

Final Project - Inverted Index and Comparative Analysis

Arthur Rabello Oliveira¹, Gabrielle Mascarelo, Eliane Moreira, Nicolás Spaniol, Gabriel Carneiro

Abstract

We present the implementation and comparative analysis of three fundamental data structures for indexing and searching textual documents: the Classic Binary Search Tree (BST), the AVL Tree, and the Red-Black Tree (RBT). Each structure was implemented with its core operations, including insertion and search. Unit tests were developed to validate the correctness and performance of these implementations. We also provide a further comprehensive comparative study of the three trees based on their time complexity, balancing efficiency, and suitability for document indexing. The results demonstrate the trade-offs between implementation complexity and query performance, offering insights into the practical considerations for choosing appropriate search tree structures in information retrieval systems.

Contents

1. Introduction	3
1.1. Context	3
1.2. Problem Statement	3
1.3. Objectives	3
2. Motivation	3
2.1. Why Inverted Index?	3
3. Implementations	3
3.1. Binary Search Tree (BST)	3
3.1.1. Algorithms	4
3.1.2. Complexity Analysis	4
3.2. AVL Tree	4
3.2.1. Algorithms	4
3.2.2. Complexity Analysis	5
3.3. Red-Black Tree (RBT)	6
3.3.1. Algorithms	6
3.3.2. Complexity Analysis	6
3.4. Inverted Index	6
3.4.1. Algorithms	6
3.4.2. Complexity Analysis	6
4. Testing and Validation	6
4.1. Unit Testing Method	6
5. Comparative Analysis	6
5.1. The Experiment	6
5.2. Graphs	6
5.2.1. AVL	6
5.2.2. BST	8
5.2.3. RBT	9
6. Conclusion	10
6.1. Raw stats	10

¹Escola de Matemática Aplicada, Fundação Getúlio Vargas (FGV/EMAp), email: arthur.oliveira.1@fgv.edu.br

6.2. Actual Analysis	10
6.2.1. Time Complexity	11
6.2.2. Memory Usage	11
7. Source code	11
8. Task Division (Required by the professor)	12
8.1. Arthur Rabello Oliveira	12
8.2. Gabrielle Mascarelo	12
8.3. Eliane Moreira	12
8.4. Nicolas Spaniol	12
8.5. Gabriel Carneiro	12
Bibliography	12

1. Introduction

1.1. Context

Humanity now produces more text in a single day than it did in the first two millennia of writing combined. Search engines must index billions of web pages, e-commerce sites hundreds of millions of product descriptions, and DevOps teams terabytes of log lines every hour. Scanning those datasets sequentially would be orders of magnitude too slow; instead, virtually all large-scale retrieval systems rely on an **inverted index**, a data structure that stores, for each distinct term, the identifiers of documents in which it appears.

1.2. Problem Statement

While the logical view of an inverted index is a simple dictionary, its physical realisation must support two conflicting workloads:

- **Bulk ingestion** of millions of (term, docID) pairs per second.
- **Sub-millisecond** queries for ad-hoc keyword combinations.

Choosing the proper data structure is therefore a trade-off between build-time cost and implementation complexity.

1.3. Objectives

1. **Implement** BST, AVL and Red-Black Tree insertion, deletion and search in C++.
2. **Build** an inverted index over a 10 k-document corpus with each tree.
3. **Measure** build time, query latency, and memory usage under identical workloads.
4. **Discuss** which structure best balances simplicity and performance for mid-scale information-retrieval tasks.

2. Motivation

2.1. Why Inverted Index?

Domain	Real-world system	Role of the inverted index
Web search	Google, Bing, DuckDuckGo	Core term → page mapping
Enterprise search	Apache Lucene & Elasticsearch	Underlying index files and query engine
Database systems	Postgres GIN & CockroachDB	Full-text and JSONB indexing
Observability / Logs	Splunk, OpenObserve	Fast filtering / aggregation of terabyte-scale logs
Bioinformatics	VariantStore, PAC	Searching billions of DNA k-mers
Operating systems	Linux schedulers & timers	Kernel subsystems use RBTs—conceptually an inverted index over time or PID keys

These examples shows the ubiquity of inverted indexes in modern era. From web search engines to bioinformatics, inverted indexes are the backbone of efficient information retrieval systems.

