

Sprawozdanie Projekt 1 - ADT		
Rafał Jasiński 259384	02.04.2022r	Dr hab. inż. Andrzej Rusiecki
Zadania na ocenę bdb (5.0)		
Link do repozytorium: https://gitlab.com/JasinskiR259384/pamsi-2022\		

1 Wstęp

Celem zadania było zaprojektowanie i zaimplementowanie odpowiedniej abstrakcyjnej struktury danych (ADT) do przechowywania oraz sortowania danych składających się z pary: liczby porządkowej oraz wiadomości (typ nie został określony). Przeanalizować czas działania - złożoność obliczeniową proponowanego rozwiązania.

2 Podejście do problemu

Ze względu na fakt, że dane, które otrzymujemy nie są posortowane, pojawia się problem jaką strukturę należy wybrać, aby łatwo się móc przemieszczać pomiędzy jej elementami w celu posortowania/wypisania posortowanej wiadomości.

Na tym etapie ograniczyłem znane mi abstrakcyjne struktury danych do:

- BST - binary search tree - binarne drzewo poszukiwań
- hash table (hash map)
- heap - kopiec

Postanowiłem rozważać moje rozwiązanie problemu na nich ze względu na wstępną możliwość posortowania dodawanych danych.

Na wstępie założyłem, aby mój algorytm miał jak najmniejszą złożoność obliczeniową.

2.1 Hash map

Wykorzystanie hash map'y było moim pierwszym pomysłem na podejście do danego problemu. Jednak napisanie własnej klasy działającej na zasadzie **std::map** jest dość nieefektywne. Pojawia się tutaj problem hashowania naszego klucza, który jest liczbą porządkową. Najprostszym sposobem było by użycie

```

1 struct MyKeyHash {
2     unsigned long (const uint32_t &key) const {
3         return key % 10;
4     }
5 };

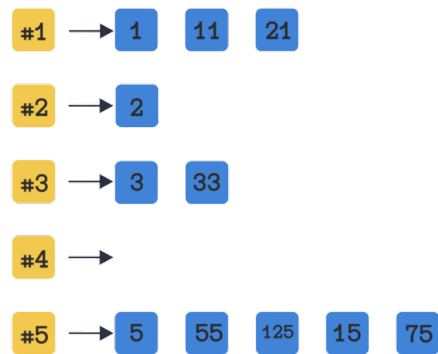
```

Listing 1: Hash function

Jednakże pojawia się tutaj problem z przyszłym poszukiwaniem naszego konkretnego elementu. Musielibyśmy przejść po tablicach z odpowiednio zahashowanym kluczem a następnie dostać się do naszego elementu. Znalezienie konkretnego elementu w tym momencie możemy wyznaczyć jako ilość kluczy hashujących (przeszukanie $O(n)$) - n w tym przypadku jest zależne od funkcji **mod**) a następnie w danym kluczu przeszukanie n elementów. Co ostatecznie daje nam $O(n^2)$.

Problem ten występuje w przypadku kiedy naszym kluczem hashującym jest liczba. Można spróbować zmniejszyć złożoność tego algorytmu stosując \sqrt{key} .

Dlatego na tym etapie porzuciłem pomysł hash map'y.



Rysunek 1: Reprezentacja przechowywania elementów w hash map'ie (jako tablica dwuwymiarowa)

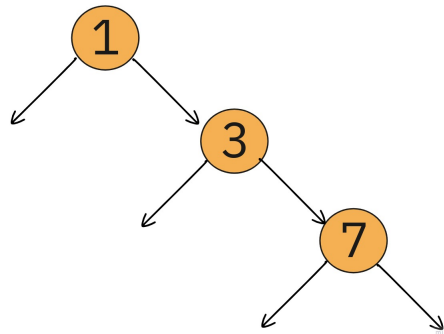
2.2 BST - binary search tree - binarne drzewo poszukiwań

W przypadku tej struktury udało mi się napisać już działający algorytm, który jako wejście przyjmował parę (liczba porządkowa, wiadomość)

```
1 std::pair<uint32_t, message>
```

Listing 2: Para danych

A na wyjściu przedstawiał nam już uporządkowaną wiadomość. Złożoność umieszczenia elementu na drzewie w przypadku tego algorytmu wynosi $O(\log(n))$. Jest to jednak wartość oczekiwana i nie gwarantuje, że zawsze będzie tyle wynosić. W najgorszym przypadku możemy otrzymać złożoność $O(n)$ - dzieje się tak w przypadku kiedy otrzymujemy na wejście już posortowane dane - w dowolną stronę. Obiekty wtedy są umieszczane na jednej gałęzi.

