

2001

A Framework for Supporting the Class of Space Partitioning

Walid G. Aref

Purdue University, aref@cs.purdue.edu

Ihab F. Ilyas

Report Number:

01-002

Aref, Walid G. and Ilyas, Ihab F., "A Framework for Supporting the Class of Space Partitioning" (2001).
Department of Computer Science Technical Reports. Paper 1500.
<https://docs.lib.purdue.edu/cstech/1500>

This document has been made available through Purdue e-Pubs, a service of the Purdue University Libraries.
Please contact epubs@purdue.edu for additional information.

**A FRAMEWORK FOR SUPPORTING THE
CLASS OF SPACE PARTITIONING TREES**

**Walid G. Aref
Ihab F. Ilyas**

**Department of Computer Sciences
Purdue University
West Lafayette, IN 47907**

**CSD TR #01-002
June 2001**

A Framework for Supporting the Class of Space Partitioning Trees

Walid G. Aref and Ihab F. Ilyas
Department of Computer Sciences, Purdue University
West Lafayette IN 47907-1398
{aref,ilyas}@cs.purdue.edu

TR 01-002

Abstract

Emerging database applications require the use of new indexing structures beyond B-trees and R-trees. Examples are the k-D tree, the trie, the quadtree, and their variants. They are often proposed as supporting structures in data mining, GIS, and CAD/CAM applications. A common feature of all these indexes is that they recursively divide the space into partitions. A new extensible index structure, termed *SP-GiST* is presented that supports this class of data structures, mainly the class of space partitioning unbalanced trees. Simple method implementations are provided that demonstrate how *SP-GiST* can behave as a k-D tree, a trie, a quadtree, or any of their variants. Issues related to clustering tree nodes into pages as well as concurrency control for *SP-GiST* are addressed. A dynamic minimum-height clustering technique is applied to minimize disk accesses and to make using such trees in database systems possible and efficient. A prototype implementation of *SP-GiST* is presented as well as performance studies of the various *SP-GiST*'s tuning parameters.

Keywords: *SP-GiST*, space-partitioning trees, *GiST*, spatial tree indexes, access methods, clustering.

1 Introduction

Emerging database applications require the use of new indexing structures beyond B+-trees. The new applications may need different index structures to suit the big variety of data being supported, e.g., video, image, and multidimensional data. Typical applications are cartography, CAD, GIS, telemedicine, and multimedia applications. For example, the quadtree [18, 29] is used in the Sloan Digital Sky Survey to build indexes for different views of the sky (a multi-terabyte database archive) [45], the linear quadtree [21] is used in the recently released Oracle spatial product [10], the trie data structure is used in [1] to index handwritten databases, and the pyramid multi-resolution data structure [46] is used in the Microsoft TerraServer [2] which is an online atlas, currently being developed that combines around eight terabytes of image data. The reader is referred to [6, 10, 14, 16, 17, 20, 24, 37, 40, 42, 44] for additional database applications that use a variety of spatial and non-traditional tree structures.

Having a single framework to cover a wide range of these tree structures, although is very attractive from the point of view of database system implementation, is hindered by two main problems. The first problem is the *storage/structure characteristics* of spatial trees. Most of the *unbalanced* spatial tree structures are not optimized for I/O, which is a crucial issue for database systems. Quadtrees, tries, and k-D trees can be so *skinny and long*. Unless the problem of appropriately clustering the tree nodes into pages is addressed properly, this would lead to many I/O accesses before getting the required query answer. Compare this to the B+-tree, that in most cases has a height of 2-3 levels, and to the R-tree [25] and its variants, the R*-tree [4] and the R+-tree [43] that play an important role as spatial database indexes, e.g., see [7, 12, 38]. The second problem is the *implementation effort* of building indexes. Hard wiring the implementation of a full fledged index structure with the appropriate concurrency and recovery mechanisms into the database engine is a non-trivial process. Repeating this process for each spatial tree that can be more appealing for a certain application requires major changes in the DBMS core code. After all, one may still need a new structure that will cause, rewriting/augmenting significant portions of the DBMS engine to add the new tree index. The *Generalized Search Tree (GiST)* [26], was introduced in order to provide single implementation for B-tree-like indexes, e.g., the B+-tree [30], the R-tree [25], and the RD-tree [27]. Although practically useful, the class of unbalanced spatial indexes, e.g., the quadtree, the trie, and the k-D tree, is not supported by GiST because of the structure characteristics mentioned.

One important common feature of the quadtree, the trie, and the k-D tree family of indexes is that at each level of the tree, the underlying space gets partitioned into disjoint partitions. For example, in the case of a two-dimensional quadtree, at each level of decomposition, the space covered by a node is decomposed into four disjoint blocks. Similarly, in the case of the trie (assuming that we store a dictionary of words), the space covered by a node in the trie gets decomposed into 26 disjoint regions (each region corresponds to one letter of the alphabet). The k-D tree exhibits similar behavior. We use the term *space-partitioning* trees to represent the class of hierarchical data structures that decomposes a certain space into disjoint partitions. The number of partitions and the way the space is decomposed differ from one tree to the other.

In this paper we study the common features among the members of the spatial space partitioning trees aiming at developing a framework that is capable of representing the different tree structures and overcoming the difficulties that prevent such useful trees from being used in database engines. The DBMS will then be able to provide a large number of index structures with simple method plug-ins. As demonstrated in the paper, for the framework of space partitioning trees, we furnish in the DBMS (only once) the common functionalities such as the insertion, deletion, and updating algorithms, concurrency control and recovery techniques and I/O access optimization. For example, in a multimedia or a data mining application, we may then freely choose the best way to index each feature depending on the application semantics. By writing the right extensions to the extensible single implementation, a quadtree, a trie, a k-D tree, or other spatial structures can be made available without messing with the DBMS internal code.

The rest of the paper is organized as follows. Section 2 presents the class of space-partitioning trees. In Section 3, the SP-GiST framework is presented. Section 3 also includes a description of SP-GiST external user interface, and illustrates the realization of various tree structures using it. This includes a realization of the k-D tree, the trie, the Patricia trie, and several variants of the quadtree. Section 4 gives the implementation of the internal methods of SP-GiST. Concurrency control and recovery for SP-GiST are discussed in Section 5. Node clustering in SP-GiST is presented in Section 6. The pseudo code of the clustering algorithm is given

in Appendix 8. Implementation and experimental results for the various tuning parameters of SP-GiST are given in Section 7. Section 8 contains some concluding remarks.

2 The Class of Space Partitioning Trees

The term *space-partitioning* tree refers to the class of hierarchical data structures that recursively decomposes a certain space into disjoint partitions. It is important to point out the difference between data-driven and space-driven decompositions of space. If the principle of decomposing the space is dependent on the input data, it is called *data-driven* decomposition, while if it is dependent solely on the space, it is called *space-driven* decomposition. Examples of the first category are the k-D tree [5] and the point quadtree [29]. Examples of the second category are the trie index [11, 19], the fixed grid [35], the universal B-tree [3], the region quadtree [18], and other quadtree variants (e.g., the MX-CIF quadtree [28], the bintree, the *PM quadtree* [39], the PR quadtree [36] and the PMR quadtree [34]).

There are common underlying features among these spatial data structures. The term *quadtrie* was introduced in [40] to reflect the structure similarity between the trie and the quadtree. Similarly, the k-D tree and the MX quadtree have many structural similarities, e.g., both structures recursively partition the space into a number of disjoint partitions. On the other hand, the two trees differ in the number of partitions to divide the space and also in the decomposition principle. The decomposition is data-driven in the case of the k-D tree, while it is space-driven in the case of the MX quadtree.

