

R-Trees

Group 6: ARR-CHEESE

Trần Tôn Minh Kỳ, Phan Anh Khoa, Lê Lương Quốc Trung,
Phạm Võ Khang Duy

March 2025

Abstract

Spatial data indexing requires multidimensional queries and modifications. The R-Tree, proposed by Antonin Guttman in 1984, attempts to solve this dilemma with worst case linearithmic time complexity. Due to the robust implementation on main memory and disk, it is extremely practical in real world applications and opens the door to additional variants in an attempt to counter technical issue the traditional R-Tree suffers from.

Contents

1	Introduction	2
2	R-Trees	2
2.1	Definitions	2
2.2	Basic implementation and operations	2
2.3	C++ implementation	6
2.4	Analysis	11
2.5	Application	12
3	Exercises	12
3.1	Quizzes	12

1 Introduction

Spatial data structures present a rather different challenge from those of one-dimensional. Suppose that we wish to seek a restaurant's xy coordinates. One may propose that we index the x -coordinate and y -coordinate in separate data structures and thus perform basic operations on one of them at a time. Because the search space is multi-dimensional, traditional one-dimensional data structures such as BSTs, B-trees are not well-suited for these purposes [1]. Thus the dilemma requires storing spatial attributes that can efficiently answers multidimensional queries. To this end, we present the R-Tree.

In the succeeding section, the first **Subsection 2.1** will introduce what is an R-Tree and its general implementation. **Subsection 2.2** and **2.3** will delve deeper into its implementation; the former will outline them in pseudocode and the latter will realize that outline using C++. **Subsection 2.4** will analyze the spacetime complexity of our implementations, whereas **Subsection 2.5** will dive in the practical purposes of R-Tree. **Subsection 2.6** will conclude this report with a few drawbacks of R-Tree and some other variants of R-Trees.

2 R-Trees

2.1 Definitions

An R-Tree is a height-balanced tree similar to a B-Tree [1] (all nodes remain sorted and all operations are in logarithmic time). R-Trees cater to spatial objects such as coordinates and thus points, rectangles, polygons, etc. As such, it is a suitable data structure to index objects in main memory and secondary storage. It is also capable of storing a multitude of data types using a concept known as Minimum Bounding Rectangle (henceforth will be referred to as MBRs) [3].

The R-Tree nodes are implemented as disk pages, where each node is an MBR bounding its children. The leaf nodes are pointers to the objects being indexed (indeed, the MBRs can overlap and contain multiple nodes). It should be noted that each MBR corresponds to each nodes and thus a search query requires traversing through many nodes before reaching its target node, if one exists in the first place [2].

2.2 Basic implementation and operations

We begin with the first fundamental operation for R-Tree, that is the search algorithm. To be more specific, the **SEARCH** algorithm returns an array of all nodes whose MBRs overlap with the query MBR.

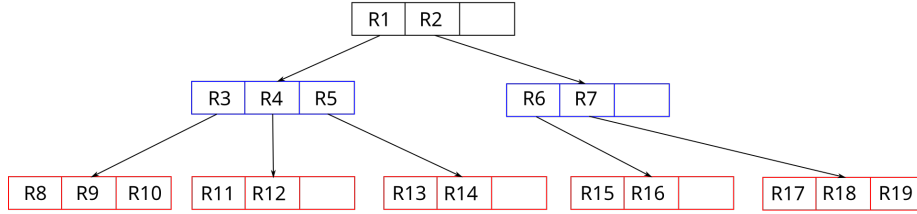
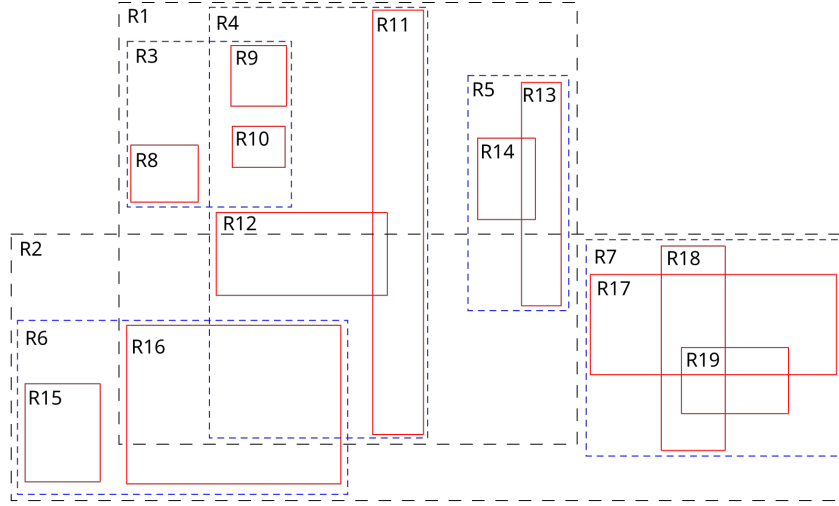


