

Comparing different data structures in merging efficiency.

Changhui (Eric) Zhou

May 28, 2025

word count: ???

Contents

1	Introduction	2
2	Theory	2
2.1	Data structure terminology	2
2.2	Optimality	3
	References	15

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 pursuit 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 affecting 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 investigating the theoretical 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 ordered data structures?

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

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

2.2 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,

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

which means

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

Using the definition of combination number,

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

which is approximately $O(n \log_2(\frac{n}{m}))$ (provided $n \leq m$).

Appendix

Listing 1: VectorSet

```

1  #include <vector>
2  #include <algorithm>
3
4  template <typename T>
5  class VectorSet {
6  private:
7      std::vector<T> elements;
8  
```

```

9 public:
10     VectorSet() = default;
11     VectorSet(const std::vector<T>& vec, bool sorted = 0){
12         elements = vec;
13         if(!sorted)
14             std::sort(elements.begin(), elements.end());
15     }
16
17     // Inserts an element into the set if not already present
18     void insert(const T& value) {
19         auto it = std::lower_bound(elements.begin(), elements.end
20             (), value);
21         if (it == elements.end() || *it != value) {
22             elements.insert(it, value);
23         }
24     }
25
26     // Removes an element from the set if present
27     void erase(const T& value) {
28         auto it = std::lower_bound(elements.begin(), elements.end
29             (), value);
30         if (it != elements.end() && *it == value) {
31             elements.erase(it);
32         }
33     }
34
35     // Checks if an element exists in the set
36     bool contains(const T& value) const {
37         return std::binary_search(elements.begin(), elements.end()
38             , value);
39     }
40
41     // Returns the number of elements in the set
42     size_t size() const {
43         return elements.size();
44     }
45
46     // Checks if the set is empty
47     bool empty() const {
48         return elements.empty();
49     }

```

```

46     }
47
48     // Allows iteration over the elements (read-only)
49     typename std::vector<T>::const_iterator begin() const {
50         return elements.begin();
51     }
52
53     typename std::vector<T>::const_iterator end() const {
54         return elements.end();
55     }
56
57     // Clear all elements from the set
58     void clear() {
59         elements.clear();
60     }
61
62     // Merge another VectorSet into this one (union operation)
63     void merge(const VectorSet<T>& other) {
64         std::vector<T> merged;
65         merged.reserve(elements.size() + other.elements.size());
66
67         auto it1 = elements.begin(), end1 = elements.end();
68         auto it2 = other.elements.begin(), end2 = other.elements.
69             end();
70
71         while (it1 != end1 && it2 != end2) {
72             if (*it1 < *it2) {
73                 merged.push_back(*it1);
74                 ++it1;
75             } else if (*it2 < *it1) {
76                 merged.push_back(*it2);
77                 ++it2;
78             } else { // Equal elements
79                 merged.push_back(*it1);
80                 ++it1;
81                 ++it2;
82             }
83         }
84
85         // Add remaining elements from either vector

```

```

85     merged.insert(merged.end(), it1, end1);
86     merged.insert(merged.end(), it2, end2);
87
88     elements = std::move(merged);
89 }
90
91 // Merge-and-assign operator
92 VectorSet<T>& operator+=(const VectorSet<T>& other) {
93     merge(other);
94     return *this;
95 }
96 };

```

Listing 2: AVLSet

```

1  #include <vector>
2  #include <algorithm>
3  #include <memory>
4  #include <list>
5  #include <functional>
6  #include <stack>
7
8  template <typename T, typename Compare = std::less<T>>
9  class AVLSet {
10 private:
11     struct Node {
12         T key;
13         Node* left;
14         Node* right;
15         int height;
16
17         template <typename... Args>
18         Node(Args&&... args)
19             : key(std::forward<Args>(args)...),
20               left(nullptr),
21               right(nullptr),
22               height(1) {}
23     };
24
25     Node* root;
26     size_t size;

