

Random Algorithms (RA-MIRI)

Assignment 3

Daniel Benedí

December 2022

During this assignment we are going to carry out an emperical study of the performance of skip lists. Skip lists are a probablistic data structe proposed as an alternative to balanced trees. [2]

1 Implementation

In this section I will explain my implementation with pseudo-code, but the whole implementation using C++ can be found in the appendix A.

A skip list can be seen as a linked list with different layers where the height of each element depends on a geometric distribution. An example of this can be found in figure 1; first of all, we observe that although the input is not ordered, the elements in the skip list will end up ordered. Also, we see that the higher we go in the skip list, the less amount number of elements there are. This height is a random variable with a geometric distribution.

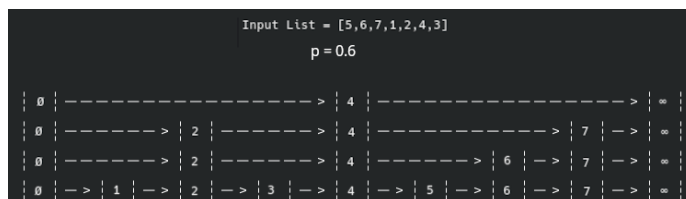


Figure 1: An example of skip list with 7 elements and $p = 0.6$.

First of all, we need to define what a node is (data structure 1). It will contain the data to be stored in the set (its data type is represented by \mathbf{T}). It also has two lists to represent the elements before and after in each level, this means that $|\text{forward}| = |\text{backward}| = \text{height}$. The list **backward** wouldn't be needed if the operation **remove** was necessary. Also, we need to have a list, called **head**, to represent the first element of each level.

Data Structure 1 Node of skiplist

```
struct Node
    list forward, backward
    T data
endStruct
```

Once the data structure is defined, we can start describing the operations that are going to be implemented: **insert**, **remove** and **contains**. Firstly, we will implement an auxiliary operation called **search**. This operation will return the node iff it is in the set or the node before in case it wasn't in the set.

Algorithm 1 Search for a node or the place where it should be

```
1: procedure SEARCH( $A$ )
2:    $h \leftarrow |\text{head}| - 1$ 
3:   while  $h > 0 \wedge \text{head}[h].\text{data} > A$  do
4:      $h \leftarrow h - 1$ 
5:   end while
6:   if  $h < 0$  then
7:     return null
8:   end if
9:    $r \leftarrow \text{head}[h]$ 
10:  while  $h > 0$  do
11:    if  $r.\text{forward}[h] \wedge r.\text{forward}[h].\text{data} \leq A$  then
12:       $r \leftarrow r.\text{forward}[h]$ 
13:    else
14:       $h \leftarrow h - 1$ 
15:    end if
16:  end while
17:  return  $r$ 
18: end procedure
```

In this algorithm, first of all we check if the searched item is bigger than any node in the head (lines 2 to 5), otherwise it means that all the items are bigger than the searched item and therefore it isn't in the set and there is no node lower than it (lines 6 to 8). In the rest of the pseudocode (lines 10 to 16) it will search if in that level there following node exists and if it does also checks whether it is lower or equal than the searched item, in that case the next result node will be the node stored in the forward list at that level. Otherwise, we decrease the level, because we cannot look for anymore in that level.

With this implementation in mind we can easily define the **contains** operation, that would be: This function is trivial because it calls the **search** procedure and it will return **true** if and only if the returned node is not **null** and it is equal to the searched element.

Algorithm 2 Returns whether the set contains a element

```

1: procedure CONTAINS( $A$ )
2:    $r \leftarrow \text{SEARCH}(A)$ 
3:   return  $r \neq \text{null} \wedge r.\text{data} = A$ 
4: end procedure

```

The **remove** operation is also quiet trivial as well, it will call the **search** procedure and if the node contains the searched item, it will assign to the nodes in the backward list the nodes in the forward list and the nodes in the forward list the nodes in the backward list.