Figure 1: A simple R-Tree of 2D data

Algorithm SEARCH [1].

1. If the root is not a leaf node, for each entry, check whether their MBR overlap the searching MBR. For each satisfied node, call **SEARCH** on the corresponding subtree.
2. If the root is a leaf node, for each entry, check whether their MBR overlap the searching MBR. Those that satisfied the condition belong to the result.

We now describe the operation of splitting nodes in an R-Tree. Let M be the largest number of children and $m \leq M/2$ be the smallest number of children a node can contain. When a node overflows, that is, has more than M children, the operation **SPLITNODE** is invoked. Here we will describe the linear splitting algorithm, which offers computational affordability but suboptimal query performance. We will now define a helper operation.

Algorithm PICKSEEDS [1].

1. Along each dimension, choose the node whose MBR has the highest low side and lowest high side and record their difference gap.
2. Divide the width of the recorded set along the corresponding dimension, effectively normalizing the gaps.
3. Select the pair with the greatest gap on any dimension.

Now we are ready to define the operation **SPLITNODE**.

Algorithm SPLITNODE [1].

1. Invoke **PICKSEEDS** to choose 2 entries as the first element of each child node.
2. If all entries have been selected, stop. If one child node has below m entries, assign other entries to it and stop.
3. Choose the next entry to assign and add it to the child node with least enlarged MBR requirement. Repeat step 2.

The algorithm **ADJUSTTREE** updates MBRs by propagating from a leaf node to the root and performing node splits if necessary.

Algorithm ADJUSTTREE [1].

1. Let cur be a copy of the query leaf node. If the query was previously split, let $cur2$ be a copy of the second leaf node.
2. If cur is the root, stop.
3. Let $parent$ be the parent of cur . Adjust the $parent$'s corresponding entry's MBR to suit the MBRs of cur .
4. Do the same to $cur2$ using $parent2$. Add its entry's MBR to the $parent2$, in case of overfull, invoke **SPLITNODE** to produce two nodes to accommodate the children.
5. Set $cur = parent$ and $cur2 = parent2$. Repeat from step 2.

We are now ready to define our second fundamental operation, that is the **INSERT** operation.

Algorithm INSERT [1].

1. Let *cur* be the root node.
2. If *cur* is a leaf node, skip to step 4.
3. Set *cur* to the child node whose MBR requires smallest enlargement to accommodate the new entry's MBR. If multiple are found, choose one with smallest area. Repeat from step 2.
4. If *cur* does not overflow when a new node is inserted, do so and stop. Otherwise, invoke **SPLITNODE** to accommodate the new node.
5. Invoke **ADJUSTTREE** on *cur*.
6. If the root is split, set a new root whose children are the product of the split root.

We now turn to our third and final fundamental operation, which is deleting a node from a tree. As per usual, we will define some helper operations to simplify the implementation. The operation **FINDLEAF** returns the leaf node containing the query node.

Algorithm FINDLEAF [1].

1. If the query node is a leaf node, check each entry of the node. If a match is found, return it.
2. Otherwise, check each entry of the node whether their MBRs overlap the query MBR. For each matching node, invoke **FINDLEAF** on them.

