# Exploring different merging algorithms for balanced trees and their time complexity optimization.

Changhui (Eric) Zhou

October 5, 2025

word count: ???

# Contents

# 1 Introduction

A *data structure* is a way to store and organize data in order to facilitate access and modifications (Cormen, Leiserson, Rivest, & Stein, 2022). Designing and choosing more efficient data structures has always been a great persuit for computer scientists, for optimal data structures can save huge amount of computing resources, especially in face of large amount of data. Basic data structures include ordered data structures like arrays, linked lists and binary search trees and unordered data structures like hashtables.

For ordered data structures, merging two or more instances of them while maintaining its ordered property may be frequently used in practice. For example, to investigate the factors affacting the school grade, data from different schools may be grouped and merged according to various factors. The efficiency of combination varies significantly based on the data structure itself and the algorithm used in the process.

This essay will focus on invesitgating the theoratical time complexity (need definitions aa) and actual performance of merging algorithms of different data structures, namely arrays, and BSTs, which are the most commonly used data structure in real life.

**Research question: How does different algorithm affect the efficiency of merging two instances of balanced search trees?**

# 2 Theory

## 2.1 Data structure terminology

When a homogeneous relation (a binary relation between two elements) $\leq$ on a set of elements $X$ satisfies:

1. Antisymmetry: $\quad \forall u, v \in X, (u \leq v \wedge v \leq u) \Leftrightarrow u = v.$
2. Totality: $\quad \forall u, v \in X, u \leq v \vee v \leq u.$
3. Transitivity: $\quad \forall u, v, w \in X, (u \leq v \wedge v \leq w) \Rightarrow u \leq w.$

We say $P = (X, \leq)$ is a total order. For example $P = (\mathbb{R}, \leq)$, where $\leq$ is numerical comparison, is a total order. The set of finite strings and

lexographical order comparison is also a total order. But $P = (\{S : S \subset \mathbb{R}\}, \subset)$ is not a total order.

Ordered data structures can store elements that satisfies a total order while maintaining their order. In C++, arrays, vectors, linked lists and sets can be ordered data structures, but unordered sets (hashtables) are not ordered data structures.

> **Definition**
>
> Ordered data structures are data structures that can store elements that satisfies a total order while maintaining their order.

## 2.2 Balanced binary search trees

- A *graph* $G = (V, E)$ is the combination of the vertex set $V$ and the edge set $E$.

- A *tree* $T = (V, E)$ is a connected acylic graph.

- A *binary tree* is a tree that has no more than two children for each node.

- A *binary search tree (BST)* is a binary tree, whose nodes contain values under a total order, that has the following property: For any node $v$, all nodes in its left subtree are less than $v$, and all nodes in its right subtree are greater than $v$ (Cormen et al., 2022).

Generally speaking, a balanced BST is a BST whose depth or the cost of iterating from the root to one specific leaf is strictly, expectedly or amortized $O(\log n)$. There are different kinds of balanced BSTs, like splay tree, treap, AVL trees and red-black trees. This essay will focus on AVL trees. AVL trees are a type of self-balancing binary search trees, which adjusts its shape through rotations and maintain the difference of the depths of two subtrees at most one (Karlton, Fuller, Scroggs, & Kaehler, 1976)[1].

In this essay, we will assume that there are two instances of AVL trees $T_1$ and $T_2$ to be merged. Without losing generality, we will assume $T_1$ has $n$ elements and $T_2$ has $m$ elements and $n \geq m$.

---

[1]In fact, this kind of BST was refered to as HB[1], but AVL trees nowadays are mainly HB[1]

## 2.3 Insertion-based merge

One of the basic operations supported by a balanced BST is insertion, where one element is added to the tree and the order of the tree is automatically maintained. In fact, merging two instances of BSTs can be reduced to a sequence of insertions to a balanced BST. To be more specific, we iterate through all the elements in $T_2$, insert them one by one into $T_1$, that would be $m$ operations with each having time complexity of $O(\log n)$, resulting in a overall complexity of $O(m \log n)$.