Algorithm 3 Removes the element from the set

```

1: procedure CONTAINS( $A$ )
2:    $r \leftarrow \text{SEARCH}(A)$ 
3:   if  $r \neq \text{null} \wedge r.\text{data} = A$  then
4:      $h \leftarrow |r.\text{forward}|$ 
5:     for  $l = h - 1$  to  $0$  do
6:       if  $r.\text{backward}[l] \neq \text{null}$  then
7:          $r.\text{backward}[l].\text{forward}[l] \leftarrow r.\text{forward}[l]$ 
8:       if  $r.\text{forward}[l] \neq \text{null}$  then
9:          $r.\text{forward}[l].\text{backward}[l] \leftarrow r.\text{backward}[l]$ 
10:      end if
11:    else ▷ This means it is in the head
12:       $\text{head}[l] \leftarrow r.\text{forward}[l]$ 
13:      if  $\text{head}[l] = \text{null}$  then
14:        remove top of head
15:      end if
16:    end if
17:  end for
18: end if
19: end procedure

```

Although the **insert** operation is really similar to the **search** operation or the remove one, we cannot use it because we need to insert the new node in the levels as we are decreasing them. This is the reason because we haven't reused it. For the simplicity of the pseudo-code we are going to assume that the element is not in the set already, but in the real pseudo-code we take it into consideration.

As we explained, first of all we insert the element to the new levels of the head if there are new levels. Then we update the levels of the head that are bigger to the new element and finally we start going down the skip list and we need to decrease a level we insert the node in that point.

We have opted to implement our own geometric distribution generator. We generate u.a.r a number between 0 and 1 (both included) and will increase a

Algorithm 4 Insertion in skip list

Pre: The set doesn't contain the element A

Post: The set contains the element A

```
1: procedure INSERT( $A$ )
2:    $h \leftarrow \text{Geom}(p)$ 
3:    $\text{node} \leftarrow \text{Node}(A, h)$ 
4:   if  $|\text{head}| < h$  then
5:     insert the node to the head until it has the size of  $h$ 
6:      $h \leftarrow |\text{head}| - |\{x \in \text{head} \mid x = A\}|$ 
7:   end if
8:    $h \leftarrow h - 1$ 
9:   if  $h < 0$  then
10:    return
11:  end if
12:   $r \leftarrow \text{head}[h]$ 
13:  while  $h \geq 0 \wedge r.\text{data} > A$  do
14:    Update head element at height  $h$ 
15:     $h \leftarrow h - 1$ 
16:     $r \leftarrow \text{head}[h]$ 
17:  end while
18:  while  $h \geq 0$  do
19:    while  $r.\text{data}[h] \wedge r.\text{data}[h] < A$  do
20:       $r \leftarrow r.\text{forward}[h]$ 
21:    end while
22:     $\text{node}.\text{backward}[h] = r$ 
23:    if  $r.\text{forward}[h]$  then
24:       $\text{node}.\text{forward}[h] = r.\text{forward}[h]$ 
25:       $r.\text{forward}[h].\text{backward}[h] = \text{node}$ 
26:    end if
27:     $\text{node}.\text{forward}[h] = \text{node}$ 
28:     $h \leftarrow h - 1$ 
29:  end while
30: end procedure
```

counter while the generated number is bigger than q . We are aware that in the `random` library there is a geometric distribution generator of integers (adding 1 to the generated number because it starts at 0 and we need to them to start at 1), but we prefer to implement our own version without any technical reason.

2 Experimentation

In order to to the experimentation we have tried different sizes of the amount of input elements, n , all of them with different probability, q , in the Geometric distribution in the construction. In order to get statistically significant result, we have run several times, M , each case shuffling the input every time.

The values for n have been taken from the assignment's statement where it is proposed to do the experiment from $n = 2000$ to $n = 20000$ in steps of 100, so we will try $n = 2000, 2100, 2200, 2300, \dots, 19700, 19800, 19900$ and 20000 . For every case we tried all possible $q = 1 - p$ with p taking values from 0.1 to 0.9, avoiding the extreme probabilities because they would have a strange behaviour; also taken from the statement. Every case of n and q have been run 100 times so the results are significant and they would have low deviation from the theoretical values.