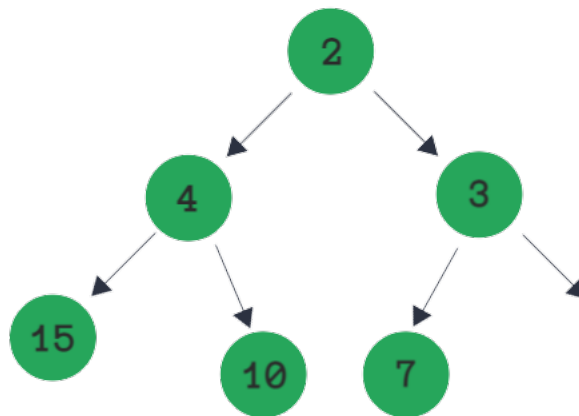


Rysunek 2: Najgorszy możliwy przypadek - otrzymujemy na wejście posortowane dane

Ze względu na te przypadki postanowiłem poszukać jeszcze innego rozwiązania - kopiec.

2.3 Heap - kopiec oparty na tablicy

Elementy w przypadku tej struktury są umieszczane od lewej do prawej, schodząc kolejno w dół do kolejnych "pokoleń". Gwarantuje nam to że rozmiar naszego kopca zawsze będzie $O(\log(n))$. Złożoność umieszczenia elementu na kopiec również wynosi $O(\log(n))$. Ze względu na wypełnianie od lewej do prawej nie mamy przypadków, w których zapełniana jest tylko jedna gałąź, co niweluje nam przypadki możliwe do wystąpienia w przypadku drzewa BST.



Rysunek 3: Na kopcu nie obowiązuje też zależność lewego - prawego dziecka. Przykład posortowanego kopca dla podanych danych

2.4 Wybrana struktura danych

Podsumowując zalety i wady rozważanych przeze mnie abstrakcyjnych struktur danych, postanowiłem rozwiązać problem przedstawiony w zadaniu wykorzystując kopiec oparty na tablicy.