The structural and behavioral similarities among many spatial trees create the class of space-partitioning trees. In contrast, the differences among these trees enable their use in a variety of emerging applications. The nature of spatial data that the application is dealing with, as well as the types of queries that need to be supported, aid in deciding which space-partitioning tree to use.

Space-partitioning trees can be differentiated on the following basis:

- **Structural differences**

- SD_1 : The type of data they represent.
- SD_2 : The decomposition fan-out (the number of partitions).
- SD_3 : The resolution of the underlying space.
- SD_4 : Allowing single-arc nodes.
- SD_5 : The use of buckets.

- **Behavioral differences**

- BD_1 : The decomposition principle (data or space driven partitioning).

The structural differences or design options can be viewed as *Shape Parameters* for the realized tree. For example, in the realization of the PR quadtree, or more precisely the *PR-quadtrie*, the represented data

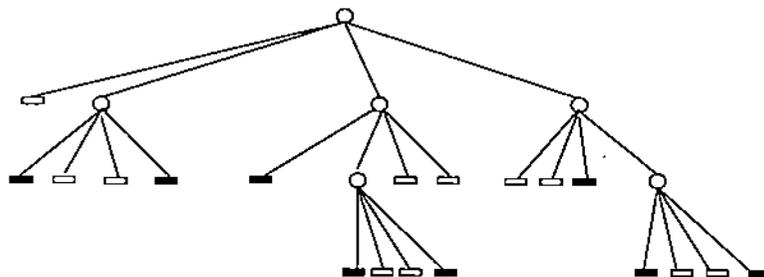


Figure 1: An example PR quadtree.

is “point” (SD_1). The decomposition depends on the space not on the data inserted (compare to the k-D tree) (BD_1). Each time a partitioning of the space quadrant into four equal quadrants (SD_2 and SD_4)

takes place to divide the quadrant that has two points so that each point is attached with one quadrant. The decomposition resolution is “variable” in the sense that the partitioning stops whenever one data point resides in the quadrant (SD_3). Figure 1 shows an example of the PR quadtree. At the leaf level, nodes can be “white” (i.e., contains no data) or “black” (i.e., contains one data point (SD_5)).

Using the same analogy, we can analyze the structure and behavior of the trie. The data represented in a trie is of type “word” (SD_1). The decomposition of the trie is space-dependent (BD_1), as we always decompose the space into 26 partitions (SD_2); one partition for each letter of the alphabet. In one variant of the trie, the resolution is “not variable” (SD_3) as we need to decompose the space until we consume all the letters of the inserted word (refer to Figure 2a for illustration). This is in contrast to stopping the decomposition only when a space partition uniquely identifies the inserted word (see Figure 2b). The same

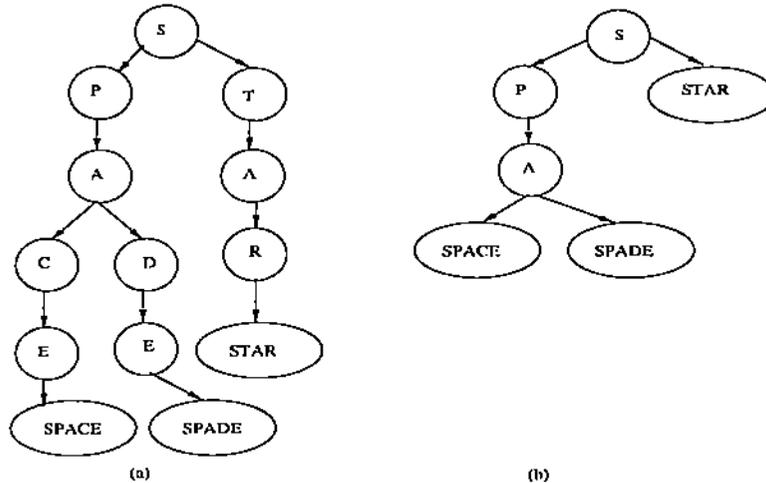


Figure 2: Two variants of the trie data structure : (a) Resolution is not variable (b) resolution is variable.

analysis can be applied to realize other quadtree and trie variants, the k-D tree, and the bin-tree.

In the following sections, we will introduce a general framework, termed *SP-GiST*, which we can use to implement a big collection of space-partitioning trees. *SP-GiST* has one core implementation as well as user plug-ins that reflect the required structural and behavioral characteristics. The existence of such a framework will facilitate the adaption of this class of space-partitioning trees into database engines.

3 SP-GiST Framework Interface

SP-GiST is a general index framework that covers a wide range of tree indexes representing the large class of *space-partitioning search trees* represented in Section 2.

The *structural* characteristics of space-partitioning trees that distinguish them from other tree classes are: (1) Space-partitioning trees decompose the space recursively. Each time, a fixed number of *disjoint* partitions is produced. (2) Space-partitioning trees are unbalanced trees (3) Space-partitioning trees suffer from limited fan-out, e.g., the quadtree has only a fan-out of four. So, space-partitioning trees can be *skinny* and *long*. (4) Two different types of nodes exist in a space-partitioning tree, namely, index nodes (internal nodes) and data nodes (leaf nodes). The framework reflects these facts by having two main parts; the *internal* tree methods that reflect the *similarities* among all members of the class of space-partitioning trees, and the *external* interface that enables us to identify the features specific to a particular tree reflecting the *differences* listed in Section 2.

By specifying user access methods as in GiST [26], SP-GiST has interface parameters and methods that allow it to represent the class of space-partitioning tree indexes and reflect the structural and behavioral differences among them.

3.1 Interface Parameters

The user can realize a particular space-partitioning tree using the following interface parameters:

- *NodePredicate*: This parameter gives the predicate to be used in the index nodes of the tree (addresses the structural difference SD_1). For example, a quadrant in a quadtree or a letter in a trie are predicates that are associated with an index node.
- *Key Type*: This parameter gives the type of the data in the leaf level of the tree. For example, "Point" is the key type in an MX quadtree while "Word" is the key type in a trie. The data type Point and the data type Word have to be pre-defined by the user.
- *NumberOfSpacePartitions*: This parameter gives the number of disjoint partitions produced at each decomposition (SD_2). It also represents the number of items in index nodes. For example, quadtrees will have four space partitions, a trie of the English alphabet will have 26 space partitions, the k-D tree will have only two space partitions at each decomposition.
- *Resolution*: This parameter gives the maximum number of space decompositions and is set depending on the space and the granularity required.
- *PathShrink*: For space-partitioning trees, recursive decomposition can lead to long sparse structures. Parameter *PathShrink* is useful in limiting the number of times the space is recursively decomposed in response to data insertion. *PathShrink* can be one of three different policies (refer to Figure 3 for an illustration of the use of *PathShrink* in the context of the trie):

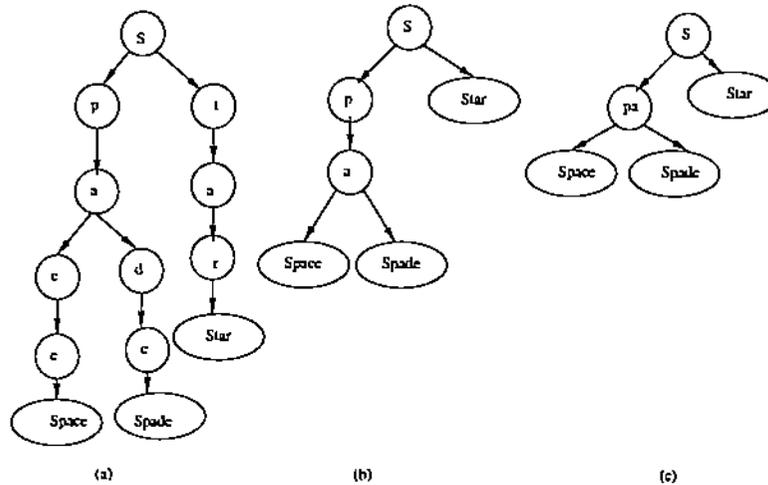


Figure 3: The effect of the parameter *PathShrink* on the trie : (a) Never Shrink, (b) Leaf Shrink, and (c) Tree Shrink.