---

**Algorithm: Insertion-Based Merge**

**Require:** Two balanced binary search trees $T_1$ and $T_2$
**Ensure:** A single balanced binary search tree containing all elements
from $T_1$ and $T_2$
1: **procedure** INSERTIONBASEDMERGE($T_1, T_2$)
2:    **for all** elements $x$ in $T_2$ (in-order traversal) **do**
3:       $T_1.$ INSERT($x$)
4:    **return** $T_1$

---

This algorithm performs well when $m/n$ is small, as the overall time complexity will be mainly $O(\log n)$. However, when $m$ and $n$ are relatively at the same scale, the overall time complexity will be close to $O(n \log n)$.

## 2.4 In-order traversal merge

Another way to merge two instances of BSTs is to utilize the property that each instance is already in-order. To combine them, we can view this process as merging two sorted subarrays into a new array, just like a merge sort. The iteration and new array construction process take $O(n + m)$ time. With proper construction function, we can create a balanced BST in linear time out of a sorted array. Therefore, the overall time complexity is $O(n + m)$.

**Algorithm: Merge-Sort-Based BST Merge**

**Require:** Two balanced binary search trees $T_1$ and $T_2$
**Ensure:** A single balanced binary search tree containing all elements
  from $T_1$ and $T_2$
1: **procedure** MERGESORTBASEDMERGE($T_1, T_2$)
2:   $A_1 \leftarrow$ INORDERTRAVERSAL($T_1$)
3:   $A_2 \leftarrow$ INORDERTRAVERSAL($T_2$)
4:   $A \leftarrow$ MERGESORTEDARRAYS($A_1, A_2$)
5:   $T \leftarrow$ BUILDBALANCEDBST($A$)
6:   **return** $T$

This algorithm performs well when $m/n$ is large, as the overall time complexity will be approximately $O(n)$. However, when $m$ is pretty negeligible compared to $n$, a full iteration over $T_1$ will be still needed and the overall time complexity will still be $O(n)$, which wastes a lot of time.

In fact, it Stockmayer and Yao has proven that in term of number of comparisions, this algorithm is optimal when $m \leq n \leq \lfloor 3m/2 \rfloor + 1$ (Stockmeyer & Yao, 1980). This algorithm, however, does not perform well outside this range.

## 2.5   Brown and Tarjan's merging algorithm (1979)

In 1979, Brown and Tarjan proposed another algorithm based on the two merging algorithm mentioned above. It utilized both the tree structure for fast insertion-place location and the ordered property to reduce redundant operations. The algorithm again chooses the $T_1$ as the base tree and view the merging process as $m$ insertions to a balanced BST of size $n$. However, the property that the inserted objects themselves are sorted helped to make the algorithm more efficient. Instead of iterating from the root, each insertion starts with the ending position of the last insertion, as it can be already told that the next insertion will happen to the right of the last insertion.

To be more specific, the algorithm keeps a stack called *path* and a stack called *successor*. The former is used to record the path from the root to the current node, while the latter records all the nodes on the *path* that is larger than the current node (that means they are on the right side of the current node their left children is visited on the *path*). Each insertion, instead of starting from the root, starts from the last node on the *successor* that is

smaller than the node to be inserted. Keep extending the *path* and *successor* during insertion. And the path shrinks back after the insertion, until a rebalance operation is triggered or we know that there need no rebalancing at all.