3. Implementations

3.1. Binary Search Tree (BST)

Since the AVL and RBT trees are both subsystems of the classic BST, we have implemented the classic BST in the `tree_utils` module, later used as based for the AVL and BST

3.1.1. Algorithms

[Here](#) are all the functions written for the classic BST, and [here](#) is the header file with the corresponding documentation. The list of functions is:

- `createNode`
- `createTree`
- `search`
- `deletionPostOrder`
- `destroy`
- `calculateHeight`
- `updateHeightUp`

3.1.2. Complexity Analysis

Below is a full complexity analysis:

createNode:

Clearly $O(1)$, the function allocates memory for a new node and initializes it, independent of the size of the tree.

createTree:

Also $O(1)$, an empty BST is allocated and initialized with a `nullptr` root.

search:

The search operation, unavoidably, has a time complexity of $O(h)$, with h being the height of the tree. The function follows a single root-to-leaf path in the BST, making at most h comparisons. No recursion, only a few local variables.

deletionPostOrder:

This function is $O(\varphi)$, where φ is the size of the subtree rooted at the node to be deleted. It is a classic post-order traversal: each node is visited & deleted precisely once.

destroy:

This function simply calls `deletionPostOrder` on the root node, so it is $O(h)$.

calculateHeight:

This is $O(k)$, with $k = \text{subtree size}$. It recursively explores both children of every node once to compute $\max(\text{left}, \text{right}) + 1$.

updateHeightUp:

This function is $O(h)$. It iterates upward, recomputing height until it stops changing or reaches the root; at most h ancestor steps, no recursion on children.

3.2. AVL Tree

3.2.1. Algorithms

The functions implemented *stritly* for the AVL can be found [here](#), and the header file with the corresponding documentation [here](#). We have used many of the BST functions, as previously stated. The list of AVL-functions is:

- `getHeight(Node*)`
- `bf(Node*)` [balance factor]
- `leftRotation(BinaryTree&, Node*)`
- `rightRotation(BinaryTree&, Node*)`
- `insert(BinaryTree& binary_tree, const std::string& word, int documentId)`
- `remove(Node*&, key)`
- `remove(BinaryTree&, key)`
- `updateHeightUp(Node*)`
- `clear(BinaryTree*)`
- `printInOrder`
- `printLevelOrder`

3.2.2. Complexity Analysis

Below is a full complexity analysis of the AVL functions:

getHeight:

$O(1)$, it simply returns the stored height of a node.

bf(Node*):

Also $O(1)$, this is nothing more than a subtraction of two integers.

leftRotation(BinaryTree&, Node*):

This is a constant quantity of pointer rewires and some height updates, so $O(1)$.

rightRotation(BinaryTree&, Node*):

Same as above, $O(1)$.

insert(BinaryTree& binary_tree, const std::string& word, int documentId):

This is $O(h)$, where h is the height of the tree. It does one BST descent $O(h)$ at most one rebalance per level. The rotations are constant.

remove(Node*&, key):

This is $O(h)$. It is a classic BST deletion + at most one structural deletion + up-to-root rebalancing.

remove(BinaryTree&, key):

Same as above, $O(h)$.

updateHeightUp(Node*):

This function iterates upward recomputing height until the value stabilises or the root is reached. Therefore it is $O(h)$.

clear(BinaryTree*):

Clearly $O(n)$, where n is the size of the tree. It removes each node once.

printInOrder:

This is a classic traversal visiting every node once, so it is $O(n)$.

printLevelOrder:

Same as above, $O(n)$.

3.3. Red-Black Tree (RBT)

3.3.1. Algorithms

3.3.2. Complexity Analysis

3.4. Inverted Index

3.4.1. Algorithms

3.4.2. Complexity Analysis

4. Testing and Validation

4.1. Unit Testing Method

We have used the Unity submodule for unit testing (espanio faz teu nome)

All trees were equivalently tested under the same principles blablabla

The testing module can be found here blablabla

5. Comparative Analysis

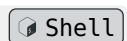
Here we develop a full comparative analysis of the three implemented data structures. We will use the following metrics:

1. Memory usage
2. Insertion time
3. Search time

5.1. The Experiment

For the graphs seen in Section 5.2, we have used the following:

```
1 ./[tree_name] stats 1000 data
```

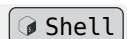


All trees were testes with the same first 1000 .txt file, which can be found here.

5.2. Graphs

For a better visualization of the statistics we will cover here, we recommend, through the repository, to run:

```
1 make bst && make avl && make rbt
2 ./(tree_name) stats 1000 data
```



Opening the generated JavaScript visualization in a browser, you will be able to see the graphs and the tree visualizations in a more interactive way.

