PhD Thesis, Princeton University, 1990.
4. Chazelle, B. *A theorem on polygon cutting with applications*, Proc. 23rd Annu. IEEE Symp. on Found. of Comput. Sci. (1982), 339-349.
5. Chazelle, B., Guibas, L.J. *Visibility and intersection problems in plane geometry*, Disc. and Comput. Geom. 4 (1989), 551-581.
6. Chazelle, B., Incerpi, J. *Triangulation and shape-complexity*, ACM Trans. on Graphics 3 (1984), 135-152.
7. Clarkson, K., Tarjan, R.E., Van Wyk, C.J. *A fast Las Vegas algorithm for triangulating a simple polygon*, Disc. and Comput. Geom. 4(5) (1989), 432-432.
8. Edelsbrunner, H., Guibas, L.J., Stolfi, J. *Optimal point location in a monotone subdivision*, SIAM J. Comput. 15 (1986), 317-340.
9. Edelsbrunner, H., Mücke, E.P. *Simulation of simplicity: a technique to cope with degenerate cases in geometric algorithms*, Proc. 4th Ann. ACM Symp. Comput. Geom. (1988), 118-133.
10. Fournier, A., Montuno, D.Y. *Triangulating simple polygons and equivalent problems*, ACM Trans. on Graphics 3 (1984), 153-174.
11. Garey, M.R., Johnson, D.S., Preparata, F.P., Tarjan, R.E. *Triangulating a simple polygon*, Inform. Process. Lett. 7 (1978), 175-180.
12. Guibas, L.J., Hershberger, J. *Optimal shortest path queries in a simple polygon*, Proc. 3rd Ann. ACM Symp. Comput. Geom. (1987), 50-63.
13. Guibas, L.J., Hershberger, J., Leven, D., Sharir, M., Tarjan, R.E. *Linear time algorithms for visibility and shortest path problems inside triangulated simple polygons*, Algorithmica 2 (1987), 209-233.
14. Guibas, L.J., Stolfi, J. *Primitives for the manipulation of general subdivisions and the computation of Voronoi diagrams*, ACM Trans. Graphics 4 (1985), 75-123.
15. Hershberger, J. *Finding the visibility graph of a simple polygon in time proportional to its size*, Algorithmica 4 (1989), 141-155.
16. Hertel, S., Mehlhorn, K. *Fast triangulation of a simple polygon*, Proc. Conf. Found. Comput. Theory, New York, Lecture Notes on Computer Science 158 (1983), 207-218.
17. Hoffman, K., Mehlhorn, K., Rosenstiehl, P., Tarjan, R. *Sorting Jordan sequences in linear time using level-linked search trees*, Inform. and Control 68 (1986), 170-184.
18. Kirkpatrick, D.G. *Optimal search in planar subdivisions*, SIAM J. Comput. 12 (1983), 28-35.

19. Kirkpatrick, D.G., Klawe, M.M., Tarjan, R.E.  *$O(n \log \log n)$  polygon triangulation with simple data structures*, Proc. 6th Ann. ACM Symp. Comput. Geom. (1990), forthcoming.
20. Klawe, M.M. *Private communication*, Jan. 1990.
21. Lipton, R.J., Tarjan, R.E. *Applications of a planar separator theorem*, SIAM J. Comput. 9 (1980), 615–627. Prelim. version in Proc. 18th Annu. IEEE Symp. on Found. of Comput. Sci. (1977), 162–170.
22. Lipton, R.J., Tarjan, R.E. *A separator theorem for planar graphs*, SIAM J. Comput. 36 (1979), 177–189.
23. Muller D.E., Preparata, F.P. *Finding the intersection of two convex polyhedra*, Theoret. Comput. Sci. 7 (1978), 217–236.
24. O'Rourke, J. *Art Gallery Theorems and Algorithms*, Oxford University Press, New York, Oxford, 1987.
25. Suri, S. *A linear time algorithm for minimum link paths inside a simple polygon*, J. Comput. Vision Graphics and Image Process. 35 (1986), 99–110.
26. Tarjan, R.E., Van Wyk, C.J. *An  $O(n \log \log n)$ -time algorithm for triangulating a simple polygon*, SIAM J. Comput. 17 (1988), 143–178.
27. Toussaint, G. *Computational Morphology*, North-Holland, 1988.
28. Toussaint, G. *An output-complexity-sensitive polygon triangulation algorithm*, Report SICS-86.3, McGill University, Montreal, 1988.
29. Yap, C.K. *A geometric consistency theorem for a symbolic perturbation scheme*, Proc. 4th Ann. ACM Symp. Comput. Geom. (1988), 134–142.