- *Never Shrink*: Data is inserted in the node that corresponds to the maximum resolution of the space. This may result in multiple recursive decompositions of the space.
- *Leaf Shrink*: Data is inserted at the first available leaf node. Decomposition will not depend on the maximum possible resolution. In this strategy, no index node will have one leaf node as we decompose only when there is no room for the newly inserted data item.
- *Tree Shrink*: The internal nodes are merged together to eliminate all single child internal nodes. This strategy is adapted from structures like the *Patricia* trie that aim at reducing the height of the tree as much as possible.

For example, in the case of *PathShrink* = "Never Shrink", when storing the word "implementation" in the trie, the word will be stored in a leaf after a 14-nodes path, one level per input character. On the

words, at level i , splitting will be according to the i^{th} character of each word in the over-full node. `PickSplit` will return the entries of the split nodes in the output parameter `splitnodes`, which is an array of buckets, where each bucket contains the elements that should be inserted in the corresponding child node. The predicates of the children are also returned in `splitpredicates`.

- `Cluster()`: This method defines how tree nodes are clustered into disk pages. The method is explained in more detail in Section 6.

The interface methods realize the *behavioral* design options listed in Section 2. Methods `Consistent` and `PickSplit` determine if the tree follows the space-driven or the data-driven partitioning. For example, in a k-D tree, which is a data-driven space partitioning tree, method `Consistent` compares the coordinates of the query point (the point to be inserted or searched for) against the coordinates of the point attached to the index node. The values of these coordinates are determined based on data that is inserted earlier into the k-D tree. On the other hand, method `Consistent` for a space-driven space partitioning tree, e.g., the trie, will only depend on the letters of the newly inserted word. The comparison is performed against the letter associated with the index node entry, which is space-dependent, and is independent of the previously inserted data.

We can also show that method `PickSplit` completes the specification of the behavioral design option by specifying the way to distribute node entries among the produced partitions. Examples of `PickSplit` for various tree structures are given in the following section.

3.3 Realization of Space-Partitioning Trees

Using the SP-GiST interface, given in the previous sections, we demonstrate how to realize some commonly used space-partitioning indexes. More specifically, we present the realization of the k-D tree, variants of the quadtree, the trie, and the Patricia trie.

The k-D tree: k-D trees [5] are a special kind of search trees, useful for answering range queries about a set of points in the k-dimensional space. The k-D tree uses a data-driven decomposition of the space (see Section 2). The tree is constructed by recursively *partitioning* the space into two sub-spaces with respect to one of the dimensions at each tree level.

The k-D tree insertion algorithm for the two-dimensional case (i.e., $k = 2$) with points in the xy plane is as follows: The algorithm selects any point and draws a line through it, parallel to the y -axis. This line partitions the plane vertically into two sub-planes. Another point is selected and is used to horizontally partition the sub-plane in which it lies. In general, a point that falls in a region created by a horizontal partition will divide this region vertically, and vice versa. This division process induces a binary tree structure, (e.g., see Figure 5).

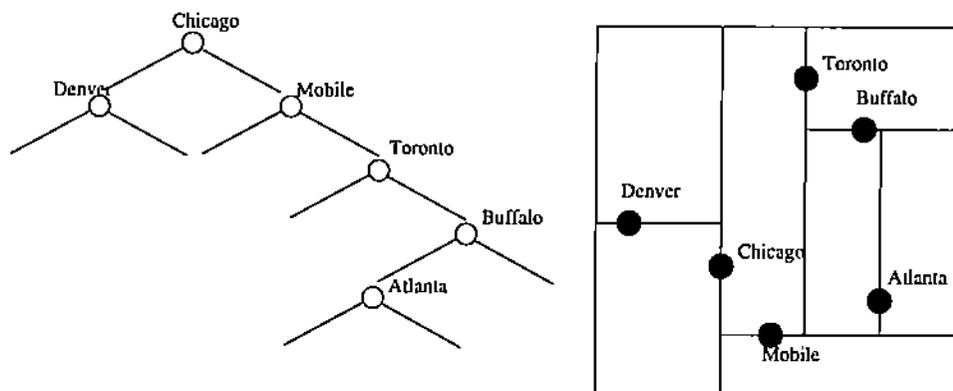


Figure 5: An example k-D Tree.

The realization of the k-D tree is given in Table 1. `PathShrink` is set to "Leaf Shrink" because we put each input point at the first available place depending on the previously inserted points. Each node will hold

Parameters	PathShrink = Leaf Shrink NodeShrink = False BucketSize = 1 NumberOfSpacePartitions = 2 Node Predicate = "left", "right", or blank. Key Type = Point
Consistent(E, q, level)	IF (level is odd AND q.x satisfies E.p.x) OR (level is even AND q.y satisfies E.p.y) RETURN TRUE ELSE RETURN FALSE
PickSplit(P, level)	Put the old point in a child node with predicate "blank" Put the new point in a child node with predicate "left" or "right". RETURN FALSE

Table 1: Realization of the k-D Tree using SP-GiST.

only one point, ($BucketSize = 1$). $NodeShrink$ is set to *false*, so each index node will have a slot for the left subtree and a slot for the right subtree. We have only two space partitions for the "right" and "left" to a point ($NumberOfSpacePartitions = 2$).

The Quadtree: The term *quadtree* describes a class of hierarchical data structures whose common property is the recursive decomposition of space into quadrants. The quadtree can be realized by SP-GiST. In the next subsections, examples of various types of quadtrees are presented for point data, rectangles, and polygonal data. Note that for all the variants, the number of space partitions is equal to four ($NumberOfSpacePartitions = 4$), with a bucket size of B items ($BucketSize = B$). $NodeShrink$ is set to *false*, so each index node will have a slot for each partition even if it is empty. Setting $NodeShrink$ to *true* would realize a quadtree with all *white* nodes eliminated (see Figure 6) [40].

The quadtree can be viewed as a trie structure in two dimensions - with only two possible characters in each dimension, in *trie* terminology, or even a one dimensional trie with only a four character alphabet set. Thus in the literature, *space-driven* quadtrees are often called *quadtries* [41].

When we treat data points as nonzero elements in a square matrix, the resulting data structure is called the *MX quadtree* (MX for matrix). In the MX quadtree, leaf nodes are black or empty (white) corresponding to the presence or absence, respectively, of data points in the appropriate position in the matrix. Each point in an MX quadtree corresponds to a 1×1 square. Figure 6 gives an example of an MX quadtree. Notice that data nodes of the MX quadtree all appear at the same level. The number of space decompositions is predefined depending on the desired space resolution.