The operation **CONDENSETREE** is similar to **ADJUSTTREE** in that it reconfigures the tree after modifying the entries, i.e. deletion in the former and insertion in the latter.

Algorithm CONDENSETREE [1].

1. Let *cur* be the query leaf node. Let *deleted* be an array of deleted nodes.
2. If *cur* is the root, skip to step 6. Otherwise, let *parent* be the parent node of *cur*.
3. If *cur* underflows, delete its corresponding entries and append *cur* to *deleted*.
4. If *cur* is not deleted, adjust its MBR.
5. Set *cur* to *parent* and repeat from step 2.
6. Re-insert all entries from *deleted*.

We are now ready to define the operation **DELETE**.

Algorithm DELETE [1].

1. Invoke **FINDLEAF** to locate the query node. If none is found, stop.
2. Remove the entry from the leaf node.
3. Invoke **CONDENSETREE** on the leaf node.
4. If the root node has only one child after tree condensation, make the child the new root.

2.3 C++ implementation

```
1 // RTree.h
2
3 #ifndef TREE_H
4 #define TREE_H
5
6 #include <vector>
7
8 struct MBR {
9     int x_min, y_min, x_max, y_max;
10     MBR();
11     MBR(int x_min, int y_min, int x_max, int y_max);
12     int area() const;
13     MBR merge(const MBR& other) const;
14     bool contains(const MBR& other) const;
15     bool intersects(const MBR& other) const;
16 };
17
18 struct Node {
```

```

19     bool is_leaf;
20     MBR mbr;
21     std::vector<Node*> children;
22     Node(bool is_leaf = false);
23 };
24
25 struct RTree {
26     Node* root;
27     int max_children;
28
29     RTree(int max_children = 4);
30     ~RTree();
31     void insert(const MBR& mbr);
32     void search(const MBR& search_mbr);
33     void find_area();
34     void display_tree();
35
36 private:
37     void insert_recursive(Node* node, const MBR& mbr, Node* parent
38     = nullptr);
39     Node* choose_best_node(Node* node, const MBR& mbr);
40     void split_node(Node* node, Node* parent);
41     void update_mbr(Node* node);
42     void search_recursive(Node* node, const MBR& search_mbr, std::
43     vector<Node*>& results);
44     void find_area_recursive(Node* node, int level);
45     void delete_node(Node* node);
46     void display_tree_recursive(Node* node, int level = 0);
47 };
48
49 #endif
50
51 // RTree-functions.cpp
52
53 #include <iostream>
54 #include <algorithm>
55 #include "RTree.h"
56
57 MBR::MBR() : x_min(0), y_min(0), x_max(0), y_max(0) {}
58 MBR::MBR(int x_min, int y_min, int x_max, int y_max) : x_min(
59     x_min), y_min(y_min), x_max(x_max), y_max(y_max) {}
60 int MBR::area() const { return (x_max - x_min) * (y_max - y_min);
61 }
62 MBR MBR::merge(const MBR& other) const {
63     return MBR(std::min(x_min, other.x_min), std::min(y_min, other.
64     y_min),
65     std::max(x_max, other.x_max), std::max(y_max, other.y_max));
66 }
67 bool MBR::contains(const MBR& other) const {
68     return x_min <= other.x_min && y_min <= other.y_min && x_max >=
69     other.x_max && y_max >= other.y_max;
70 }
71 bool MBR::intersects(const MBR& other) const {
72     return !(x_max < other.x_min || x_min > other.x_max || y_max <
73     other.y_min || y_min > other.y_max);
74 }
75
76 // Node constructor