All statistics are also available in the `(tree_name).csv` file, generated by the `stats` command.

5.2.1. AVL

These graphs were generated by thye JavaScript visualizer.

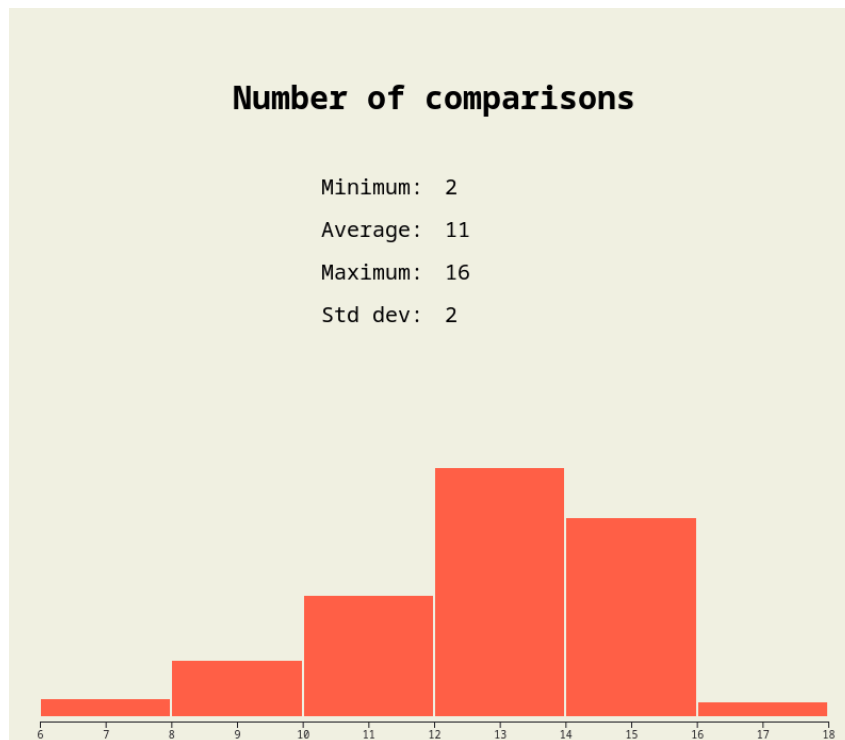


Figure 1: Number of Comparisons done by the AVL Tree