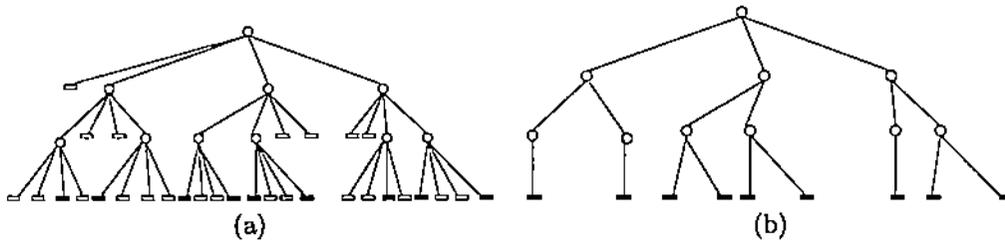


Figure 6: An example MX quadtree: (a) $NodeShrink = false$, and (b) $NodeShrink = true$.

For the MX quadtree, realized in Table 2, $PathShrink$ is set to "Never Shrink". Therefore, the tree is expanded through successive splitting to the maximum space resolution. Thus, $PickSplit$ will not be invoked as each point will fall in one node.

The MX quadtree is applicable as long as the domain of data points is discrete and finite. If this is not the case, the data points cannot be represented using an MX quadtree since the minimum separation between the data points will be unknown. This leads to the idea of associating data points with quadrants

Parameters	PathShrink = Never Shrink NodeShrink = False BucketSize = B NumberOfSpacePartitions = 4 Node Predicate = Quadrant represented by (x1,y1,x2,y2) where (x1,y1) are the values of the coordinates of the top left corner and (x2,y2) are the values of the coordinates of the bottom right corner Key Type = Point
Consistent(E, q, level)	IF (q coordinates inside E.quadrant) RETURN TRUE ELSE RETURN FALSE
PickSplit(P, level)	RETURN FALSE

Table 2: Realization of the MX quadtree using SP-GiST.

Parameters	PathShrink = Leaf Shrink NodeShrink = False BucketSize = B NumberOfSpacePartitions = 4 Node Predicate = Quadrant represented by (x1,y1,x2,y2) where (x1,y1) are the values of the coordinates of the top left corner and (x2,y2) are the values of the coordinates of the bottom right corner Key Type = Point
Consistent(E, q, level)	IF (q coordinates inside E.quadrant) RETURN TRUE ELSE RETURN FALSE
PickSplit(P, level)	Partition and allocate data points into quadrants according to the locations of the data points IF any partition is still over-full RETURN TRUE ELSE RETURN FALSE

Table 3: Realization of the PR quadtree using SP-GiST.

and hence realizing the *PR quadtree* (P for point and R for region) [36]. Now, each data point maps to a quadrant and not to a 1×1 square as in the MX quadtree. Figure 1 gives an example of the PR quadtree.

The PR quadtree can be realized using SP-GiST by setting *PathShrink* to "Leaf Shrink" as we put each input point at the first available leaf node. The leaf node is not necessarily of size 1×1 . Realization of the PR quadtree using SP-GiST is given in Table 3.

The *MX-CIF quadtree* is a quadtree variation for storing rectangles. It associates each rectangle, say *R*, with the quadtree node corresponding to the smallest block that contains *R* in its entry. Rectangles can be associated with both leaf and non-leaf nodes. The subdivision ceases whenever a node's block contains no rectangles. Figure 7 gives an example MX-CIF quadtree. Notice that more than one rectangle can be associated with a given node.

The MX-CIF quadtree can be realized by SP-GiST, as given in Table 4. *PickSplit* is not applicable here, because according to the MX-CIF insertion algorithm, there is not much choice as to where a rectangle gets inserted.

Another quadtree variant is the *PMR quadtree* [34] that is used to store polygonal maps. The key is of type line segment in the PMR quadtree, where line segments serve as the building block to construct polygons.

The PMR quadtree is an edge-based data structure. A line segment is stored in a PMR quadtree by inserting the line segment into the nodes corresponding to all the blocks that it intersects. If the bucket capacity is exceeded, the node's block is split *once, and only once*, into four equal quadrants. Thus, bucket

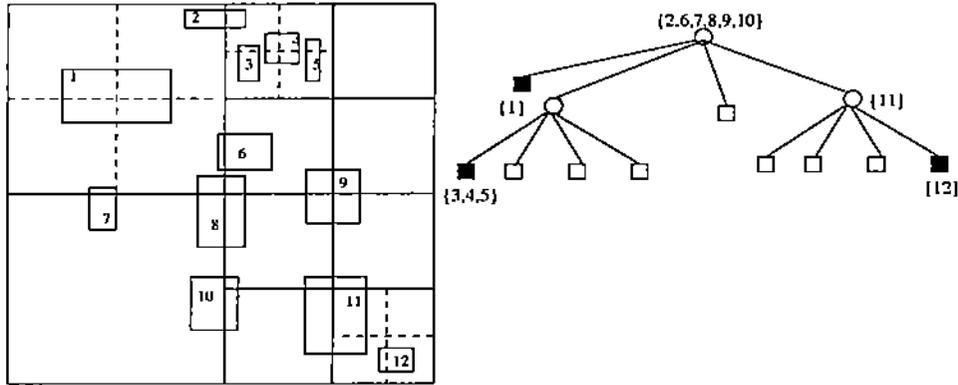


Figure 7: An example MX-CIF quadtree.

Parameters	PathShrink = Leaf Shrink NodeShrink = False BucketSize = B NumberOfSpacePartitions = 4 Node Predicate = Quadrant represented by $(x1,y1,x2,y2)$ where $(x1,y1)$ are the values of the coordinates of the top left corner and $(x2,y2)$ are the values of the coordinates of the bottom right corner Key Type = Rectangle
Consistent (E, q, level)	IF (Node predicate is the minimum bounding quadrant of q AND the E.p is Blank) RETURN TRUE IF (E.p contains q) RETURN TRUE ELSE RETURN FALSE
PickSplit(P, level)	RETURN FALSE

Table 4: Realization of the MX-CIF quadtree using SP-GiST.

capacity is really a splitting threshold. The PMR quadtree can be realized using SP-GiST, as given in Table 5.

The Trie: A *trie* [11, 19] is a tree in which the branching at any level is determined by only a portion of the key as in Figure 4 (a). The trie contains two types of nodes; index and data nodes. In the trie of Figure 4 (a), each index node contains 27 link fields. In the Figure, index nodes are represented by rectangles, while data nodes are represented by ovals.

All characters in the key values are assumed to be one of the 26 letters of the alphabet. A blank is used to terminate a key value. At level 1, all key values are partitioned into 27 disjoint classes depending on their first character. Thus, $LINK(T,i)$ points to a subtree containing all key values beginning with the i^{th} character (T is the root of the trie). On the j^{th} level the branching is determined by the j^{th} character. When a subtree contains only one key value, it is replaced by a node of type data. This node contains the key value, together with other relevant information such as the address of the record with this key value, etc. The trie can be realized using SP-GiST, as given in Table 6. Notice that *PathShrink* is set to "Leaf Shrink" (refer to Section 3.1). *NodeShrink* is set to *false* in this realization of the trie. Another option is to set the *NodeShrink* to *true* to realize the forest trie [30], as discussed in Section 3.1.

The regular trie suffers from the problem of long skinny paths of single arc nodes. For example, for a trie with a bucket size of 2, inserting the three words "abate", "abacus", and "abort" will cause the node to split. Since we are at the first level, the split will depend on the first character in each word. Since all the words have "ab" as their first and second characters, splitting must continue until the third character, resulting in a skinny trie (see Figure 8).