2.1 Results

Once we have done the experiments, we can plot the experimental total search cost to compare with the theoretical one. In figures 2 and 3 we observe different means of the experiment total search cost compared with the theory. We observe that they are really similar to the theoretical cost. To obtain this results, the total search cost of each experiment has been gathered and then plotted the mean of the 100 experiments for each n and q . The theoretical cost $C_{n,q}$ has been explained in the lectures and in the assignment's statement and it is:

$$C_{n,q} = \frac{1}{q} n \log_{1/q} n + n \left(\frac{1/q}{\ln 1/q} (\gamma - 1) + \frac{1}{\ln 1/q} - \frac{1/q}{2} \right), \gamma = 0.5772156649$$

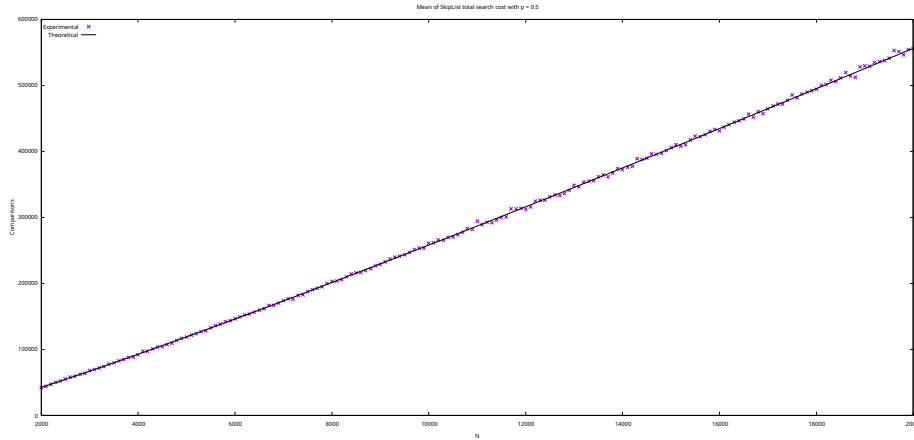


Figure 2: Theoretical and experimental total search cost with $q = 0.5$

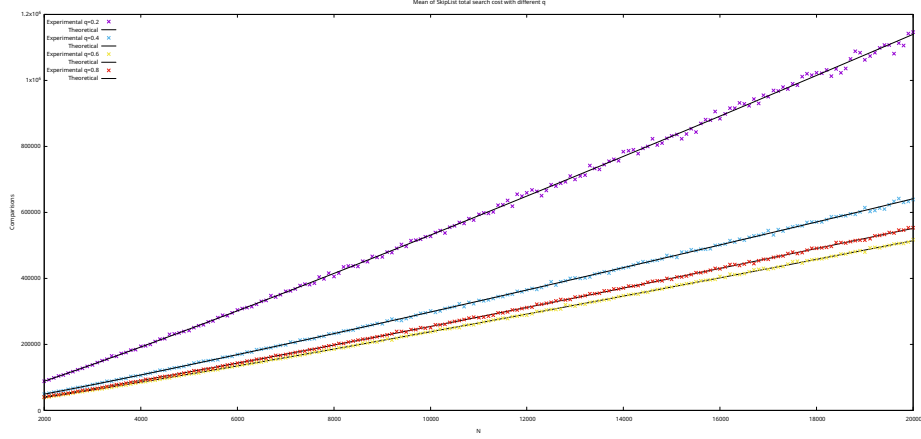


Figure 3: Theoretical and experimental total search cost with $q = 0.2, 0.4, 0.6$ and 0.8