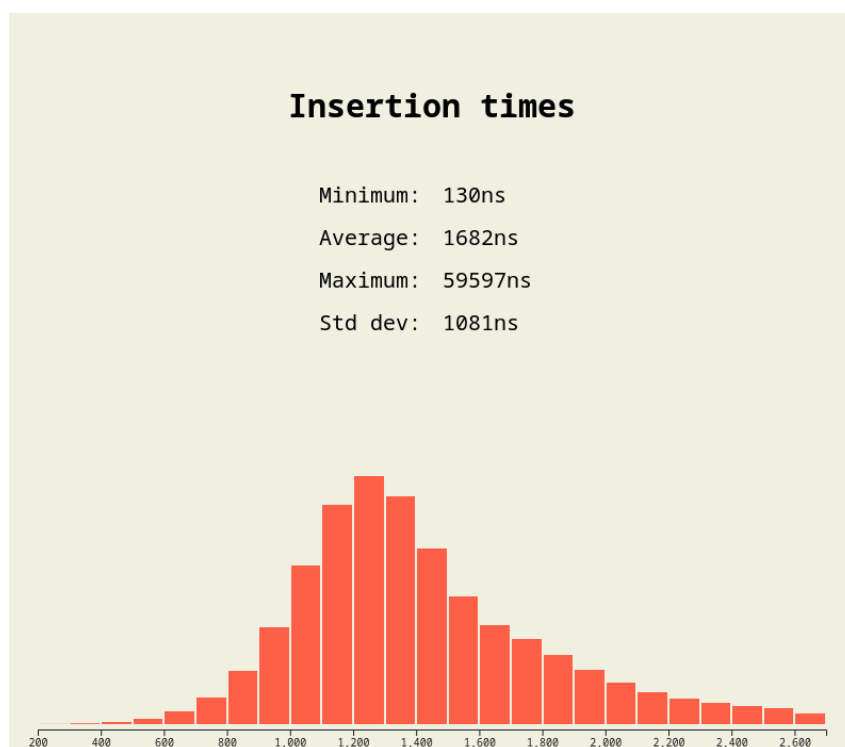


Figure 2: Insertion time of the AVL Tree

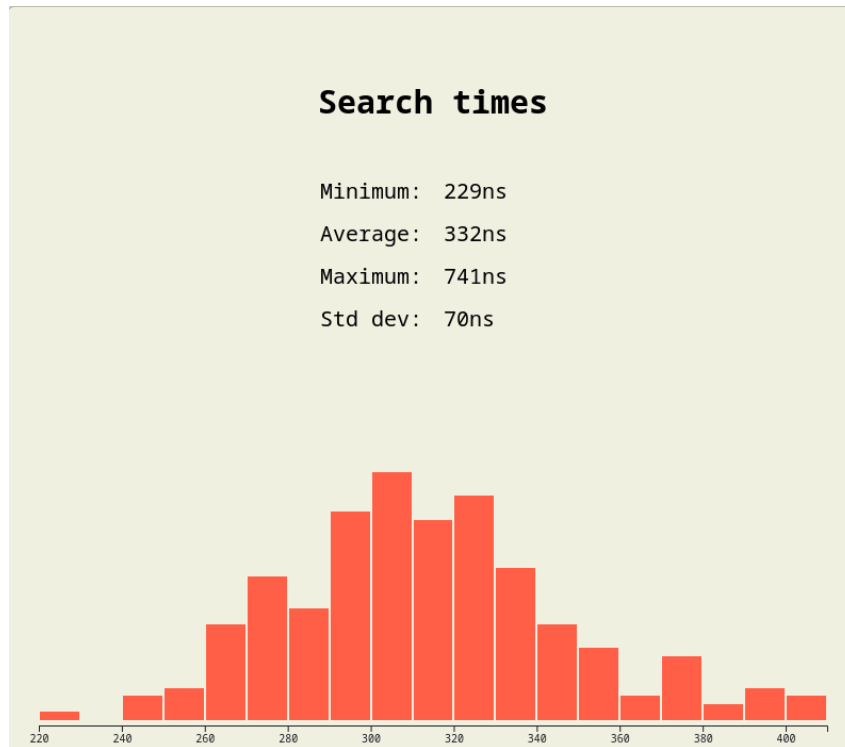


Figure 3: Search time of the AVL Tree

5.2.2. BST

These graphs were generated by thye JavaScript visualizer.