Parameters	PathShrink = Leaf Shrink NodeShrink = False BucketSize = B NumberOfSpacePartitions = 4 Node Predicate = Quadrant represented by (x1,y1,x2,y2) where (x1,y1) are the values of the coordinates of the top left corner and (x2,y2) are the values of the coordinates of the bottom right corner Key Type = Line Segment represented by end points
Consistent (E, q, level)	IF (inserted line intersects E.quadrant) RETURN TRUE ELSE RETURN FALSE
PickSplit(P, level)	Partition the line segments according to their intersections with quadrants RETURN FALSE

Table 5: Realization of the PMR quadtree using SP-GiST.

Parameters	PathShrink = Leaf Shrink NodeShrink = False BucketSize = B NumberOfSpacePartitions = 26 Node Predicate = letter or Blank Key Type = String
Consistent(E, q, level)	IF (q[level] == E.letter) OR (E.letter == BLANK AND level > length(q)) RETURN TRUE ELSE RETURN FALSE
PickSplit(P, level)	Partition the data strings in P according to the character values at position "level" IF any data string has length < level, insert data string in Partition "blank" IF any of the partitions is still over-full RETURN TRUE ELSE RETURN FALSE

Table 6: Realization of the Trie using SP-GiST.

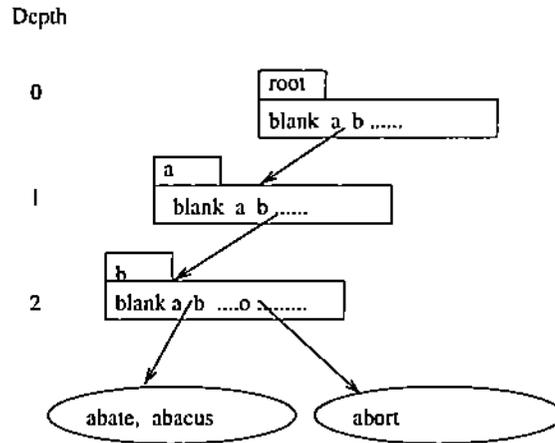


Figure 8: An example trie with *BucketSize* = 2: the final tree after inserting the three words.

The *Patricia trie* [33, 30] is a special trie structure that addresses this problem. It has the property that all nodes that have only one arc are merged with their parent nodes. To avoid false matches, each node in the Patricia trie must have either a count of the number of eliminated nodes or a pointer to the eliminated symbols. In the previous example, (refer to Figure 9), the Patricia trie will split only once, thus eliminating the single arc nodes and storing the eliminated symbols ("ab") in the parent node.

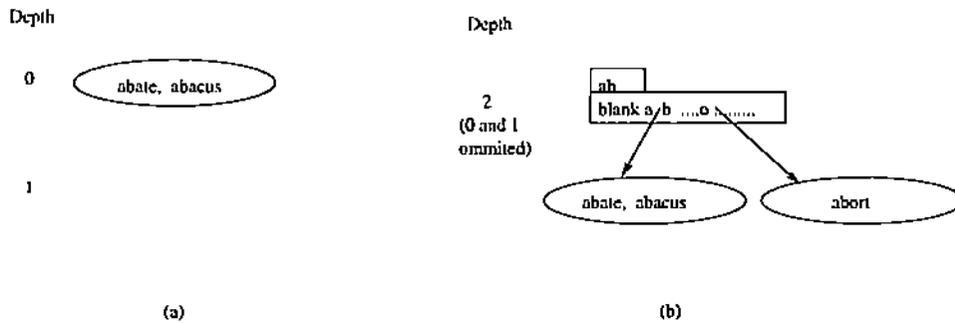


Figure 9: An example Patricia trie with *BucketSize* = 2: (a) after inserting the first two words, and (b) after inserting the third word.

In the Patricia trie, *PathShrink* is set to "Tree Shrink". When splitting a node, we search for the common prefix of all words. The common prefix is returned as the predicate of the parent node, while splitting is performed based on the next letter after that prefix. The realization of the Patricia trie using SP-GiST is given in Table 7.

4 SP-GiST Internal Methods

The methods for insertion, deletion, and search in SP-GiST are internal operations that are implemented inside the SP-GiST index engine. These methods are used in conjunction with the external methods to realize specific space-partitioning trees. The user of SP-GiST provides only the external methods, while the internal methods are hard coded into the SP-GiST index engine. The internal methods are general for the class of space-partitioning trees, and their behavior is tuned by making use of the user-defined external methods and parameters.

The internal methods are designed to accommodate for the space-partitioning, recursive decomposition, bucket sizes, insertion resolution, and node clustering (refer to the structural and behavioral characteristics of space-partitioning trees, given in Section 2).

Parameters	PathShrink = Tree Shrink NodeShrink = False BucketSize = B NumberOfSpacePartitions = 26 Node Predicate = letter or Blank Key Type = String
Consistent(E, q, level)	IF (q[level] == E.letter) OR (E.letter == BLANK AND level > length(q)) RETURN TRUE ELSE RETURN FALSE
PickSplit(P, level)	Find a Common prefix among words in P Update level = level+length of the common prefix Let P predicate = the common prefix Partition the data strings in P according to the character values at position "level" IF any data string has length < level, insert data string in Partition "blank" IF any of the partitions is still over full RETURN TRUE ELSE RETURN FALSE

Table 7: Realization of the Patricia Trie using SP-GiST.

Recall that unlike the GiST structure, SP-GiST has to support two distinct types of nodes; index and data nodes. **Index nodes (non-leaf nodes)** hold the various space partitions at each level. Each entry in an index node is a root of a subtree that holds all the entries that lie in this partition. The space partitions are disjoint. Besides having a slot for each space partition, the index node contains an extra *blank* slot to point to data nodes attached to the partition represented by this node. On the other hand, **data nodes (leaf nodes)** hold the key data and other pointer information to physical data records. We can think of data nodes as *Buckets* of data entries. Thus, a splitting strategy determined by *PickSplit* will be applied to split over-full data nodes.

The insert algorithm, given in Table 8, depends on the following interface parameters and external methods:

1. Parameter *PathShrink* specifies how deep we should proceed with the space decomposition.
2. Method *Consistent* specifies which branch to follow.
3. Method *PickSplit* to split over-full nodes. The return value of *PickSplit* tells us when we should stop the splitting process.

Method *Insert* begins by checking Parameter *PathShrink*. If *PathShrink* is set to "Never Shrink", method *Insert* performs a successive creation of index nodes to the maximum space resolution. If the parameter is set to "Leaf Shrink" or "Tree Shrink", the insertion algorithm searches for the first leaf node with a predicate that is *Consistent* with the key to be inserted. In the case of "Tree Shrink", some eliminated index nodes may be needed while locating the leaf. Hence, an internal split is performed to "expand" the eliminated index nodes. If the leaf node is over-full, then method *PickSplit* will be invoked continuously to distribute the entries among non over-full children or until it reaches the maximum resolution of the underlying space. Notice that method *Insert* invokes method *cluster* to dynamically re-cluster the nodes properly after insertion. Node clustering is further explained in Section 6.

Method *Search* in SP-GiST is exactly similar to that of the GiST scheme (see [26]), and is given in Table 9 for completeness. Method *Search* uses method *Consistent* as the main navigation guide. Starting from the root, the algorithm will check the search item against all branches using the method *Consistent* till reaching leaf nodes (data nodes in SP-GiST).