```

```

22 Node::Node(bool is_leaf) : is_leaf(is_leaf) {}
23
24 // RTree methods
25 RTree::RTree(int max_children) : max_children(max_children) {
26     root = new Node(true);
27 }
28
29 RTree::~RTree() {
30     delete_node(root);
31 }
32
33 void RTree::delete_node(Node* node) {
34     if (!node) return;
35     for (Node* child : node->children) {
36         delete_node(child);
37     }
38     delete node;
39 }
40
41 void RTree::insert(const MBR& mbr) {
42     insert_recursive(root, mbr, nullptr);
43 }
44
45 void RTree::insert_recursive(Node* node, const MBR& mbr, Node*
parent) {
46     if (node->is_leaf) {
47         if (node->children.size() < max_children) {
48             Node* new_node = new Node(true);
49             new_node->mbr = mbr;
50             node->children.push_back(new_node);
51             update_mbr(node);
52             if (parent) update_mbr(parent);
53         }
54         else {
55             if (!parent) { // Root case
56                 Node* new_parent = new Node(false);
57                 split_node(node, new_parent);
58                 root = new_parent;
59                 Node* target = choose_best_node(new_parent, mbr);
60                 insert_recursive(target, mbr, new_parent);
61             }
62             else { // Non-root case
63                 split_node(node, parent);
64                 Node* target = choose_best_node(parent, mbr);
65                 insert_recursive(target, mbr, parent);
66             }
67         }
68     }
69     else {
70         Node* target = choose_best_node(node, mbr);
71         insert_recursive(target, mbr, node);
72     }
73 }
74
75 Node* RTree::choose_best_node(Node* node, const MBR& mbr) {
76     if (node->children.empty()) return node;
77     Node* best = node->children[0];

```



```

78     int min_enlargement = (node->children[0]->mbr.merge(mbr)).area()
    - node->children[0]->mbr.area();
79     for (size_t i = 1; i < node->children.size(); ++i) {
80         int enlargement = (node->children[i]->mbr.merge(mbr)).area()
    - node->children[i]->mbr.area();
81         if (enlargement < min_enlargement) {
82             min_enlargement = enlargement;
83             best = node->children[i];
84         }
85     }
86     return best;
87 }
88
89 void RTree::split_node(Node* node, Node* parent) {
90     int median = node->children.size() / 2;
91     Node* new_node = new Node(node->is_leaf);
92
93     new_node->children = std::vector<Node*>(node->children.begin()
    + median, node->children.end());
94     node->children = std::vector<Node*>(node->children.begin(),
    node->children.begin() + median);
95
96     update_mbr(node);
97     update_mbr(new_node);
98
99     if (parent->children.empty()) {
100         parent->children.push_back(node);
101         parent->children.push_back(new_node);
102     }
103     else {
104         for (size_t i = 0; i < parent->children.size(); ++i) {
105             if (parent->children[i] == node) {
106                 parent->children[i] = node;
107                 parent->children.insert(parent->children.begin() + i + 1,
    new_node);
108                 break;
109             }
110         }
111     }
112     update_mbr(parent);
113 }
114
115 void RTree::update_mbr(Node* node) {
116     if (node->children.empty()) return;
117     node->mbr = node->children[0]->mbr;
118     for (size_t i = 1; i < node->children.size(); ++i) {
119         node->mbr = node->mbr.merge(node->children[i]->mbr);
120     }
121 }
122
123 void RTree::search(const MBR& search_mbr) {
124     std::vector<Node*> results;
125     search_recursive(root, search_mbr, results);
126     if (results.empty()) {
127         std::cout << "No MBR found in the search range." << std::endl
    ;
128     }