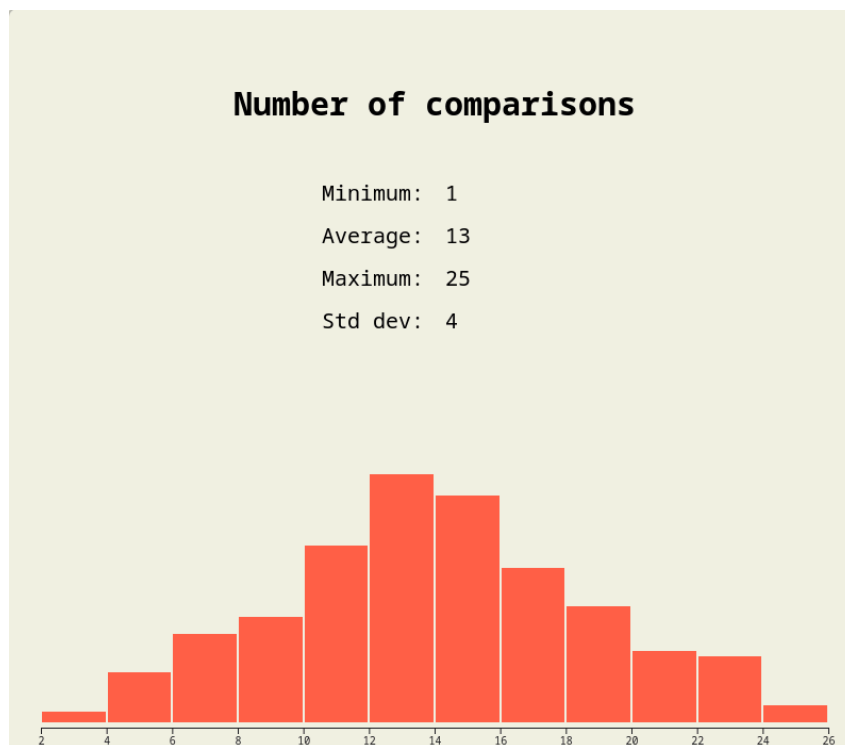


Figure 4: Number of Comparisons done by the BST Tree

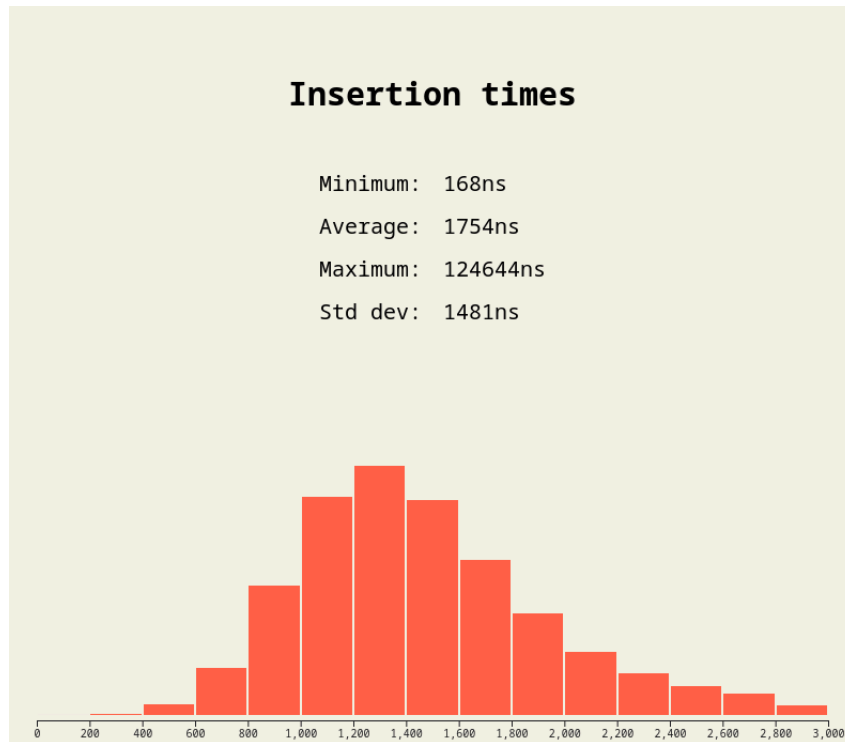


Figure 5: Insertion time of the BST Tree

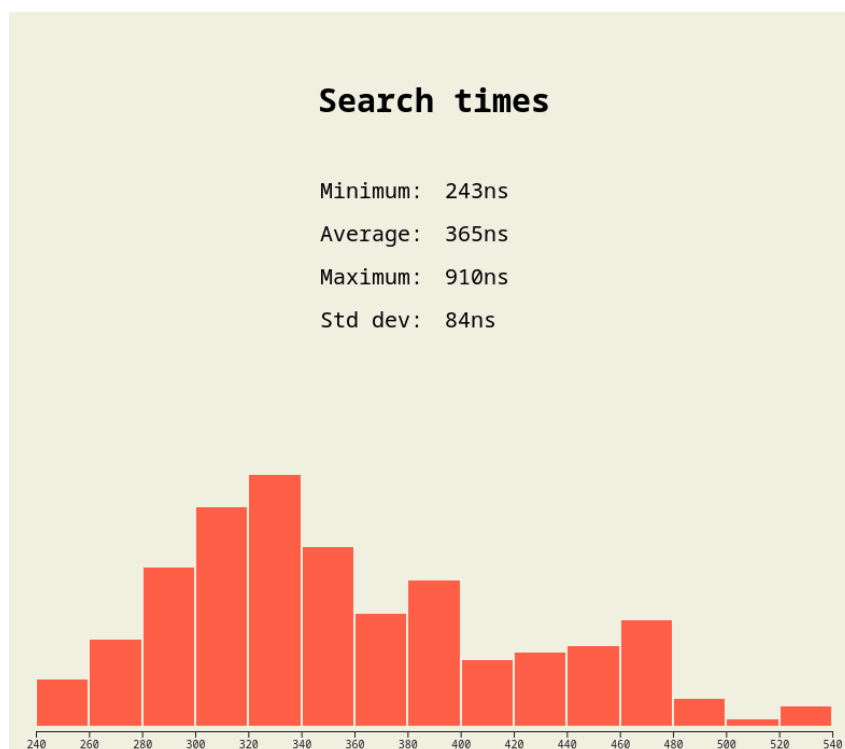


Figure 6: Search time of the BST Tree

5.2.3. RBT

These graphs were generated by thye JavaScript visualizer.