The algorithm for method *Delete* in SP-GiST uses logical deletion. Deleted items are marked deleted and are not physically removed from the tree. This will save the effort of reorganizing the tree after each deletion, specially for data-driven space-partitioning trees. A rebuild is used from time to time as a *clean* procedure.

```

ALGORITHM INSERT (TreeNode root, Key, level)
CurrentNode =root /* Initially root is null */
IF PathShrink is "Never Shrink" THEN
  LOOP WHILE level < SpaceResolution AND level < Key length
  IF node is NULL THEN E = Create a new node of type INDEX
  FOR each slot i in the index node LOOP
    IF (Consistent(E[i],key,level)) THEN index=i
  IF None is consistent /*due to NodeShrink*/
  THEN Create the missing index slot w.r.t level
    index = the position of the new slot
    CurrentNode = E[index].ptr /* the child pointed by entry E[index]*/
    level = level +1
IF CurrentNode is INDEX node /* pick a child to go */
  Compare the key with the CurrentNode predicate
  IF no match AND PathShrink is "Tree Shrink"
  THEN get the common prefix between the two
    Change CurrentNode predicate to the common prefix
    Create a new INDEX node with the rest of the old node predicate
    Let CurrentNode be the new index node
  FOR each slot i in the index node LOOP
    IF (Consistent(E[i], key, level)) THEN index=i
    IF None is consistent /*due to NodeShrink*/
    THEN Create the missing index slot w.r.t level
      index = the position of the new slot
      CurrentNode=CurrentNode[index].ptr
  INSERT (CurrentNode, key, level+1) /* recursive */
IF CurrentNode is full THEN /* DATA node and may need to be split*/
  LOOP WHILE PickSplit(node,level)
  n=Create new node of type INDEX
  Create Children for the split entries
  Parent(n) = Parent(CurrentNode)
  Adjust branches of 'n' to point to the new children
  level = level +1
ELSE insert the key in CurrentNode /* not a full node */
Cluster() /* to recluster the tree nodes in pages */

```

Table 8: SP-GiST Insertion Algorithm.

```

SEARCH (TreeNode root, Key, level)
Found = false
CurrentNode =root /* Initially root is null */
LOOP WHILE level < SpaceResolution AND CurrentNode is an index node
  Compare the key with the CurrentNode predicate
  IF no match AND PathShrink is "Tree Shrink"
  THEN Found = FALSE
    break
  FOR each slot i in the index node LOOP
    IF (Consistent(E[i], key, level)) THEN index=i
    IF None is consistent /*due to NodeShrink*/
    THEN Found = FALSE
      break
    CurrentNode = E[index].ptr /* the child pointed by entry E[index]*/
    level = level +1
IF CurrentNode is NOT NULL /* leaf node */
  Search for the key among leaf node entries
  IF Key is in the leaf node THEN Found = TRUE
RETURN Found

```

Table 9: SP-GiST Search Algorithm.

5 Concurrency and Recovery in SP-GiST

Concurrency and recovery in GiST have been addressed in [9, 32]. In [32], the authors provide general algorithms for concurrency control in tree-based access methods as well as a recovery protocol and a mechanism for ensuring repeatable read isolation [22]. They suggest the use of *Node Sequence Number* (NSN) for concurrency control, first introduced in [31].

For SP-GiST, a split (only at the leaf level) transforms a data node into an index node. Data is then distributed among new leaf nodes *rooted at that split node*. This fact simplifies the concurrency control mechanism significantly. For example, consider the case when a search for a key is interleaved with an insertion that causes the splitting of the target node. By the time the search reaches the target node, it can not falsely conclude the non-existence of the searched key, e.g., in contrast to a B-Tree scenario, because the new node is an index node. In that case, no right links need to be maintained between leaves as the search will need to continue *deeper in the tree* not on the siblings level. Thus, no special sequence number is needed for the concurrent operation to know that the node in question has been split. The operation will directly continue working with the child nodes.

Phantom protection in GiST has also been addressed in two different techniques. Predicate locking [15] is used in [32] while the authors in [9], propose a dynamic granular locking approach (GL/GiST) to phantom protection. We adopt the granular locking technique since it is more preferable and less expensive than predicate locking. The fact that a "Containment Hierarchy" exists in space-partitioning trees, represented by SP-GiST, makes the algorithm introduced in [8, 9] highly applicable and much simpler. Hence, in SP-GiST, because the node predicates form a containment hierarchy, we simply use the node predicates for granular locks.

The main difference in SP-GiST is that a page may contain multiple SP-GiST nodes. A clustering algorithm will hold the mapping between nodes and pages. In this context, we assume that the node size is smaller than or equal to the page size. Hence the problem transforms to locking at a finer granularity. Treating nodes clustered in pages as records, granular locks [23] are used. The recovery technique used in [32] is directly applicable to SP-GiST.

6 Node Clustering in SP-GiST

Node clustering means choosing the group of nodes that will reside together in the same disk page. Considering physical storage of the tree nodes, a direct and simple implementation of a node is to assign a disk page for each node. However, for very sparse nodes, this simple assignment will not be efficient for database use. We provide to the user a default node clustering method that is shown to perform well in the dynamic case [13]. However, we allow the user to override the default clustering method and provide a different node clustering policy that is more suitable for the type and nature of the operations to be performed on the constructed index. This will enhance the query response time of SP-GiST. We propose the interface method *Cluster* for this purpose.

Introducing new nodes in the tree structure, e.g., due to insertions, will internally invoke the dynamic clustering algorithm defined in *Cluster* to reconstruct the tree disk page structure and reflect the change. However, for unexperienced users or for typical database applications, SP-GiST has a default node clustering algorithm that achieves minimum height and hence minimum I/O access. The dynamic clustering algorithm in [13] is a good clustering algorithm and we use it as our default in SP-GiST. The pseudo code and a brief outline of the clustering algorithm is given in Appendix 8.

The user can choose other clustering algorithms that reflect the application semantics, specially for non-traditional data types as in multimedia or video databases. Some alternatives are: (1) **Fill-Factor Clustering:** Tries to keep each page half-full for space utilization efficiency. (2) **Deep Clustering:** Chooses the longest linked subtree from the collection of page nodes to be stored together in the same page. This clustering method will enhance performance for depth-first traversal of trees. (3) **Breadth Clustering:** Chooses the maximum number of siblings of the same parent to be stored together in the same page.

7 Implementation and Experimental Results

We implemented SP-GiST using C++ on SunOS 5.6 (Sparc). As a proof of concept, using SP-GiST, we implemented the extensions for some data structures namely, the MX quadtree, the PR quadtree, the trie, and the Patricia trie. The implementation has proven the feasibility of representing space-partitioning trees using the interface proposed by SP-GiST and the settings in the tables in Section 3.3. We performed experiments on various settings of the tunable interface parameters; *BucketSize* and *PathShrink*. In our implementation we adopt the *minimal height* clustering technique in [13]. Results show that applying this clustering technique reduces the path length in terms of pages significantly.

As explained in Section 3.1, the interface parameter *PathShrink* can take one of three values; "Never Shrink", "Leaf Shrink", or "Tree Shrink". For the trie, setting *PathShrink* to "Never Shrink" will have the effect of realizing the original trie, where splitting is performed to the maximum resolution of the space, leading to a long sparse tree. Setting *PathShrink* to "Leaf Shrink" will realize a common variant of the trie where data can be put in the first available node. On the other hand, if *PathShrink* is set to "Tree Shrink", it will realize the Patricia implementation of the trie where no single-arc nodes are allowed.

Figure 10 gives the effect of this parameter on the trie data structure for various settings of *BucketSize* for a dataset of 10000 records with "string" keys. As expected, for the trie and the Patricia trie, the path length and the number of pages improve as the bucket size increases since less splitting takes place. On the other hand, the bucket size does not have an effect on the original trie. In this case, splitting will take place not because of the bucket size limit but to decompose the space to the maximum resolution. In the case of the original trie, each record will fall in a single node regardless of the setting of the bucket size.