---

**Algorithm: Brown-Tarjan Fast Height-Balanced Tree Merge**

**Require:** Two balanced binary search trees $T_1$ (size $n$) and $T_2$ (size $m$), where $n \geq m$
**Ensure:** A single balanced BST containing all elements from $T_1$ and $T_2$

1: **procedure** FASTMERGE($T_1, T_2$)
2:     Initialize stack $path \leftarrow \{\text{root}(T_1)\}$
3:     Initialize empty stack $successor$
4:     $height \leftarrow \text{height}(T_1)$
5:     **for all** nodes $x$ in $T_2$ (in-order traversal) **do**
6:         Detach $x$ from $T_2$
            ▷ — *Step 1: Adjust path to maintain PathPredicate* —
7:         **while** $successor$ not empty **and** $key(x) > key(\text{top}(successor))$ **do**
8:             **repeat**
9:                 pop($path$)
10:            **until** top($path$) = top($successor$)
11:            pop($successor$)
        ▷ — *Step 2: Search down from last successor and insert x* —
12:        $p \leftarrow \text{top}(path)$
13:        **while** True **do**
14:            **if** $key(x) < key(p)$ **then**
15:                **if** $p.left = Nil$ **then**
16:                    $p.left \leftarrow x$; **break**
17:                **else**
18:                    push($p$, $successor$); $p \leftarrow p.left$
19:            **else**
20:                **if** $p.right = Nil$ **then**
21:                    $p.right \leftarrow x$; **break**
22:                **else**
23:                    $p \leftarrow p.right$

---

```
24:   │   │   └   push(p, path)
        ▷ — Step 3: Adjust balance factors and rebalance if needed —
25:   │   │   while path not empty do
26:   │   │   │   s ← pop(path)
27:   │   │   │   if tree at s is unbalanced then
28:   │   │   │   │   REBALANCE(s); break
29:   │   │   │   else
30:   │   │   │   └   Update balance factor of s
31:   │   │   │   if top(successor) = s then
32:   │   │   └   └   pop(successor)
33:   └   return root of T₁ (now merged)
```

It is worth noticing that the rebalance operation may make the initial path unusable. In this case, we can simply dispose of the path under the rotated node start next insertion there(Brown & Tarjan, 1979).

## 2.6   Optimality

When merging two instances of size $n$ and $m$ respectively, there are in total $\binom{n+m}{n}$ possible outcomes. According to the decision tree theory, each of them corresponds to a decision tree leaf node. Since the merging algorithm is comparison based, the decision tree has to be a binary tree (i.e. Each node has at most two children). The height of the decision tree is therefore no lower than $O(\log_2(\binom{n+m}{n}))$.

According to Sterling's approximation,

$$n! \approx \sqrt{2\pi n}(\frac{n}{e})^n \tag{1}$$

which means

$$O(\log(n!)) = O(\frac{1}{2}\log(2\pi n) + n\log n - n\log e) \quad = O(n\log n - n + O(logn)) \tag{2}$$

Using the definition of combination number,

$$\binom{n+m}{m} = \frac{(n+m)!}{n!m!} \tag{3}$$

$$\log\left(\binom{n+m}{n}\right) = \log\left(\frac{(n+m)!}{n!}\right) \tag{4}$$

$$= (n+m)\log(n+m) - n\log n - m\log m + O(\log(n+m)) \tag{5}$$

$$= n\log(1+\frac{m}{n}) + m\log(1+\frac{n}{m}) + O(\log(n+m)) \tag{6}$$

Since $m \leq n$ we have $n\log(1+m/n) \leq n \cdot \left(\frac{m}{n}\right) = m$, therefore the first term is $O(m)$. This means the first term should be neglected as $m\log(1+\frac{n}{m})$ is the dominant term compared to $O(m)$.

Since $m\log(1+\frac{n}{m})$ can bwe written as $m\log(n+m) - m\log(m)$, where the first term is more dominant than $O(\log(n+m))$, the third term should be neglected as well.

We can get the overall expression

$$\boxed{\log\left(\binom{n+m}{n}\right) = O(m\log(1+\frac{n}{m}))} \tag{7}$$

**Theorem**

The optimal time complexity of merging two instances of ordered data structures is $O(m\log(1+\frac{n}{m}))$, multiplied by the comparision cost with is assumed to be $O(1)$ in this case.