6. Conclusion

6.1. Raw stats

The raw statistics for the trees are:

Metric	BST	AVL	RBT
Docs indexed	1000	1000	0
Words indexed	213099	213099	0
Total insertion time (ns)	373834912	393149216	0
Avg insertion time (ns)	1754	1844	0
Max insertion time (ns)	124644	189938	0
Min insertion time (ns)	168	147	0
Total search time (ns)	90034	81607	0
Avg search time (ns)	365	331	0
Max search time (ns)	910	773	0
Min search time (ns)	243	232	0
Total comparisons (search)	3356	2942	0
Avg comparisons (search)	13	11	0
Max comparisons (search)	25	16	0
Min comparisons (search)	1	2	0
Node count	16986	16986	0
Tree height	32	16	0
Min depth	4	10	0
Balance diff	28	6	0
Relative balance	8	1.6	0

6.2. Actual Analysis

Let $h(n)$ be the height of node n . Because every AVL is a BST with an additional constraint, every theorem about BST search order also applies to AVL unless it contradicts the balance rule.

For the BST, If the input keys arrive in sorted order, each new key is inserted as the right child of the previous one, yielding a linear chain of n nodes, so:

$$h_{\max}(n) = n - 1 \quad 1$$

For the AVL, the maximum height is logarithmical [1], as it is a balanced tree. The height of an AVL tree with n nodes is at most:

$$h_{\max}(n) = O(\log n) \quad 2$$

Which might make the AVL interesting.

6.2.1. Time Complexity

When analyzing time complexity for these trees, we look at the cost of each operation:

Operation	BST	AVL	Proof Idea
Insertion	$O(h(n))$	$O(h) + O(1)$ rotations	Keys are compared on a root-to-leaf path, so the height of the tree determines the number of comparisons. For BST, this is linear in the worst case, while AVL guarantees logarithmic height due to balancing rules.
Search	$O(h(n))$	$O(h(n))$	Shown in Section 3.2.2 and Section 3.1.2
Deletion	$O(h(n))$	$O(h) + O(h)$ rotations worst-case	Each descent to find the node to delete is $O(h)$, and rebalancing may require up to $O(h)$ rotations in the worst case. For BST, this is linear in the worst case, while AVL guarantees logarithmic height due to balancing rules.

6.2.2. Memory Usage

All implementations use:

```

1 struct Node {
2     std::string word;
3     std::vector<int> docIds;
4     Node *left, *right, *parent;
5     int height;           // both codes already keep this
6     bool isRed;           // reserved for RBT
7 };

```

Field counts are identical, therefore heap consumption is $O(n)$ for the AVL and BST. Recursive algorithms allocate one activation record per visited level.

- BST worst-case stack depth: $h = n - 1 \rightarrow O(n)$ extra bytes; may overflow for large n .
- AVL stack depth: $h = O(\log n) \rightarrow$ asymptotically optimal.

No additional global buffers are required; rotations operate with $O(1)$ local variables.

7. Source code

All implementations and tests are available in the public repository at <https://github.com/arthurabello/dsa-final-project>

8. Task Division (Required by the professor)

8.1. Arthur Rabello Oliveira

Contributed with:

- Keeping the repository civilized (main-protecting rules, enforcing code reviews)
- Writing and documenting the `Makefile` for building the project
- Writing the `README.md` and `report.typ`
- Writing and documenting functions for the classic BST in `bst.cpp`

8.2. Gabrielle Mascarelo

Contributed with:

- Writing and documenting functions to read files in `data.cpp`
- Structuring statistics in the CLI

8.3. Eliane Moreira

Contributed with:

- Testing all function related to the BST
- Writing and documenting functions for `tree_utils.cpp`
- Fixing bugs in `data.h`

8.4. Nicolás Spaniol

Contributed with:

- Code reviews and suggestions for improvements in the codebase.
- Built from scratch the JavaScript visualization of all trees

8.5. Gabriel Carneiro

Contributed with:

- Writing and documenting functions for the classic BST
- Writing and documenting functions for the AVL Tree
- Testing functions for `tree_utils.cpp`

Bibliography

- [1] A. Tomescu, “AVL Trees Height Proof.” [Online]. Available: <https://people.csail.mit.edu/alinush/6.006-spring-2014/avl-height-proof.pdf>