**FIGURE 1.1**

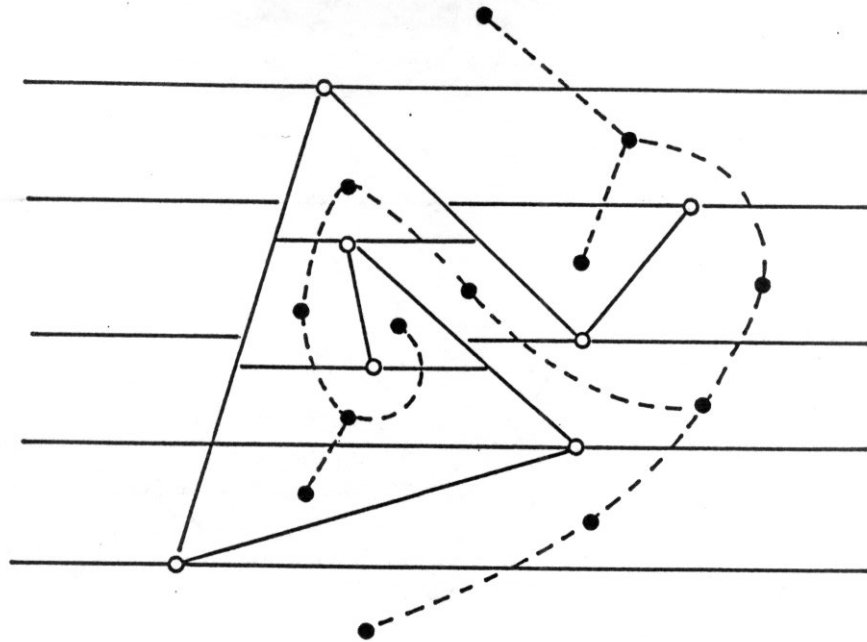


FIGURE 2.1

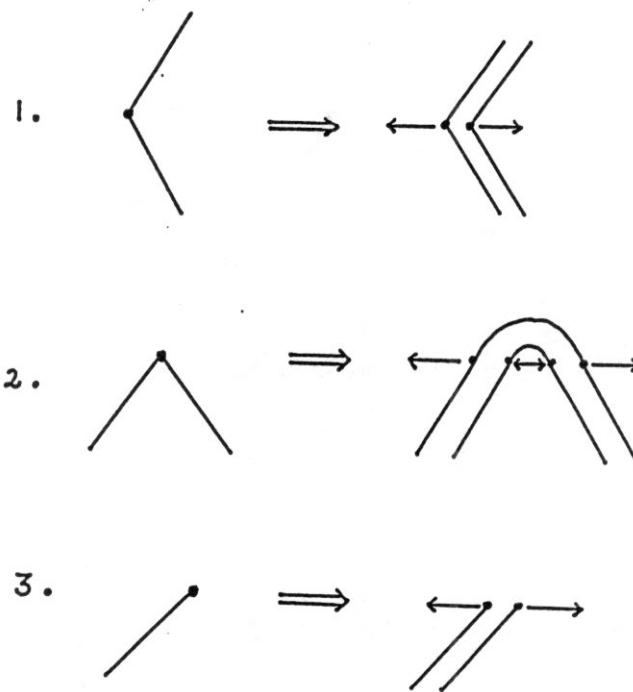


FIGURE 2.2

