## САНКТ-ПЕТЕРБУРГСКИЙ НАЦИОНАЛЬНЫЙ ИССЛЕДОВАТЕЛЬСКИЙ УНИВЕРСИТЕТ

# ИНФОРМАЦИОННЫХ ТЕХНОЛОГИЙ, МЕХАНИКИ И ОПТИКИ ФАКУЛЬТЕТ ИНФОКОММУНИКАЦИОННЫХ ТЕХНОЛОГИЙ

# Отчет по лабораторной работе №6 по курсу «Алгоритмы и структуры данных» Тема: Хеширование. Хеш-таблицы Вариант 15

Выполнил:

Скворцов Д.А.

K3140

Проверил:

Афанасьев А.В.

Санкт-Петербург 2024 г.

### Содержание отчета

Содержание отчета	2
Задачи по варианту	3
Задача №2. Телефонная книга	3
Задача №4. Прошитый ассоциативный массив	7
Дополнительные задачи	12
Задача №6. Фибоначчи возвращается	12
Задача №8. Почти интерактивная хеш-таблица	15
Вывод	19

#### Задачи по варианту

#### Задача №2. Телефонная книга

В этой задаче ваша цель - реализовать простой менеджер телефонной книги. Он должен уметь обрабатывать следующие типы пользовательских запросов:

- add number name это команда означает, что пользователь добавляет в телефонную книгу человека с именем name и номером телефона number. Если пользователь с таким номером уже существует, то ваш менеджер должен перезаписать соответствующее имя.
- del number означает, что менеджер должен удалить человека с номером из телефонной книги. Если такого человека нет, то он должен просто игнорировать запрос.
- find number означает, что пользователь ищет человека с номером телефона number. Менеджер должен ответить соответствующим именем или строкой «not found» (без кавычек), если такого человека в книге нет.
- Формат ввода / входного файла (input.txt). В первой строке входного файла содержится число N (1 ≤ N ≤ 10<sup>5</sup>) количество запросов. Далее следуют N строк, каждая из которых содержит один запрос в формате, описанном выше. Все номера телефонов состоят из десятичных цифр, в них нет нулей в начале номера, и каждый состоит не более чем из 7 цифр. Все имена представляют собой непустые строки из латинских букв, каждая из которых имеет длину не более 15. Гарантируется при проверке, что не будет человека с именем «not found».
- Формат вывода / выходного файла (output.txt). Выведите результат каждого поискового запроса find имя, соответствующее номеру телефона, или «not found» (без кавычек), если в телефонной книге нет человека с таким номером телефона. Выведите по одному результату в каждой строке в том же порядке, как были заданы запросы типа find во входных данных.
- Ограничение по времени. 6 сек.
- Ограничение по памяти. 512 мб.

#### Листинг кода.

```
class Node:
   def init_ (self, key, value):
       self.key = key
       self.value = value
       self.next = None
       self.size = 0
       self.max size = max size
   def hash(self, k):
   def insert(self, k, v):
       if self.table[index] is None:
           self.table[index] = Node(k, v)
           self.size += 1
           while cur node is not None:
                   cur node.value = v
               cur node = cur node.next
           new node = Node(k, v)
           new node.next = self.table[index]
           self.table[index] = new node
           self.size += 1
   def search(self, k):
       cur node = self.table[index]
           if cur node.key == k:
           cur node = cur node.next
```

```
def remove(self, k):
       prev = None
                if prev is None:
                    prev.next = cur_node.next
            prev = cur node
       return self.size
def solution(commands):
   dct = HashTable(len(commands))
       cmd = command.split()
           dct.remove(int(cmd[1]))
           ret = dct.search(int(cmd[1]))
           ans.append("not found" if (ret is None) else ret)
    return ans
```

Обработчик запросов реализован при помощи самописного красса Хеш-таблицы. Для каждого номера генерируется хэш, по которому происходит гораздо более быстрый поиск в словаре. Аналогично происходит удаление и добавление.