```

```

129     else {
130         for (const Node* node : results) {
131             std::cout << "Found MBR: (" << node->mbr.x_min << ", " <<
node->mbr.y_min
132             << ", " << node->mbr.x_max << ", " << node->mbr.y_max <<
")\n";
133         }
134     }
135 }
136
137 void RTree::search_recursive(Node* node, const MBR& search_mbr,
std::vector<Node*>& results) {
138     if (node->is_leaf) {
139         for (Node* child : node->children) {
140             if (child->mbr.intersects(search_mbr)) {
141                 results.push_back(child);
142             }
143         }
144     }
145     else {
146         for (Node* child : node->children) {
147             if (child->mbr.intersects(search_mbr)) {
148                 search_recursive(child, search_mbr, results);
149             }
150         }
151     }
152 }
153
154 void RTree::find_area() {
155     if (root->children.empty()) {
156         std::cout << "Tree is empty." << std::endl;
157         return;
158     }
159     std::cout << "Areas of all MBRs in the R-tree:\n";
160     find_area_recursive(root, 0);
161 }
162
163 void RTree::find_area_recursive(Node* node, int level) {
164     std::string indent(level * 2, ' '); // Indentation based on
level
165
166     // Print the current node's MBR area
167     std::cout << indent << (node->is_leaf ? "Leaf" : "Internal") <<
" Node MBR ("
168     << node->mbr.x_min << ", " << node->mbr.y_min << ", "
169     << node->mbr.x_max << ", " << node->mbr.y_max << "): "
170     << node->mbr.area() << std::endl;
171
172     // If leaf, print areas of all child MBRs
173     if (node->is_leaf) {
174         for (size_t i = 0; i < node->children.size(); ++i) {
175             Node* child = node->children[i];
176             std::cout << indent << " Child " << i + 1 << " MBR ("
177             << child->mbr.x_min << ", " << child->mbr.y_min << ", "
178             << child->mbr.x_max << ", " << child->mbr.y_max << "): "
179             << child->mbr.area() << std::endl;
180         }

```

```

181     }
182     // If internal, recurse into children
183     else {
184         for (Node* child : node->children) {
185             find_area_recursive(child, level + 1);
186         }
187     }
188 }
189
190 void RTree::display_tree() {
191     if (root->children.empty()) {
192         std::cout << "Tree is empty." << std::endl;
193         return;
194     }
195     std::cout << "R-tree Structure:\n";
196     display_tree_recursive(root, 0);
197 }
198
199 void RTree::display_tree_recursive(Node* node, int level) {
200     if (!node) return;
201
202     // Indentation based on depth
203     for (int i = 0; i < level; ++i) {
204         std::cout << " ";
205     }
206
207     // Print the node's MBR
208     std::cout << "|-- MBR (" << node->mbr.x_min << ", " << node->
mbr.y_min
209         << ") to (" << node->mbr.x_max << ", " << node->mbr.y_max
<< ")";
210
211     if (node->is_leaf) {
212         std::cout << " [Leaf]";
213     } else {
214         std::cout << " [Internal]";
215     }
216     std::cout << std::endl;
217
218     // Recursively print children
219     for (Node* child : node->children) {
220         display_tree_recursive(child, level + 1);
221     }
222 }

```

2.4 Analysis

Let m and n be the minimum number of children a node can contain and the number of nodes in the tree respectively. The operation **SEARCH** has a worst time complexity of $\Theta(n)$, where all nodes must be traversed. In the best case where there are no overlap, the best time complexity is $\Omega(\log_m n)$. One can observe that the average time complexity depends on the degree of overlapping nodes; that is, the node splitting method. The space required to store MBRs is $O(n)$ and since utilize recursion, the worst case space complexity is $\Theta(n)$ due

to recursion stack.

As shown in the implementation, the operation **PICKSEEDS** has a time complexity of $O(\log_m n)$. In the case where splitting is required, the operation **SPLITNODE** is invoked, taking $\Theta(m)$ time in the worst case. We conclude that the operation **INSERT** has a time complexity of $O(\log_m n)$ on average and $\Theta(m \log_m n)$ in the worst case. It also uses $O(n)$ space to store n objects. In a similar line of reasoning, the operation **DELETE** also has a time complexity of $\Theta(m \log_m n)$ in the worst case and a space complexity of $O(n)$.

2.5 Application

The R-Tree has a myriad of applications in the real world. A couple of questions posed regarding this is finding the closest restaurants in one's current location, or locating which bus crossed a region at some time interval. This in turn has raised an emergence to additional variants of the R-Tree, including but not limited to:

- R*-tree
- R⁺-tree
- 2+3-tree
- RT-tree.

3 Exercises

We conclude this report with a handful of mental gymnastics for the readers.

3.1 Quizzes

Question 1. What is the primary purpose of an R-tree data structure?