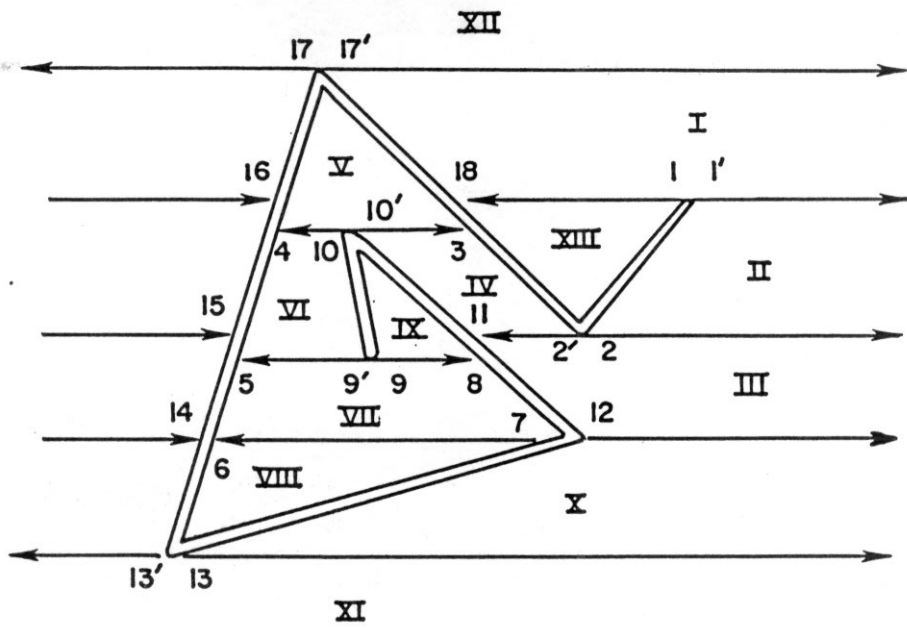


FIGURE 2.3

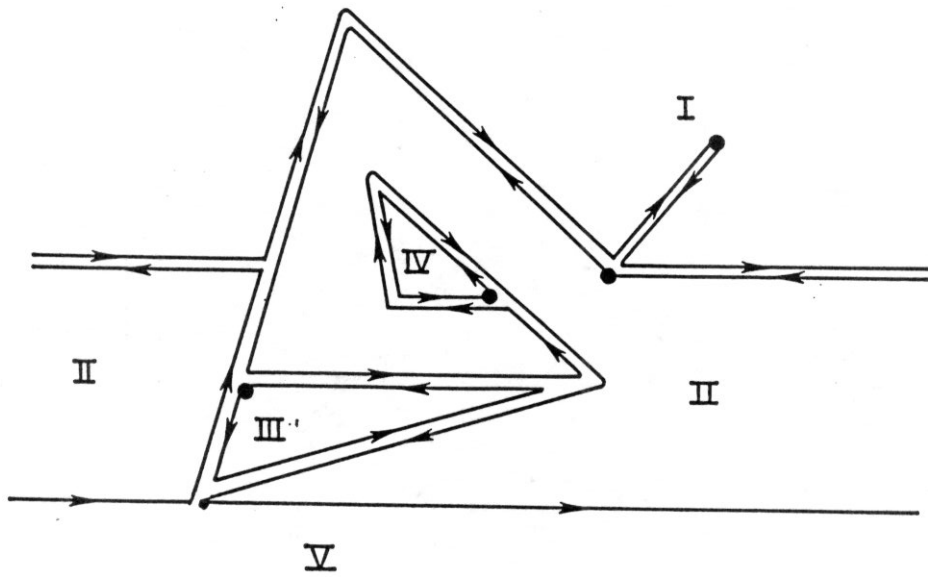


FIGURE 2.4



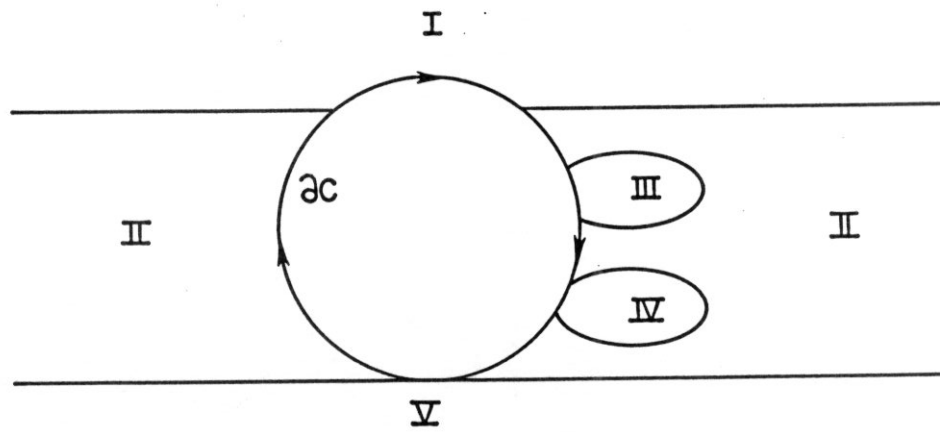