Apparently, the implementation seems to fit the theoretical cost. We propose the R^2 metric, coefficient of determination, to measure how well it fits the theoretical cost. We will assume that our observed data y is the mean cost of our experiments for each input size (we will fix the q parameter) and \hat{y} will be the theoretical values. This metric consist on the proportion of the variation in the dependent variable that is predictable from the independent variable:

$$R^2 = 1 - \frac{\sum_i (y_i - \hat{y}_i)^2}{\sum_i (y_i - \frac{1}{n} \sum_j y_j)^2}$$

It is used to know how well it approximates to the theoretical value compared with the worst possible least-squares predictor (the mean). The higher is the metric the better it approximates. In table 1, we observe that all the values are really close to 1 which means that our experiments seems to corroborate the theoretical expectation of the total search cost.

q	0.1	0.2	0.3	0.4	0.5	0.6	0.7	0.8	0.9
R^2	0.99676	0.99761	0.99837	0.99926	0.99968	0.99983	0.99960	0.99911	0.99774

Table 1: R^2 of the experiments with respect to the theoretical values

Similarly, to what we have done with the expected total search cost, we can do an analysis of its variance. Peter Kirschenhofer and Conrado Martinez did a detailed analysis of the total search cost (succeeded and unsucceded) of the skiplist including its expectation and **variance** [1]. In their article, they give a complex formula for the variance, but also some factors K to approximate the variance assintotically with:

$$Var(C_{n,q}) = K * n^2$$

These values can be found in their original work [1]. In our experiments we have obtained results that are really close to this theoretical analysis. In figures 4 and 5 we observe that effectively they follow this theoretical analysis as expected, but it is important to notice that as n grows, the standard deviation is more dispersed around the expected value, this probably indicates that we should have done more experiments at least as n was getting bigger.

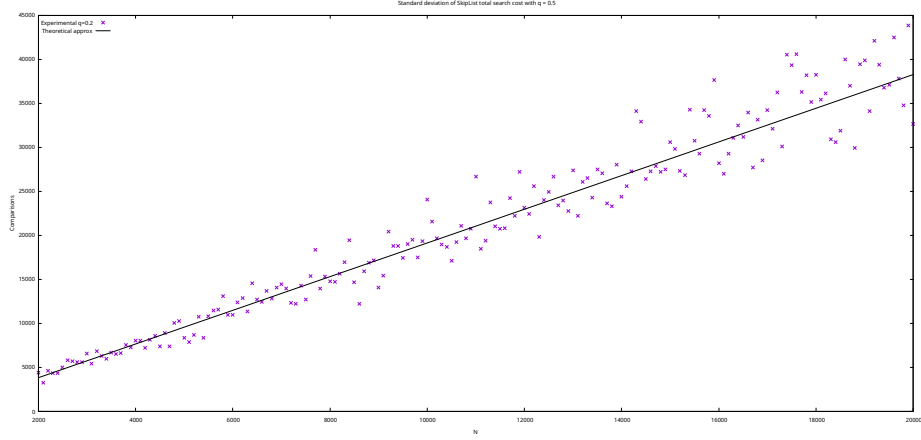


Figure 4: Theoretical and experimental total search cost standard deviation with $q = 0.5$

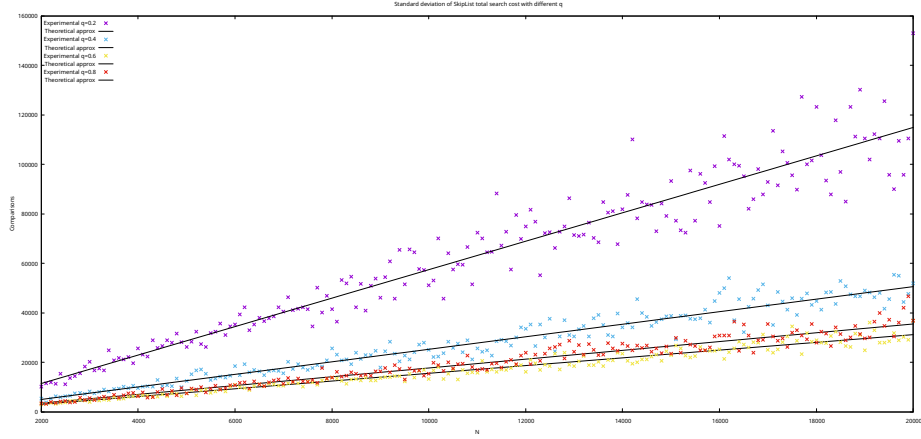


Figure 5: Theoretical and experimental total search cost with $q = 0.2, 0.4, 0.6$ and 0.8