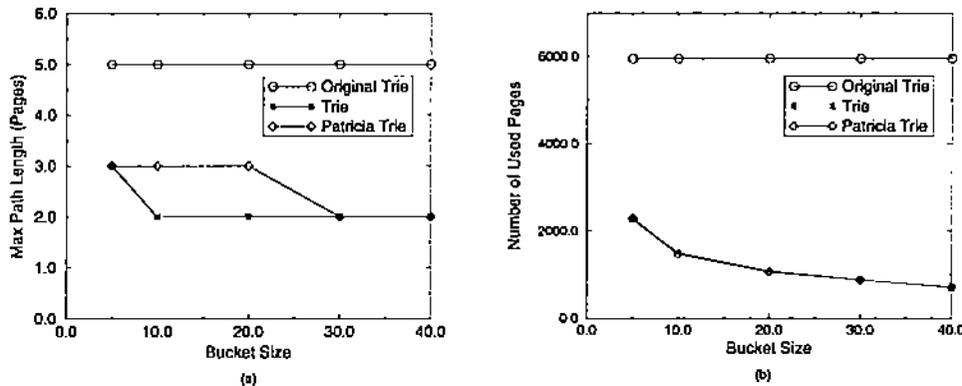


Figure 10: Effect of *BucketSize* on the maximum path length and number of used pages in the trie for different settings of *PathShrink*.

For the quadtree, the same argument holds. Experimental results for point datasets of 10000 points are given in Figure 11. In this case, setting *PathShrink* to "Never Shrink" will have the effect of realizing the MX quadtree while setting it to "Leaf Shrink" will realize the PR quadtree where data can be put in the first available node. Experiments with setting *PathShrink* to "Tree Shrink" show the realization of another variant of quadtree, where all *white nodes* are eliminated [41], making it more attractive for databases and solving the problem of long degenerate quadtrees when the workload is highly skewed.

8 Conclusions

SP-GiST is a generalized space-partitioning tree implementation of a wide range of tree data structures that are not I/O-optimized for databases. This makes it possible to have single tree index implementation to cover various types of trees that suit different applications. Emerging database applications will require the availability of various index structures due to the heterogeneous collection of data types they deal with. SP-GiST is an interesting choice for multimedia databases, spatial databases, GIS, and other modern database

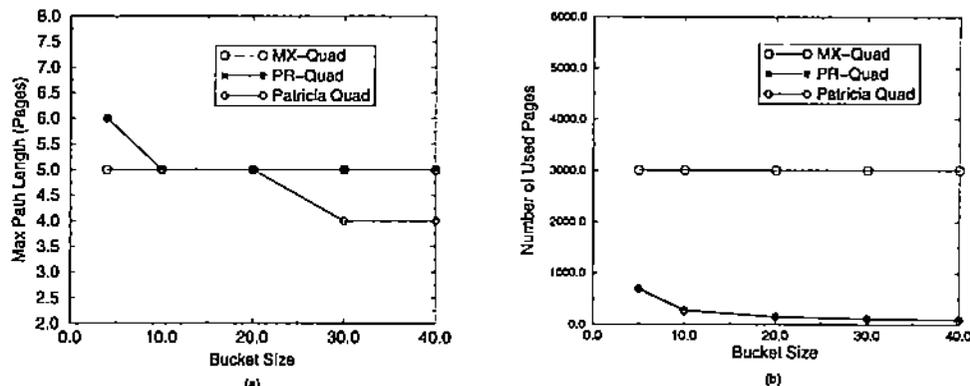


Figure 11: Effect of *BucketSize* on the maximum path length and number of used pages in the quadtree for different settings of *PathShrink*.

systems. We have shown how to augment SP-GiST with parameters and methods that will enable the coverage of this class of space-partitioning trees. Clustering methods were also addressed to realize the use of these structures in practice in non-traditional database applications.

Recovery and concurrency for SP-GiST are addressed to enable the realization of SP-GiST in commercial database systems. Experiments proved the concept of SP-GiST and provided insight on the effect of the tunable interface parameters on the tree structure and performance.

APPENDIX A - Minimal Height Clustering Algorithm

The clustering algorithm in [13] re-clusters the tree nodes into disk pages after updates to an already clustered state, and hence is dynamic, and guarantees *minimal height* mapping after deleting or inserting in the tree. The pseudo code of the algorithm is given in Figure 12. The algorithm begins by removing all deleted nodes from the disk pages. Now, all new nodes or affected roots of subtrees are kept in a set of affected nodes. Processing all the affected nodes starts bottom-up (no node is processed until all its children are processed). The algorithm tries to put the longest path of nodes together in the same page. The authors in [13] have shown that the algorithm achieves minimal height mapping. They suggested some heuristics of merging sparse pages to achieve minimum fill factor of at least 50%.

References

- [1] Walid G. Aref, Daniel Barbará, and Padmavathi Vallabhaneni. The handwritten trie: Indexing electronic ink. In *Proceedings of SIGMOD'95*, San Jose, California, May 1995.
- [2] Tom Barclay, Jim Gray, and Don Slutz. Microsoft terraserver: A spatial data warehouse. In *Proceedings of SIGMOD'00*, Dallas TX, May 2000.
- [3] R. Bayer. The universal B-tree for multidimensional indexing: General concepts. *Lecture Notes in Computer Science*, 1274, 1997.
- [4] N. Beckmann, H. P. Kriegel, R. Schneider, and B. Seeger. The R* -tree: an efficient robust access method for points and rectangles. *SIGMOD Record*, 19(2), 1990.
- [5] Jon L. Bentley. Multidimensional binary search trees used for associative searching. *Communications of the ACM*, 19:509–517, 1975.
- [6] S. Berchtold, C. Boehm, and H.-P. Kriegel. Improving the query performance of high-dimensional index structures by bulk load operations. *Lecture Notes in Computer Science*, 1377, 1998.