# 3 Hypothesis

# 4 Experiment Design

## 4.1 Variables

### 4.1.1 Control Variables

# Appendix

> **Test environment**
>
> - Device: Laptop
>
> - CPU: Intel(R) Core(TM) Ultra 7 155H (1.40 GHz)
>
> - Memory: 32GB
>
> - OS: Windows 11 24H2 26100.6584
>
> - Compiler (C++): g++ 15.2.0 (for algorithm efficiency test and timing)
>
> - Interpreter: Python 3.12.2 (for batch test)

Listing 1: AvlSet

```cpp
#include <vector>
#include <algorithm>
#include <memory>
#include <list>
#include <functional>
#include <stack>
#ifdef DEBUG
#include <assert.h>
#endif

template <typename T, typename Compare = std::less<T>>
class AVLSet {
private:
    struct Node {
```

```
15        T key;
16        Node* left;
17        Node* right;
18        int height;
19
20        template <typename... Args>
21        Node(Args&&... args)
22            : key(std::forward<Args>(args)...),
23              left(nullptr),
24              right(nullptr),
25              height(1) {}
26    };
27
28    Node* root;
29    size_t size;
30    Compare comp_;
31
32    // Helper functions for merging
33    int height(Node* node) const {
34        return node ? node->height : 0;
35    }
36
37    Node* create_node(const T& key) {
38        ++size;
39        return new Node(key);
40    }
41
42    Node* create_node(T&& key) {
43        ++size;
44        return new Node(std::move(key));
45    }
46
47    // Generates a balanced subtree in O(n) time out from a
            ordered sequence.
48    // Returns the root node pointer.
49    // *Preconditions
50    // keys have to be ordered
51    // _RandAccIt is the random access iterator
52    // (*bg) and (*ed) should be of type T
53
```

```cpp
template<typename _RandAccIt>
Node* build(const _RandAccIt& bg, const _RandAccIt& ed){
    if(bg == ed) return nullptr;
    auto it = bg;
    if(++it == ed) return create_node(*bg);
    it = bg + (ed - bg) / 2; // The same as (bg + ed)/2 but
        avoids overflow problems
    auto cur = create_node(*it);
    cur->left = build(bg, it);
    cur->right = build(it + 1, ed);
    update_height(cur);
    return cur;
}

void update_height(Node* node) {
    node->height = 1 + std::max(height(node->left), height(
        node->right));
}

void delete_tree(Node* node) {
    if (!node) return;
    delete_tree(node->left);
    delete_tree(node->right);
    --size;
    delete node;
}

Node* rotate_right(Node* y) {
    Node* x = y->left;
    Node* T2 = x->right;

    x->right = y;
    y->left = T2;

    update_height(y);
    update_height(x);

    return x;
}
```

```
92      Node* rotate_left(Node* x) {
93          Node* y = x->right;
94          Node* T2 = y->left;
95
96          y->left = x;
97          x->right = T2;
98
99          update_height(x);
100         update_height(y);
101
102         return y;
103     }
104
105     int balance_factor(Node* node) const {
106         return node ? height(node->left) - height(node->right) :
                    0;
107     }
108
109     Node* balance(Node* node) {
110         update_height(node);
111         int bf = balance_factor(node);
112
113         // Left Heavy
114         if (bf > 1) {
115             if (balance_factor(node->left) < 0)
116                 node->left = rotate_left(node->left);
117             return rotate_right(node);
118         }
119         // Right Heavy
120         if (bf < -1) {
121             if (balance_factor(node->right) > 0)
122                 node->right = rotate_right(node->right);
123             return rotate_left(node);
124         }
125         return node;
126     }
127
128     template <typename Func>
129     void traverse_in_order(Node* node, Func f) const {
130         if (!node) return;
```

```cpp
131            std::stack<Node*> stack;
132            Node* current = node;
133            while (current || !stack.empty()) {
134                while (current) {
135                    stack.push(current);
136                    current = current->left;
137                }
138                current = stack.top();
139                stack.pop();
140                f(current->key);
141                current = current->right;
142            }
143        }
144
145        Node* insert(Node* node, const T& key) {
146            if (!node) {
147                return create_node(key);
148            }
149
150            if (comp_(key, node->key)) {
151                node->left = insert(node->left, key);
152            } else if (comp_(node->key, key)) {
153                node->right = insert(node->right, key);
154            } else {
155                return node;
156            }
157
158            return balance(node);
159        }
160
161 public:
162     AVLSet(Compare Comp = Compare()) : root(nullptr), size(0),
163            comp_(Comp) {}
163     Compare comparator() const { return comp_; }
164     // Move operations
165     AVLSet(AVLSet&& other) noexcept
166            : root(other.root),
167              size(other.size){
168            other.root = nullptr;
169            other.size = 0;
```