When analyzing data structures not only their efficiency should be taken into account, but also their memory usage, because usually it is a trade off. From

the assignment's statements we have that:

$$S_{n,q} = \frac{n}{1-q} + \log_{1/q}(n) = E[\text{sum of the height of all nodes}] + E[\text{size of head}]$$

Note that the size of the head will be equal to the maximum of the heights of all nodes. We could do a theoretical analysis of the variance of the memory usage, taking into account that the nodes' heights are independent and identically distributed geometric variables, X , and we will define the size of the head as Y . Although we know that $Var(X) = n^2 \frac{1-p}{p^2}$ and could use an approximation to find the $Var(Y)$ using an article by Szpankowski and Rego [3], we are not going to do this analysis because it is beyond the analysis of this work.

In figure 6, we observe in the first two rows the mean memory use of the skip list with respect to the parameter q and, as we expected, the mean memory usage follows the theoretical analysis. We know that the standard deviation should be lineal because $Var(X) = O(n^2)$ and $Var(Y)$ seems to be $O(n^2)$ as well according to the previous article [3], and of course our standard deviation is lineal in the experiments as we expected, but as before it is quite disperse probably because we should have had to increase the number of experiments as n was growing.

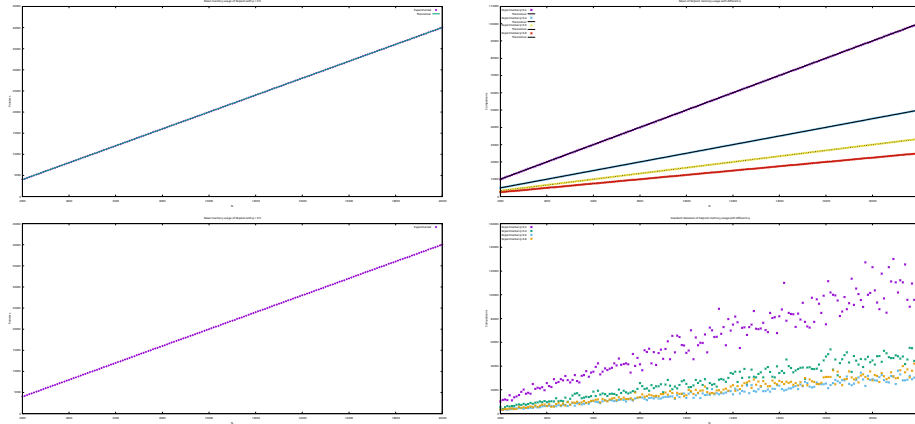


Figure 6: Mean and standard deviation of memory usage of Skip List.

3 Conclusion

Overall, we can conclude that the Skip List performs really nice with respect their average case and with a good performance. The total search cost seems to be linear in the number of items which is a really good performance for a data structure. The implementation was really simple, compared with self-balancing binary search tree. It is remarkable that as the number of elements in the skip list increased, the memory usage and the performance got more disperse around the theoretical values and therefore, one could think that when using it with a big amount of elements its good performance might not be guaranteed for single use cases.