FIGURE 2.5

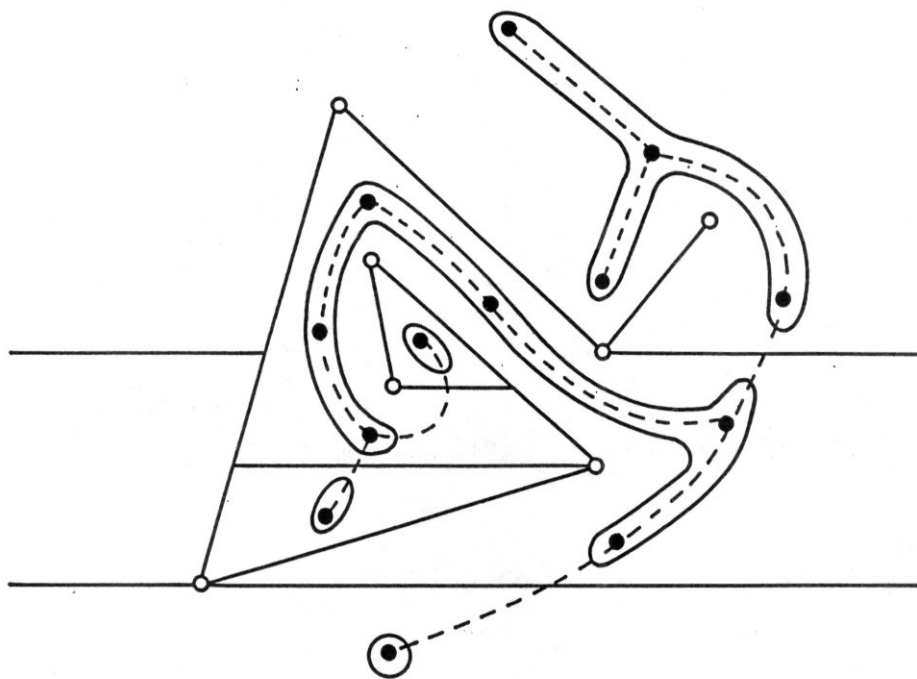


FIGURE 2.6

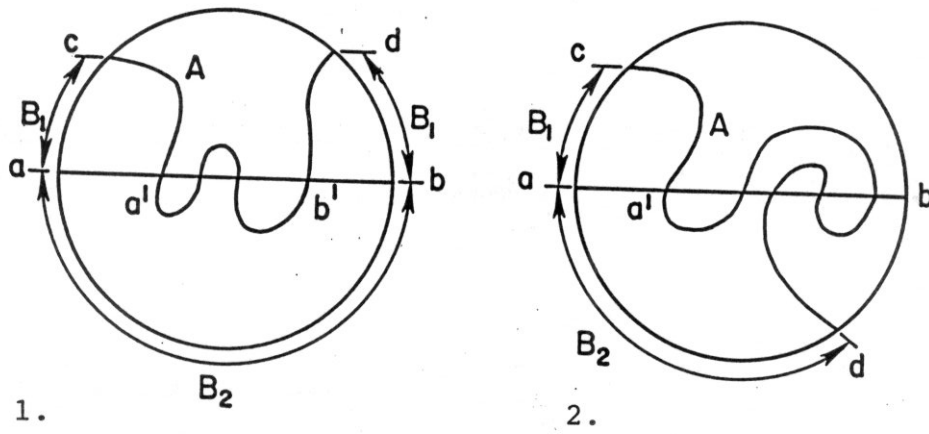


FIGURE 2.7