13

```
170        }
171
172        AVLSet& operator=(AVLSet&& other) noexcept {
173            if (this != &other) {
174                clear();
175                root = other.root;
176                size = other.size;
177                other.root = nullptr;
178                other.size = 0;
179            }
180            return *this;
181        }
182
183        // Disable copy operations
184        AVLSet(const AVLSet&) = delete;
185        AVLSet& operator=(const AVLSet&) = delete;
186
187        void clear() {
188            delete_tree(root);
189            root = nullptr;
190        }
191
192        bool empty() const {
193            return size == 0;
194        }
195
196        size_t get_size() const {
197            return size;
198        }
199
200        template <typename Func>
201        void traverse_in_order(Func f) const {
202            traverse_in_order(root, f);
203        }
204
205        std::vector<T> items() const {
206            std::vector<T> result;
207            traverse_in_order([&result](const T& key) {
208                result.push_back(key);
209            });
```

```cpp
210        return result;
211    }
212
213
214    void swap_with(AVLSet& other) {
215        std::swap(root, other.root);
216        std::swap(size, other.size);
217    }
218
219    void insert(const T& val){
220        root = insert(root, val);
221    }
222    void remove(const T& val){
223        root = remove(root, val);
224    }
225    template<typename RandAccIt>
226    void construct(RandAccIt bg, RandAccIt ed){
227        // Clear the current tree
228        delete_tree(root);
229        root = nullptr;
230        // Build new tree
231        root = build(bg, ed);
232        size = static_cast<size_t>(ed - bg);
233    }
234
235
236    /** Merge two sets in O(N+M) time.
237     * Some additional space may be costed.
238     * But it does not affect the result of the experiment.
239     * @param other The AVLSet to be merged into this set.
240     */
241
242    void linearmerge(AVLSet&& other) {
243        if (other.empty()) return;
244
245        // Get elements from both trees
246        std::vector<T> q1 = items();
247        std::vector<T> q2 = other.items();
248        std::vector<T> all_elements;
249        all_elements.reserve(q1.size() + q2.size());
```

```
250
251         // Merge sorted vectors
252         std::merge(q1.begin(), q1.end(), q2.begin(), q2.end(),
253                    std::back_inserter(all_elements), comp_);
254
255         // Recycle both trees
256         delete_tree(root);
257         delete_tree(other.root);
258         root = nullptr;
259         other.root = nullptr;
260
261         // Build new tree
262         root = build(all_elements.begin(), all_elements.end());
263         size = all_elements.size();
264     }
265
266     /**
267      * Merge two sets in O(N log(M)) time.
268      * @param other The AVLSet to be merged into this set.
269      */
270     void simplemerge(AVLSet&& other) {
271         if (other.empty()) return;
272
273         if (size < other.size) { swap_with(other); }
274
275         // Insert all elements from the smaller tree (now 'other')
                 into this
276         other.traverse_in_order([this](const T& key) {
277             this->insert(key);
278         });
279
280         other.clear();
281     }
282
283     void brownmerge(AVLSet&& other) {
284         if (other.empty()) return;
285         if (size < other.size) swap_with(other);
286
287         std::vector<T> elems = other.items();
288         other.clear();
```

```
289
290         // stacks of pointers-to-links (Node**). Each points to
                some parent->left or parent->right or &root
291         std::vector<Node**> path;
292         std::vector<Node**> successor;
293
294         path.push_back(&root);
295
296         for (const T& x : elems) {
297             // ==== CLIMB/RETRACT: pop successors while their key
                    <= x (i.e. not (x < succ_key)) ====
298             while (!successor.empty() && !comp_(x, (*successor.
                    back())->key)) {
299                 Node** succLink = successor.back();
300                 // pop path entries until top == succLink
301                 while (!path.empty() && path.back() != succLink)
                        path.pop_back();
302                 successor.pop_back();
303             }
304
305             // ==== DESCEND from current finger (path.back()) to
                    insertion point ====
306             Node** curLink = path.empty() ? &root : path.back();
307             Node* p = *curLink;
308             if (!p) {
309                 // empty subtree (rare because root existed),
                        insert directly
310                 *curLink = create_node(x);
311                 path.push_back(curLink);
312             } else {
313                 for (;;) {
314                     path.push_back(curLink); // link that points to
                            p
315                     if (comp_(x, p->key)) {
316                         // go left; mark this link as a successor (
                                we turned left here)
317                         if (p->left == nullptr) {
318                             p->left = create_node(x);
319                             path.push_back(&(p->left));
320                             break;
```