Appendices

A C++ Implementation

```
1  //*****
2  // File: Set.h
3  // Author: Daniel Benedi Garcia
4  // Date: 8th December 2022
5  // Coms: Set class that implements a skip list
6  //          and a lot of associated utilities
7  //*****
8  #pragma once
9  #include <vector>
10
11 template <typename T>
12 class Set {
13 public:
14     /**
15      * Constructor.
16      *
17      * Creates an empty Set, sets the skip list
18      parameter q to the
19      * given value
20      * @param q double. Represents the probability
21      * of an element
22      * not having
23      */
24     Set(double q) : q(q) {}
25
26     /**
27      * Destructor.
28      *
29      * Clear all the extra structures used.
30      */
31     ~Set(){
32         Node* elem = head[0];
33         while(elem != nullptr){
34             Node* next = elem->forward[0];
35             delete elem;
36             elem = next;
```

```

36     }
37 }
38
39 /**
40  * Insert element x in the set.
41  *
42  * @param x    Type depends on template. The
43  *              element will be added
44  *              to the Set.
45  */
46 void insert(const T& x){
47     if(contains(x)) return;
48
49     int height = randomHeight();
50     Node* new_elem = new Node(x, height);
51
52     // If it is the first element in that height,
53     // add it to the head
54     while(head.size() < new_elem->height){
55         head.push_back(new_elem);
56         height--;
57     }
58     height--;
59
60     // It might happen that it is the first one
61     // added
62     if(height < 0) return;
63
64     Node* tmp = head[height];
65     while(height >= 0 && tmp->data > x){
66         new_elem->forward[height] = tmp;
67         tmp->backward[height] = new_elem;
68         head[height] = new_elem;
69         height--;
70         if(height >= 0)
71             tmp = head[height];
72     }
73     while(height >= 0){
74         // While for that level there is a smaller
75         // element
76         // and exist and following node, go forward
77         while(tmp->forward[height] && tmp->forward[
78             height]->data < x)
79             tmp = tmp->forward[height];
80
81         new_elem->backward[height] = tmp;

```

```

77         if(tmp->forward[height]) {
78             new_elem->forward[height] = tmp->forward[
                height];
79             tmp->forward[height]->backward[height] =
                new_elem;
80         }
81         tmp->forward[height] = new_elem;
82         height--;
83     }
84 }
85
86 /**
87  * Removes the element x from the set if present.
88  *
89  * @param x    Type depends on template. The
                element will be removed
90  *             from the set.
91  */
92 void remove(const T& x){
93     Node* elem = search(x);
94
95     if(elem && elem->data == x){
96         for(int lvl = elem->height-1; lvl >= 0; lvl--)
97             {
98                 if(elem->backward[lvl]){
99                     elem->backward[lvl]->forward[lvl] = elem->
100                        forward[lvl];
101                     if(elem->forward[lvl]){
102                         elem->forward[lvl]->backward[lvl] = elem
103                            ->backward[lvl];
104                     }
105                 }else{
106                     //If there is no previous element, that
                        means it is in the head
107                     head[lvl] = elem->forward[lvl];
108                     if(!head[lvl]){
109                         while(lvl + 1 > head.size())
110                             head.pop_back();
111                     }
112                 }
113             }
114         delete elem;
115     }
116 }
117
118 /**

```

```

116      * It will return wheter the element is in the set
117      *
118      * @param x    Type depends on template. The
119                  element will be searched
120                  in the set.
121      * @return    It will return true if the element
122                  is in the set, false
123                  otherwise.
124      */
125      bool contains(const T& x) const{
126          Node* res = search(x);
127          return res && res->data == x;
128      }
129
130      /**
131       * Returns the total search cost.
132       *
133       * @return    It will return the toal search cost.
134       */
135      int total_search_cost() const{
136          int cost = 0;
137          Node* next = head[0];
138          while(next != nullptr){
139              int height = head.size() - 1;
140              while(height > 0 && head[height]->data > next
141                  ->data){
142                  cost++;
143                  height--;
144              }
145
146              Node* res = head[height];
147              while(height >= 0){
148                  cost++;
149                  if(res->forward[height] && res->forward[
150                      height]->data <= next->data){
151                      res = res->forward[height];
152                  }else{
153                      height--;
154                  }
155              }
156
157              next = next->forward[0];
158          }
159
160          return cost;