Lab6. Task 2	
Input	Output
   8   find 3839442   add 123456 me   add 0 granny   find 0   find 123456   del 0   del 0	not found   granny   me   not found

Тест 100 элементов:

Время работы: 0.000103 секунд Затрачено памяти: 7.2 Килобайт

Тест 10е3 элементов:

Время работы: 0.00188 секунд Затрачено памяти: 67.1 Килобайт

Тест 10е5 элементов:

Время работы: 0.14511 секунд Затрачено памяти: 6.5 Мегабайт

	Время выполнения	Затраты памяти
Нижняя граница диапазона значений входных данных из текста задачи	0.000103 секунд	7.2 Килобайт
Медиана диапазона значений входных данных из текста задачи	0.00188 секунд	67.1 Килобайт
Верхняя граница диапазона значений входных данных из текста задачи	0.14511 секунд	6.5 Мегабайт

Вывод по задаче:

Алгоритм работает за линейное время, каждая операция по отдельности занимает константное время. В зависимости от качества алгоритма хеширования время может отличаться.

#### Задача №4. Прошитый ассоциативный массив

Реализуйте прошитый ассоциативный массив. Ваш алгоритм должен поддерживать следующие типы операций:

- $\bullet$  get x если ключ x есть в множестве, выведите соответствующее ему значение, если нет, то выведите .
- prev x вывести значение, соответствующее ключу, находящемуся в ассоциативном массиве, который был вставлен позже всех, но до x, или, 6 если такого нет или в массиве нет x.
- next x вывести значение, соответствующее ключу, находящемуся в ассоциативном массиве, который был вставлен раньше всех, но после x, или, если такого нет или в массиве нет x.
- put x y поставить в соответствие ключу x значение y. При этом следует учесть, что
  - если, независимо от предыстории, этого ключа на момент вставки в массиве не было, то он считается только что вставленным и оказывается самым последним среди добавленных элементов – то есть, вызов next с этим же ключом сразу после выполнения текущей операции put должен вернуть
  - если этот ключ уже есть в массиве, то значение необходимо изменить, и в этом случае ключ не считается вставленным еще раз, то есть, не меняет своего положения в порядке добавленных элементов.
- delete x удалить ключ x. Если ключа в ассоциативном массиве нет, то ничего делать не надо.
- Формат входного файла (input.txt). В первой строке входного файла находится строго положительное целое число операций N, не превышающее 5 · 10<sup>5</sup>. В каждой из последующих N строк находится одна из приведенных выше операций. Ключи и значения операций строки из латинских букв длиной не менее одного и не более 20 символов.

- Формат выходного файла (output.txt). Выведите последовательно результат выполнения всех операций get, prev, next. Следуйте формату выходного файла из примера.
- Ограничение по времени. 4 сек.
- Ограничение по памяти. 256 мб.

#### Листинг кода.

```
class Node:
       self.key = key
       self.value = value
       self.next = None
       self.next arr = None
       self.prev arr = None
       self.size = 0
       self.table = [None] * max size
   def hash(self, k):
       index = self. hash(k)
       if self.table[index] is None:
           self.size += 1
           return self.table[index]
       while cur node is not None:
                return cur node
           cur node = cur node.next
```

```
new node = Node(k, v)
       new node.next = self.table[index]
       self.size += 1
       return new node
   def remove(self, k) -> Node:
       prev = None
            if cur node.key == k:
                if prev is None:
                    self.table[index] = cur node.next
                   prev.next = cur node.next
                self.size -= 1
           prev = cur node
       return self.size
   def contains (self, k):
       return self.search(k) is None
class AssociationArray(HashTable):
```

```
self.head = None
    def insert(self, k, v):
        node = super().insert(k, v)
        if self.head is not None:
            connect(self.head, node)
        self.head = node
   def remove(self, k):
        node = super().remove(k)
           reconnect (node)
   def prev(self, k):
        node = super().search(k, mode="Node")
        if node is not None and node.prev arr is not None:
            return node.prev arr.value
        node = super().search(k, mode="Node")
            return node.next arr.value
def connect(node1, node2):
   if node1 is not None:
    if node2 is not None:
        node2.prev arr = node1
def reconnect(node):
    connect(node.prev_arr, node.next_arr)
def solution(commands):
   dct = AssociationArray(len(commands))
    for command in commands:
        cmd = command.split()
            dct.insert(cmd[1], cmd[2])
        elif cmd[0] == 'delete':
            dct.remove(cmd[1])
```

```
elif cmd[0] in ('get', 'prev', 'next'):
    if cmd[0] in 'get':
        ret = dct.search(cmd[1])
    elif cmd[0] == 'prev':
        ret = dct.prev(cmd[1])
    else:
        ret = dct.next(cmd[1])
    ans.append("<none>" if (ret is None) else ret)
return ans
```

Для каждого элемента при добавлении также запоминаем следующий и предыдущий, чтобы при помощи связного списка производить быстрое удаление и поиск следующего/предыдущего.