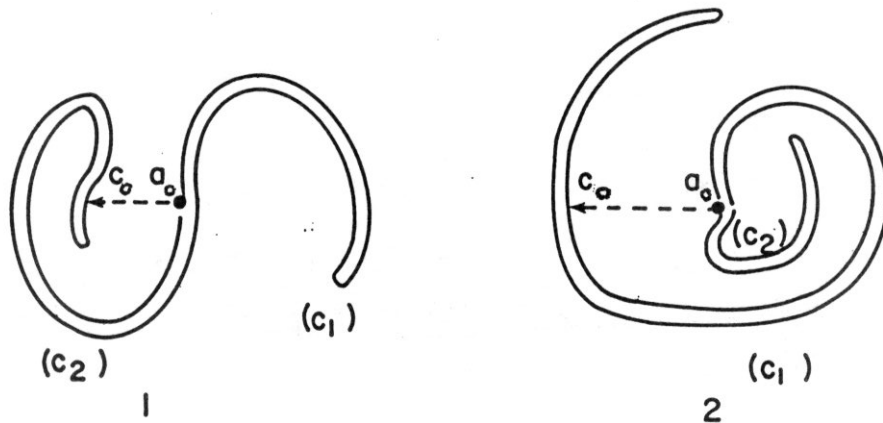


FIGURE 2.8

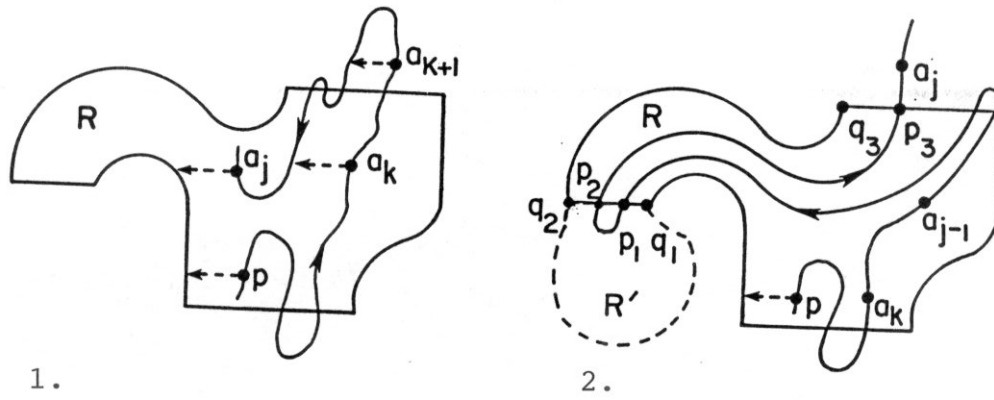


FIGURE 2.9

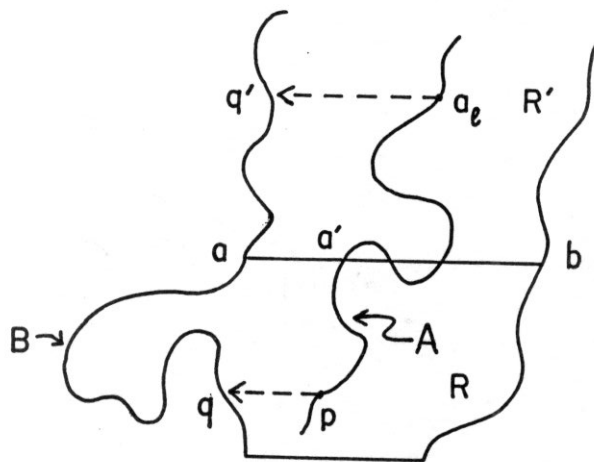