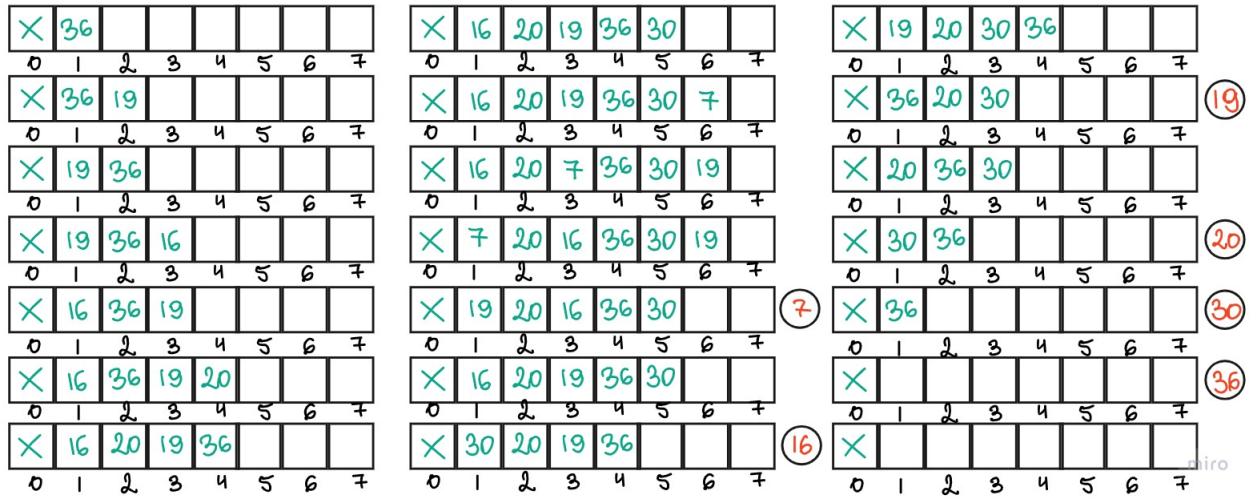
3 Rozwiązanie

Aby móc w jakikolwiek sposób testować algorytm, posłużyłem się dodatkową klasą *Packets*. Działa ona w dwóch trybach:

1. generuje parę liczb (liczba porządkowa, liczba - która jest równa $n-lp$ [gdzie n ilość liczb, lp liczba porządkowa]), a następnie wykonywana jest operacja pomieszania danych.
2. wczytuje dane z pliku, który musi zostać przygotowany w odpowiedni sposób tj. liczba porządkowa wiadomości, dowolny separator będący znakiem białym oraz sama wiadomość do znaku enter `"\n"`

3.1 Zasada działania

Dane te są następnie wczytywane przez funkcję klasy *Heap*, która czyta każdą parę z osobna i wykonuje operację dodawania elementu do struktury. Zazwyczaj w kopcu elementy największe będą znajdować się na samej górze. Jednak w moim przypadku postanowiłem zrobić odwrotny kopiec, a więc to wartość najmniejsza będzie na szczycie. Sortowanie kopca polega na tak zwanym kopcowaniu - heapifying, czyli zamianie elementu aktualnie dodawanego z jego rodzicem do momentu, aż nie natkniemy się na wartość mniejszą (działa na zasadzie rekurencji). Kiedy zakończy się proces dodawania elementów na kopiec, odczytanie posortowanych wartości odbywa się poprzez odczytanie pierwszego elementu (który w tym przypadku będzie najmniejszy), a następnie zamiany pierwszego elementu z ostatnim. Wyzerowaniu/usunięcie ostatniego elementu. A następnie odbywa się operacja "spychania" naszego największego elementu na dół. Po tej operacji nasza nowa najmniejsza wartość znowu pojawi się na pierwszym miejscu.



Rysunek 4: Pokazanie zasady działania algorytmu kopca (dla przykładowych danych)

3.2 Przechowywanie danych

Przechowywanie danych wykorzystując dynamiczną alokację pamięci może wydawać się trochę mało optymalne obliczeniowo. Ponieważ algorytm wykorzystuje operacje na tablicy posiada dostęp do i -tego elementu (w naszym przypadku będziemy wykorzystywać głównie operacje $i/2$). Gdybyśmy chcieli zadbać o optymalizację pamięciową moglibyśmy użyć listy dwukierunkowej z iteracją by móc przemieszczać się po elementach. Jednak stwarza to problem potrzeby przemieszczania się po liście, w celu dostępu do naszego elementu. Dlatego postanowiłem zastosować operację alokacji statycznej z automatyczną opcją rozszerzania tablicy (podwojenie dostępnej pamięci)

3.2.1 Początek iteracji od 1 - tab[1]

Dane, które są dodawane na kopiec są zapisywane w tablicy dopiero od pierwszego indeksu w celu łatwiejszej implementacji dostępu do rodziców oraz dzieci.

4 Analiza złożoności obliczeniowej oraz pamięciowej

```

1 template <typename message>
2 void Heap<message>::heapify(uint32_t node) {
3     if (node == 1 || packetArr[node / 2] < packetArr[node]) return; // 0(1)
4     swap(packetArr[node / 2], packetArr[node]); // 0(1)
5     heapify(node / 2); // 0(log(n))
6 }
7
8 template <typename message>
9 std::pair<uint32_t, message> Heap<message>::minValue() { // 0(1)
10     return packetArr[1];
11 }
12
13 template <typename message>
14 void Heap<message>::pushDown(uint32_t node) {
15     if (node * 2 > size) return; // 0(1)
16     if (node * 2 + 1 > size) { // 0(1)
17         if (packetArr[node * 2] < packetArr[node]) { // 0(1)
18             swap(packetArr[node], packetArr[node * 2]); // 0(1)
19             this->pushDown(node * 2); // 0(log(n))
20         }
21         return; // 0(1)
22     }
23     if (packetArr[node * 2] < packetArr[node * 2 + 1]) { // 0(1)
24         if (packetArr[node * 2] < packetArr[node]) { // 0(1)
25             swap(packetArr[node], packetArr[node * 2]); // 0(1)
26             this->pushDown(node * 2); // 0(log(n))
27         }
28     }
29 }