- [7] Thomas Brinkhoff, Hans-Peter Kriegel, and Bernhard Seeger. Parallel processing of spatial joins using R-trees. In *Proceedings of ICDE'96*, New Orleans, Louisiana, February 1996.
- [8] Kaushik Chakrabarti and Sharad Mehrotra. Dynamic granular locking approach to phantom protection in R-trees. In *Proceedings of ICDE'98*, pages 446–454, Orlando, Florida, USA, February 1998.
- [9] Kaushik Chakrabarti and Sharad Mehrotra. Efficient concurrency control in multidimensional access methods. In *Proceedings of SIGMOD'99*, pages 25–36, Philadelphia, Pennsylvania, USA, June 1999.
- [10] Oracle Corporation. Oracle spatial (data sheet). <http://www.oracle.com/database/documents/spatial-ds.pdf>, March 1999.
- [11] Rene de la Briandais. File searching using variable length keys. In *Proceedings of the Western Joint Computer Conference*, pages 295–298, 1959.
- [12] David J. DeWitt, Navin Kabra, Jun Luo, Jignesh M. Patel, and Jie-Bing Yu. Client-Server Paradise. In *Proceedings of VLDB'94*, pages 558–569, Santiago, Chile, 1994.
- [13] A. A. Diwan, Sanjeeva Rane, S. Seshadri, and S. Sudarshan. Clustering techniques for minimizing external path length. In *Proceedings of VLDB'96*, pages 342–353, Mumbai (Bombay), India, September 1996.
- [14] Claudio Esperanca and Hanan Samet. Spatial database programming using sand. *Proceedings of the Seventh International Symposium on Spatial Data Handling*, May 1996.
- [15] K. P. Eswaran, J. N. Gray, R. A. Lorie, and I. L. Traiger. The notions of concurrency and predicate locks in a data base system. *Communications of the ACM*, 19(11), 1976.
- [16] Christos Faloutsos and Volker Gaede. Analysis of n-dimensional quadtrees using the Hausdorff fractal dimension. In *Proceedings of VLDB'96*, pages 40–50, 3–6 September 1996.
- [17] Christos Faloutsos, H. V. Jagadish, and Yannis Manolopoulos. Analysis of the n-dimensional quadtree decomposition for arbitrary hyperrectangles. *TKDE*, 9(3):373–383, 1997.
- [18] R. A. Finkel and J. L. Bentley. Quad trees: a data structure for retrieval on composite key. *Acta Informatica*, 4(1):1–9, 1974.
- [19] E. Fredkin. Trie memory. *Communications of the ACM*, 3:490–500, 1960.
- [20] Volker Gaede and Oliver Gunther. Multidimensional access methods. In *ACM Computer Surveys*, 30,2, pages 170–231, June 1998.
- [21] I. Gargantini. An effective way to represent quadtrees. *Communications ACM*, 1982, 25(12):905–910, 1982.
- [22] J. N. Gray. Notes on data base operating systems. In *Springer Verlag (Heidelberg, FRG and New York NY, USA) LNCS, 'Operating Systems, an Advanced Course', Bayer, Graham, Seegmuller(eds)*, volume 60. 1978.
- [23] J. N. Gray and A. Reuter. *Transaction Processing: concepts and techniques*. Data Management Systems. Morgan Kaufmann Publishers, Inc., San Mateo (CA), USA, 1993.
- [24] Ralf Hartmut Güting. An introduction to spatial database systems. *VLDB Journal*, 3(4):357–399, 1994.
- [25] A. Guttman. R-trees: a dynamic index structure for spatial searching. *Proceedings of SIGMOD'84*, pages 47–57, June 1984.
- [26] Joseph M. Hellerstein, Jeffrey F. Naughton, and Avi Pfeffer. Generalized search trees for database system. *Proceedings of VLDB'95*, 1995.

- [27] Joseph M. Hellerstein and Avi Pfeffer. The RD-tree: An index structure for sets. Technical report, University of Wisconsin Computer Science, 1994.
- [28] G. Kedem. The quad-CIF tree: A data structure for hierarchical on-line algorithms. In *ACM IEEE Nineteenth Design Automation Conference Proceedings*, pages 352–357, Los Alamitos, Ca., USA, June 1982.
- [29] A. Klinger. Pattern and search statistics. In *S. RUSTAGI Ed., Optimizing Methods in Statistics*, pages 303–337, 1971.
- [30] Donald E. Knuth. *The Art of Computer Programming, Vol. 3*. Addison-Wesley, Reading, 1973.
- [31] Marcel Kornacker and Douglas Banks. High-concurrency locking in R-trees. In *Proceedings of VLDB'95*, Zurich, Switzerland, Sept. 1995.
- [32] Marcel Kornacker, C. Mohan, and Joseph M. Hellerstein. Concurrency and recovery in generalized search trees. *Proceedings of SIGMOD'98*, pages 62–72, May 1998.
- [33] D. R. Morrison. PATRICIA - practical algorithm to retrieve coded in alphanumeric. *J. Assoc. Comput. Mach.*, 15(4):514–534, 1968.
- [34] R. C. Nelson and H. Samet. A consistent hierarchical representation for vector data. In *Computer Graphics (SIGGRAPH '86 Proceedings)*, volume 20(4), August 1986.
- [35] J. Nievergelt, H. Hinterberger, and K. Sevcik. The grid file: an adaptable symmetric multi-key file structure. *ACM Transactions on Database Systems*, 9(1):38–71, 1984.
- [36] Jack A. Orenstein. Multidimensional tries used for associative searching. *Information Processing Letters*, 14(4):150–157, June 1982.
- [37] Jack A. Orenstein and F.A. Manola. PROBE spatial data modeling and query processing in an image database application. *IEEE Transactions on Software Engineering*, 14(5):611–629, May 1988.
- [38] Dimitris Papadias, Nikos Mamoulis, and Vasilis Delis. Algorithms for querying by spatial structure. In *Proceedings of VLDB'98*, pages 546–557, New York City, New York, USA, August 1998.
- [39] H. Samet and R. E. Webber. Storing a collection of polygons using quadtrees. *ACM Transactions on Graphics, Volume 4, Issue 3*, 1985.
- [40] Hanan Samet. *Applications of Spatial Data Structures*. Addison-Wesley, 1990.
- [41] Hanan Samet. *The Design and Analysis of Spatial Data Structure*. Addison-Wesley, 1990.
- [42] Bernhard Seeger and Hans-Peter Kriegel. The Buddy-tree: An efficient and robust access method for spatial data base systems. In *Proceedings of VLDB'90*, Queensland, Australia, 1990.
- [43] T. Sellis, N. Roussopoulos, and C. Faloutsos. The R+ -tree: A dynamic index for multi-dimensional objects. In *Proceedings of VLDB'87*, Brighton, UK, September 1987.
- [44] Timos K. Sellis, Nick Roussopoulos, and Christos Faloutsos. Multidimensional access methods: Trees have grown everywhere. In *Proceedings of VLDB'97*, pages 13–14, 1997.
- [45] A. Szalay, P. Kunszt, A. Thakar, J. Gray, D. Slutz, and R. Brunner. Designing and mining multi-terabyte astronomy archives: The sloan digital sky survey. In *Proceedings of SIGMOD'00*, pages 451–462, Dallas TX, May 2000.
- [46] S. Tanimoto and T. Pavlidis. A hierarchical data structure for picture processing. *Computer Graphics and Image Processing*, 4(2):104–119, June 1975.

```

PROCEDURE ReCluster-Bottom Up(TreeNode root)
BEGIN
  S={};
  FOR each node n in delete-list (List of deleted nodes) DO
    Remove n from its current page
    IF n is the last node in the page
      THEN delete the page.
      decluster(root);
    END FOR
  /* S is now the set of nodes that are affected */
  LOOP WHILE there are nodes in S that are not yet processed
    Choose an affected node P that is either a leaf or
    all of whose children are either not in S or have been processed
    process-node(P);
  END LOOP
END

PROCEDURE decluster(node n)
BEGIN
  add n to S
  IF (n is not a new inserted node) THEN
    remove n from its current page.
    IF n is the last node in the page
      THEN delete the page.
  FOR each child n1 of n DO
    IF (n1 is a new inserted node or if the subtree from n1 is modified)
      THEN decluster(n1)
    ELSE
      IF (n1 is in the same page as n)
        THEN move n1 and all its descendants in the same
        page as n to a new page
    END FOR
  END

PROCEDURE process-node(TreeNode P)
BEGIN
  IF P is a leaf node
    THEN create a new page C containing node P.
  ELSE Let P1 . . . Pn, be the children of P.
    Let C1 . . . Cn be the pages containing P1 . . . Pn, respectively.
    Let P11 . . . P1m, be the children among the above
    whose page height is the greatest.
    IF node P and the contents of the pages G11 . . . G1m can be merged in 1 page
    THEN merge the contents of C11 . . . C1m, into a
    new page C and delete (C11 . . . C1m).
    ELSE create a new page C containing only P
  END

```

Figure 12: Minimum height clustering algorithm