17

```
321                    } else {
322                        successor.push_back(curLink);
323                        curLink = &(p->left);
324                        p = *curLink;
325                    }
326                } else {
327                    // go right
328                    if (p->right == nullptr) {
329                        p->right = create_node(x);
330                        path.push_back(&(p->right));
331                        break;
332                    } else {
333                        curLink = &(p->right);
334                        p = *curLink;
335                    }
336                }
337            }
338        }
339
340        // ==== REBALANCE upward using the link-stack; attach
                 rotated subtree via *link ====
341        while (!path.empty()) {
342            Node** link = path.back();
343            Node* s = *link;
344            path.pop_back();
345
346            if (!successor.empty() && successor.back() == link
                 ) successor.pop_back();
347
348            update_height(s);
349            int bf = balance_factor(s);
350
351            if (std::abs(bf) > 1) {
352                Node* newsub = balance(s); // returns new root
                     of this subtree
353                *link = newsub; // reattach correctly via the
                     link
354
355                // retract path until the remaining entries are
                     consistent with this rotation
```

18

```
356              while (!path.empty() && path.back() != link)
                     path.pop_back();
357              break; // stop climbing after performing
                     rotation
358          }
359
360          if (bf == 0) {
361              // height didn't increase => stop climbing
362              break;
363          }
364          // else continue climbing
365      }
366  }
367  #ifdef DEBUG
368  auto items_after = items();
369  for (size_t i = 1; i < items_after.size(); ++i)
370      assert(!comp_(items_after[i], items_after[i-1]));
371  #endif
372  }
373
374 };
```

# References

Brown, M. R., & Tarjan, R. E. (1979). A fast merging algorithm. *Journal of the ACM (JACM)*, *26*(2), 211–226.

Cormen, T. H., Leiserson, C. E., Rivest, R. L., & Stein, C. (2022). *Introduction to algorithms* (Fourth ed.). MIT Press.

Karlton, P. L., Fuller, S. H., Scroggs, R., & Kaehler, E. (1976). Performance of height-balanced trees. *Communications of the ACM*, *19*(1), 23–28.

Stockmeyer, P. K., & Yao, F. F. (1980). On the optimality of linear merge. *SIAM Journal on Computing*, *9*(1), 85–90.