```

```

157     }
158
159     /**
160      * Returns the number of pointers used, that is
161      * the sum of the height
162      * of all nodes in the skip list.
163      *
164      * @return Total number of pointers.
165      */
166     int number_pointers() const {
167         Node* next = head[0];
168         int size = head.size();
169         while(next != nullptr){
170             size += next->height;
171             next = next->forward[0];
172         }
173
174         return size;
175     }
176
177     double get_q() const {
178         return q;
179     }
180
181     int nodes_lvl(int lvl) const {
182         if(lvl >= head.size()) return 0;
183
184         int nodes = 0;
185         Node* node = head[lvl];
186         while(node != nullptr){
187             nodes++;
188             node = node->forward[lvl];
189         }
190
191         return nodes;
192     }
193
194     friend std::ostream& operator<< (std::ostream&
195         stream, const Set& set) {
196         // Print HEAD
197         for(int i = 0; i < set.head.size(); i++)
198             stream << "——"; stream << std::endl;
199         for(int i = 0; i < set.head.size(); i++)
200             stream << "  "; stream << std::endl;
201         for(int i = 0; i < set.head.size(); i++)
202             stream << "——"; stream << std::endl;

```

```

201
202 // Print BODY
203 Node* node = set.head[0];
204 while(node != nullptr){
205     for(int i = 0; i < set.head.size(); i++)
206         stream << "└┘"; stream << std::endl;
207     for(int i = 0; i < node->height; i++)
208         stream << "└v┘";
209     for(int i = node->height; i < set.head.size();
210         i++)
211         stream << "└┘"; stream << std::endl;
212     for(int i = 0; i < node->height; i++) stream
213         << "——┘";
214     for(int i = node->height; i < set.head.size();
215         i++)
216         stream << "└┘"; stream << std::endl;
217     for(int i = 0; i < node->height; i++)
218         stream << "——┘";
219     for(int i = node->height; i < set.head.size();
220         i++)
221         stream << "└┘"; stream << std::endl;
222
223     node = node->forward[0];
224 }
225
226 // Print END
227 for(int i = 0; i < set.head.size(); i++)
228     stream << "└┘"; stream << std::endl;
229 for(int i = 0; i < set.head.size(); i++)
230     stream << "└v┘"; stream << std::endl;
231 for(int i = 0; i < set.head.size(); i++)
232     stream << "——┘"; stream << std::endl;
233 for(int i = 0; i < set.head.size(); i++)
234     stream << "└┘"; stream << std::endl;
235 for(int i = 0; i < set.head.size(); i++)
236     stream << "——┘"; stream << std::endl;
237
238 return stream;
239 }
240 private:
241 struct Node{

```



```

242     const T& data;
243     Node** forward;
244     Node** backward;
245     int height;
246     Node(const T& data, int height) : data(data),
        height(height) {
247         forward = new Node*[height];
248         backward = new Node*[height];
249
250         for(int i = 0; i < height; i++){
251             forward[i] = nullptr;
252             backward[i] = nullptr;
253         }
254     }
255
256     ~Node(){
257         delete [] forward;
258         delete [] backward;
259     }
260 };
261
262
263 // data members
264 std::vector<Node*> head; // First element for
        every height
265 double q; // Probability of not having another
        level
266
267 // Internal procedures
268 Node* search(const T& x) const{
269     int height = head.size() - 1;
270     while(height > 0 && head[height] -> data > x)
        height--;
271     if(height < 0) return nullptr;
272
273     Node* res = head[height];
274     while(height >= 0){
275         if(res -> forward[height] && res -> forward[height]
        -> data <= x){
276             res = res -> forward[height];
277         } else{
278             height--;
279         }
280     }
281
282     return res;

```

```

283     }
284
285     int randomHeight(){
286         static std::random_device rd;
287         static std::mt19937 gen(rd());
288         static std::uniform_real_distribution<> dis
            (0.0,1.0);
289
290         int height = 1;
291         while(dis(gen) > q) height++;
292
293         return height;
294     }
295 };