- (a) To store one-dimensional numerical data.
- (b) To manage multi-dimensional spatial data.
- (c) To store text-based data.
- (d) To perform fast searching on strings.

Answer: (b) To manage multi-dimensional spatial data.

Question 2. In an R-tree, the leaf nodes contain

- (a) only pointers to the child nodes.
- (b) bounding boxes that cover regions.
- (c) actual spatial data or references to them.
- (d) pairs of keys and values.

Answer: (c) actual spatial data or references to them.

Question 3. What type of data is most commonly indexed using an R-tree?

- (a) Linear numerical data.
- (b) Multi-dimensional spatial data (e.g., geographic coordinates, rectangles).
- (c) Text strings.
- (d) Time series data.

Answer: (b) Multi-dimensional spatial data (e.g., geographic coordinates, rectangles).

Question 4. Which of the following is a key characteristic of an R-tree?

- (a) It is always unbalanced.
- (b) It can only handle one-dimensional data.
- (c) It uses bounding boxes to represent groups of spatial objects.
- (d) It stores data as linked lists.

Answer: (c) It uses bounding boxes to represent groups of spatial objects.

Question 1. Count MBRs in an R-Tree

In R-Tree, a node's MBR is the smallest axis-aligned rectangle that encloses all the data points or sub-regions within that node.

Given a root of an R-Tree, you are asked to return the total number of MBRs.

Constraint: The R-Tree has a maximum fanout (number of children per node) of 4, and you must count the total number of MBRs without using recursion (e.g., use an iterative approach like a stack or queue).

```
1  int RTree::count_mbrs() {
2      return count_mbrs_recursive(root);
3  }
4
5  int RTree::count_mbrs_recursive(Node* node) {
6      if (!node) return 0;
7      if (node->is_leaf) {
8          return node->children.size(); //Leaf nodes store MBRs
          directly
9      }
10     int total=0;
11     for (Node* child : node->children) {
12         total += count_mbrs_recursive(child);
13     }
14     return total;
15 }
```

Question 2. Find overlapping MBRs in a Range

Two MBRs are called overlapped when they enclose the same space in the tree.

Given a R-Tree, its root and a range (a rectangle with coordinates), you are asked to return the number of overlapping MBRs in that range, and coordinates of each MBR.

Constraint: The R-Tree contains at least 10,000 nodes, and you must optimize the solution to run in $O(\log n + k)$ time complexity, where k is the number of overlapping MBRs, assuming a balanced R-Tree.

```
1  void RTree::count_overlaps(const MBR& range) {
2      std::vector<Node*> overlaps;
3      search_recursive(root, range, overlaps);
4      std::cout << "Number of MBRs overlapping with (" << range.x_min
5      << ", " << range.y_min << ", " << range.x_max << ", " << range
6      .y_max << "): " << overlaps.size() << std::endl;
7      for (const Node* node : overlaps) {
8          std::cout << " Overlapping MBR: (" << node->mbr.x_min << ",
9          " << node->mbr.y_min << ", " << node->mbr.x_max << ", " << node
10         ->mbr.y_max << ")\n";
11     }
12 }
```

References

- [1] Antonin Guttman. “R-Trees: A Dynamic Index Structure for Spatial Searching”. In: *SIGMOD’84, Proceedings of Annual Meeting, Boston, Massachusetts, USA, June 18-21, 1984*. Ed. by Beatrice Yormark. ACM Press, 1984, pp. 47–57. DOI: 10.1145/602259.602266.
- [2] Yannis Manolopoulos et al. *R-Trees: Theory and Applications*. Advanced Information and Knowledge Processing. Springer, 2006. ISBN: 978-1-85233-977-7. DOI: 10.1007/978-1-84628-293-5.
- [3] Apostolos N. Papadopoulos et al. “R-Tree (and Family)”. In: *Encyclopedia of Database Systems, Second Edition*. Ed. by Ling Liu and M. Tamer Özsu. Springer, 2018. DOI: 10.1007/978-1-4614-8265-9_300.