```

```

27 Compare comp;
28
29 // Memory pool management
30 std::list<Node> node_storage;
31 std::vector<Node*> free_nodes;
32
33 // Helper functions for merging
34 int height(Node* node) const {
35     return node ? node->height : 0;
36 }
37
38 // Memory pool operations
39 Node* create_node(const T& key) {
40     if (!free_nodes.empty()) {
41         Node* node = free_nodes.back();
42         free_nodes.pop_back();
43         *node = Node(key);
44         return node;
45     }
46     node_storage.emplace_back(key);
47     return &node_storage.back();
48 }
49
50 Node* create_node(T&& key) {
51     if (!free_nodes.empty()) {
52         Node* node = free_nodes.back();
53         free_nodes.pop_back();
54         *node = Node(std::move(key));
55         return node;
56     }
57     node_storage.emplace_back(std::move(key));
58     return &node_storage.back();
59 }
60
61 // Generates a balanced subtree in O(n) time out from a
62 // ordered sequence.
63 // Returns the root node pointer.
64 // *Preconditions
65 // keys have to be ordered
66 // _RandAccIt is the random access iterator

```



```

66 // (*bg) and (*ed) should be of type T
67
68 template<typename _RandAccIt>
69 Node* create_node(const _RandAccIt& bg, const _RandAccIt& ed)
70 {
71     if(bg == ed) return nullptr;
72     auto it = bg;
73     if(++it == ed) return create_node(*bg);
74     it = bg + (ed - bg) / 2; // The same as (bg + ed)/2 but
75     // avoids overflow problems
76     auto cur = create_node(*it);
77     cur->left = create_node(bg, it);
78     cur->right = create_node(++it, ed);
79     update_height(cur);
80     return cur;
81 }
82
83 void recycle_node(Node* node) {
84     free_nodes.push_back(node);
85 }
86
87 void update_height(Node* node) {
88     node->height = 1 + std::max(height(node->left), height(
89         node->right));
90 }
91
92 Node* rotate_right(Node* y) {
93     Node* x = y->left;
94     Node* T2 = x->right;
95
96     x->right = y;
97     y->left = T2;
98
99     update_height(y);
100     update_height(x);
101
102     return x;
103 }
104
105 Node* rotate_left(Node* x) {

```

```

103     Node* y = x->right;
104     Node* T2 = y->left;
105
106     y->left = x;
107     x->right = T2;
108
109     update_height(x);
110     update_height(y);
111
112     return y;
113 }
114
115 int balance_factor(Node* node) const {
116     return node ? height(node->left) - height(node->right) :
117         0;
118 }
119
120 Node* balance(Node* node) {
121     update_height(node);
122     int bf = balance_factor(node);
123
124     // Left Heavy
125     if (bf > 1) {
126         if (balance_factor(node->left) < 0)
127             node->left = rotate_left(node->left);
128         return rotate_right(node);
129     }
130     // Right Heavy
131     if (bf < -1) {
132         if (balance_factor(node->right) > 0)
133             node->right = rotate_right(node->right);
134         return rotate_left(node);
135     }
136     return node;
137 }
138
139 template <typename Func>
140 void traverse_in_order(Node* node, Func f) const {
141     if (!node) return;
142     traverse_in_order(node->left, f);

```

```

142     f(node->key);
143     traverse_in_order(node->right, f);
144 }
145
146 void insert_merge(const T& key) {
147     if (!root) {
148         root = create_node(key);
149         size++;
150         return;
151     }
152
153     std::vector<Node*> path;
154     std::vector<Node*> successor;
155     Node* current = root;
156     bool inserted = false;
157
158     // Climb up to find the insertion path
159     while (true) {
160         path.push_back(current);
161         if (comp(key, current->key)) {
162             if (!current->left) break;
163             successor.push_back(current);
164             current = current->left;
165         } else if (comp(current->key, key)) {
166             if (!current->right) break;
167             current = current->right;
168         } else {
169             // Duplicate, do not insert
170             return;
171         }
172     }
173
174     // Insert the new node
175     Node* newNode = create_node(key);
176     if (comp(key, current->key)) {
177         current->left = newNode;
178     } else {
179         current->right = newNode;
180     }
181     size++;