```

```

1  //*****
2  // File:    assignment3.cpp
3  // Author:  Daniel Benedi Garcia
4  // Date:    8th December 2022
5  // Coms:    Main file of assignment 3 with functions
6  //           to do the experiments and debug.
7  //*****
8
9  #include <iostream>
10 #include <sstream>
11 #include <vector>
12 #include <fstream>
13 #include <random>
14 #include <algorithm>
15 #include "Set.h"
16
17 #define M 100
18
19 void debug(){
20     std::random_device rd;
21     std::mt19937_64 g(rd());
22     int n = 7, i = 60;
23
24     // Initialize vector with increasing numbers
25     // and then shuffle it
26     std::vector<int> elems(n, 0);
27     std::iota(elems.begin(), elems.end(), 1);
28     std::shuffle(elems.begin(), elems.end(), g);
29
30

```

```

31     std::cout << "Input_=" << elems[0];
32     for(int i = 1; i < elems.size(); i++)
33         std::cout << "," << elems[i];
34     std::cout << "]" << std::endl;
35
36     double q = i/100.;
37     Set<int> set(q);
38
39     // Insert the numbers
40     for(int i = 0; i < elems.size(); i++)
41         set.insert(elems[i]);
42
43     int nodes;
44     std::cout << "levels_=";
45     for(int lvl = 0; nodes=set.nodes_lvl(lvl); lvl++)
46         std::cout << nodes << ",";
47     std::cout << std::endl;
48
49     std::cout << "Skiplist" << std::endl;
50     std::cout << set << std::endl;
51
52     std::cout << "Total_search_cost_=" << set.
53         total_search_cost() << std::endl;
54 }
55 void experiment(){
56     std::random_device rd;
57     std::mt19937_64 g(rd());
58
59     std::ofstream cost("total_search_cost.csv");
60     std::ofstream mem("total_memory_use.csv");
61
62     #pragma omp parallel for collapse(2)
63     for(int n = 2000; n <= 20000; n += 100){
64         for(int i = 10; i < 100; i += 10){
65             double q = i/100.;
66
67             std::stringstream buf1, buf2;
68             buf1 << n << "," << q;
69             buf2 << n << "," << q;
70
71             std::vector<int> elems(n, 0);
72             std::iota(elems.begin(), elems.end(), 1);
73             for(int m = 0; m < M; m++){
74                 Set<int> set(q);
75

```

```

76         std::shuffle(elems.begin(), elems.end(), g);
77
78         for(int i = 0; i < elems.size(); i++)
79             set.insert(elems[i]);
80
81         buf1 << "," << set.total_search_cost();
82         buf2 << "," << set.number_pointers();
83     }
84
85     buf1 << std::endl;
86     buf2 << std::endl;
87
88     #pragma omp critical
89     cost << buf1.str();
90
91     #pragma omp critical
92     mem << buf2.str();
93 }
94 }
95 }
96
97 int main(int argc, char* argv[]) {
98     #ifdef DEBUG
99         debug();
100     #else
101         experiment();
102     #endif
103 }

```

References

- [1] Peter Kirschenhofer, Conrado Martínez, and Helmut Prodinger. “Analysis of an optimized search algorithm for skip lists”. In: *Theoretical Computer Science* 144.1 (1995), pp. 199–220. ISSN: 0304-3975. DOI: 10.1016/0304-3975(94)00296-U. URL: <https://www.sciencedirect.com/science/article/pii/030439759400296U>.
- [2] William Pugh. “Skip Lists: A Probabilistic Alternative to Balanced Trees”. In: *Commun. ACM* 33.6 (June 1990), pp. 668–676. ISSN: 0001-0782. DOI: 10.1145/78973.78977. URL: <https://doi.org/10.1145/78973.78977>.
- [3] W. Szpankowski and V. Rego. “Yet another application of a binomial recurrence order statistics”. In: *Computing* 43.4 (1990), pp. 401–410. DOI: 10.1007/bf02241658.