FIGURE 2.10

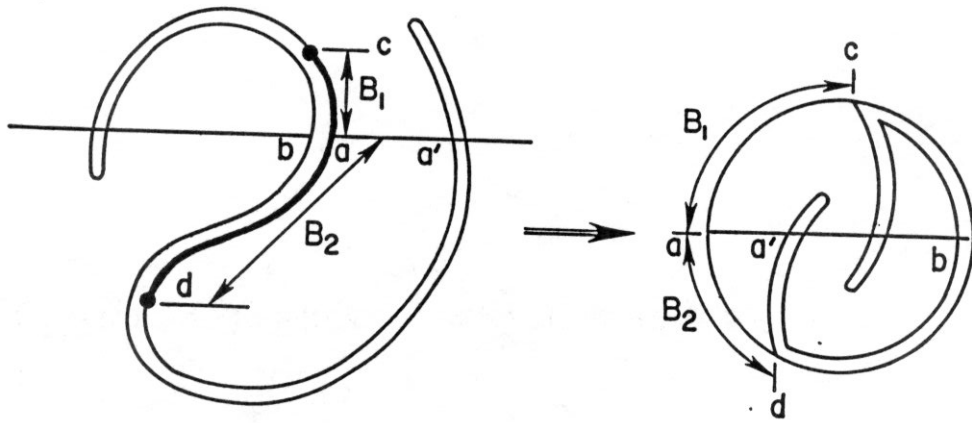


FIGURE 2.11

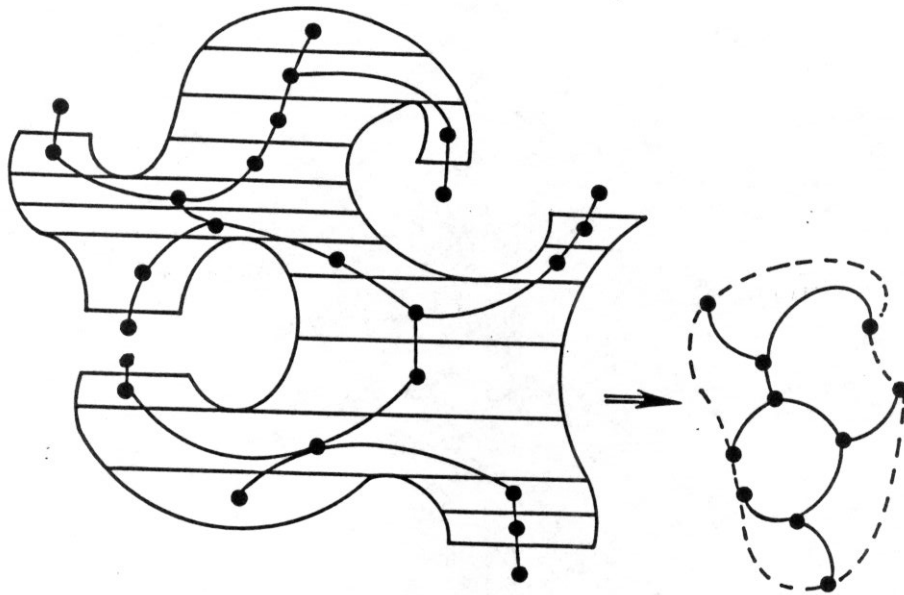


FIGURE 2.12