```

```

28     return; // 0(1)
29 }
30 if (packetArr[node * 2 + 1] < packetArr[node]) { // 0(1)
31     swap(packetArr[node], packetArr[node * 2 + 1]); // 0(1)
32     this->pushDown(node * 2 + 1); // 0(log(n))
33 }
34 }
35
36 template <typename message>
37 void Heap<message>::pop() {
38     swap(packetArr[1], packetArr[size]); // 0(1)
39     packetArr[size--] = std::make_pair(NULL, NULL); // 0(1)
40     this->pushDown(1); // 0(log(n))
41 }
42
43 template <typename message>
44 void Heap<message>::insertKey(std::pair<uint32_t, message> element) {
45     if (size + 1 >= sizeofTab) { // 0(1)
46         sizeofTab *= 2; // 0(1)
47         resize(); // 0(1) - amortyzowany
48     }
49     packetArr[++size] = element; // 0(1)
50     heapify(size); // 0(log(n))
51 }
52
53 template <typename message>
54 void Heap<message>::resize() {
55     std::pair<uint32_t, message> *tmp =
56         new std::pair<uint32_t, message>[sizeofTab]; // 0(1)
57     for (uint32_t i = 0; i <= size; ++i) {
58         tmp[i] = packetArr[i]; // 0(1)
59     } // 0(n)
60     delete[] packetArr; // 0(1)
61     packetArr = tmp; // 0(1)
62 } // 0(n)

```

Listing 3: Kod z analizą złożoności

Powiększanie tablicy odbywa się poprzez podwojenie dostępnej pamięci. Początkowo funkcja ta będzie wykonywać się znacznie częściej, niż gdybyśmy ustawili stały przyrost. Jednakże wraz z rosnącą liczbą elementów ilość wykonywanych operacji powiększania maleje. Jej liczba wywołań wyniesie więc $O(\log(n))$. Natomiast liczba wszystkich przepisania nie przekroczy $2n$. Stąd amortyzowana złożoność wywołania funkcji *resize()* w funkcji *insertKey(element)* będzie $O(1)$. Podsumowując, złożoność całego algorytmu dla n elementów wyniesie $O(n \log(n))$

5 Testy i uwagi

Na GitLabie w folderze *test* znajdują się przykładowe pomieszczone pliki tekstowe jako forma testów. Dostępny jest również skrypt bash'owy wykonujący operację sprawdzania poprawności pliku docelowego oraz wyjściowego programu.

Aby móc uruchomić program należy go skompilować za pomocą funkcji *make* w terminalu.

Pliki o nazwach *Tree.** nie podlegają pod rozwiązanie zadania (jest to kompletny algorytm, który był rozważany, lecz został porzucony na rzecz innego)

Folder *fml* zawiera pliki klas szablonowych i zostały one wydzielone w celu posegregowania plików.

6 Wnioski

1. Patrząc na przedstawione algorytmy, nie istnieje szybsza metoda "posortowania" wiadomości niż $O(n \log(n))$
2. Biorąc pod uwagę koszt zamortyzowany funkcji *resize()* pozwala to na osiągnięcie końcowej złożoności *insertKey(element)* jako $O(\log(n))$
3. Złożoność obliczeniowa $O(\log(n))$ jest na tyle wydajna czasowo, że przy próbie sporządzenia wykresu nawet przy ilości elementów wynoszących 50 000 000, charakterystyka czasowa niewiele odbiegała od funkcji liniowej.
4. Porównując rozważane przeze mnie abstrakcyjne struktury, kopiec wydaje się mieć jak najmniej wad.