```

```

182
183     // Retrace the path to update heights and balance
184     while (!path.empty()) {
185         Node* node = path.back();
186         path.pop_back();
187         node = balance(node);
188
189         if (!path.empty()) {
190             if (path.back()->left == node) {
191                 path.back()->left = node;
192             } else {
193                 path.back()->right = node;
194             }
195         } else {
196             root = node;
197         }
198     }
199 }
200
201 public:
202     AVLSet() : root(nullptr), size(0), comp(Compare()) {}
203
204     // Move operations
205     AVLSet(AVLSet&& other) noexcept
206         : root(other.root),
207           size(other.size),
208           node_storage(std::move(other.node_storage)),
209           free_nodes(std::move(other.free_nodes)) {
210         other.root = nullptr;
211         other.size = 0;
212     }
213
214     AVLSet& operator=(AVLSet&& other) noexcept {
215         if (this != &other) {
216             clear();
217             root = other.root;
218             size = other.size;
219             node_storage = std::move(other.node_storage);
220             free_nodes = std::move(other.free_nodes);
221             other.root = nullptr;

```

```

222         other.size = 0;
223     }
224     return *this;
225 }
226
227 // Disable copy operations
228 AVLSet(const AVLSet&) = delete;
229 AVLSet& operator=(const AVLSet&) = delete;
230
231 void clear() {
232     node_storage.clear();
233     free_nodes.clear();
234     root = nullptr;
235     size = 0;
236 }
237
238 bool empty() const {
239     return size == 0;
240 }
241
242 size_t get_size() const {
243     return size;
244 }
245
246 template <typename Func>
247 void traverse_in_order(Func f) const {
248     traverse_in_order(root, f);
249 }
250
251 std::vector<T>* to_vector(Node *node = nullptr){
252     if(node == nullptr)
253         return nullptr;
254     std::vector<T> *lson = to_vector(node->left), *rson =
        to_vector(node->right);
255     if(lson == nullptr)
256         lson = new std::vector<T>;
257     lson->push_back(node->key);
258     if(rson != nullptr)
259         for(T it: *rson)
260             lson->push_back(it);

```

```

261     return lson;
262 }
263
264 void merge(AVLSet&& other) {
265     if (other.empty()) return;
266
267     if (size < other.size) {
268         // Swap to merge smaller into larger
269         std::swap(root, other.root);
270         std::swap(size, other.size);
271         std::swap(node_storage, other.node_storage);
272         std::swap(free_nodes, other.free_nodes);
273     }
274
275     // Insert all elements from the smaller tree (now 'other')
276     // into this
277     other.traverse_in_order([this](const T& key) {
278         this->insert_merge(key);
279     });
280
281     other.clear();
282 }
283
284 void linearmerge(AVLSet&& other){
285 }
286
287 private:
288 // Example of modified insert implementation
289 Node* insert(Node* node, const T& key) {
290     if (!node) {
291         size++;
292         return create_node(key);
293     }
294
295     if (comp(key, node->key)) {
296         node->left = insert(node->left, key);
297     } else if (comp(node->key, key)) {
298         node->right = insert(node->right, key);
299     } else {

```

```

300         return node;
301     }
302
303     return balance(node);
304 }
305
306 // Example of modified remove implementation
307 Node* remove(Node* node, const T& key) {
308     if (!node) return nullptr;
309
310     if (comp(key, node->key)) {
311         node->left = remove(node->left, key);
312     } else if (comp(node->key, key)) {
313         node->right = remove(node->right, key);
314     } else {
315         // Node deletion with recycling
316         if (!node->left || !node->right) {
317             Node* temp = node->left ? node->left : node->right
318                 ;
319             recycle_node(node);
320             size--;
321             node = temp;
322         } else {
323             Node* temp = findMin(node->right);
324             node->key = std::move(temp->key);
325             node->right = remove(node->right, temp->key);
326         }
327     }
328
329     return node ? balance(node) : nullptr;
330 }
331
332 public:
333 void insert(const T& val){
334     insert(root, val);
335 }
336 void remove(const T& val){
337     remove(root, val);
338 }

```

```
339     friend AVLSet<T>* AVLSet_from_ordered(std::vector<T> data){};
340 };
341
342 template <typename T>
343 AVLSet<T>* AVLSet_from_ordered(std::vector<T> data){
344     AVLSet<T>* ret = new AVLSet<T>;
345     //ret->root = ret->create_node()
346 }
```

References

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