Lab6. Task 4	
Input	Output
14   put zero a   put one b   put two c   put three d   put four e   get two   prev two   next two   delete one   delete three   get two   prev two	c
next two	

```
Lab6. Task 4

Тест 100 элементов:
Время работы: 0.000121 секунд
Затрачено памяти: 8.9 Килобайт

Тест 10е3 элементов:
Время работы: 0.001698 секунд
Затрачено памяти: 83.7 Килобайт

Тест 5*10е5 элементов:
Время работы: 0.518511 секунд
Затрачено памяти: 41.2 Мегабайт
```

	Время выполнения	Затраты памяти
Нижняя граница диапазона значений входных данных из текста задачи	0.000121 секунд	8.9 Килобайт
Медиана диапазона значений входных данных из текста задачи	0.001698 секунд	83.7 Килобайт
Верхняя граница диапазона значений входных данных из текста задачи	0.518511 секунд	41.2 Мегабайт

Вывод по задаче:

Алгоритм работает.

#### Дополнительные задачи

#### Задача №6. Фибоначчи возвращается

Вам дается последовательность чисел. Для каждого числа определите, является ли оно числом Фибоначчи. Напомним, что числа Фибоначчи определяются, например, так: F0 = F1 = 1 (1) Fi = Fi-1 + Fi-2 для  $i \ge 2$ .

- Формат ввода / входного файла (input.txt). Первая строка содержит одно число N ( $1 \le N \le 10^6$ ) количество запросов. Следующие N строк содержат по одному целому числу. При этом соблюдаются следующие ограничения при проверке: 1. Размер каждого числа не превосходит 5000 цифр в десятичном представлении. 2. Размер входа не превышает  $1 \, \text{Mб}$ .
- Формат вывода / выходного файла (output.txt). Для каждого числа, данного во входном файле, выведите «Yes», если оно является числом Фибоначчи, и «No» в противном случае.
- Ограничение по времени. 2 сек.
- Ограничение по памяти. 128 мб.

#### Листинг кода.

```
def fill_fib(x, last_n):
   global fibs by num
   while True:
        fibs.add(new fib)
       fibs by num.append(new fib)
        elif new fib == x:
def solution(lst):
   global fibs
   global fibs by num
   fibs by num = [1, 1]
   res = []
       x = int(i)
       if x not in fibs:
            if x > fibs by num[last n]:
                res.append("Yes" if verdict else "No")
                res.append("No")
            res.append("Yes")
    return res
```

Для решения функция fill\_fib находит все новые числа фибоначчи, начиная с последнего известного, и возвращает результат, как только доходит до проверяемого числа или перескакивает его.

Lab6. Task 6		
Input	Output	
	+	
8	Yes	
1	Yes	
2	Yes	
] 3	No	
4	Yes	
5	No	
6	No i	
j 7	Yes	
8	i i	

Гест 100 элементов:

Время работы: 5.4e-05 секунд Затрачено памяти: 3.8 Килобайт

Тест 10е6 элементов:

Время работы: 0.252006 секунд Затрачено памяти: 8.1 Мегабайт

Гест 5000 разрядов:

Время работы: 0.052138 секунд Ватрачено памяти: 28.1 Мегабайт

	Время выполнения	Затраты памяти
Нижняя граница диапазона значений входных данных из текста задачи	0.000054 секунд	3.8 Килобайт
Медиана диапазона значений входных данных из текста задачи	0.252006 секунд	8.1 Мегабайт
Верхняя граница диапазона значений входных данных из текста задачи	0.052138 секунд	28.1 Мегабайт

#### Вывод по задаче:

Время работы алгоритма - O(n), но он расходует память на сохранение всех найденных чисел Фибоначчи. Каждое последующее число считается

итеративно и при помощи последних известных, что позволяет не экономить память.

#### Задача №8. Почти интерактивная хеш-таблица

В данной задаче у Вас не будет проблем ни с вводом, ни с выводом. Просто реализуйте быструю хеш-таблицу. В этой хеш-таблице будут храниться целые числа из диапазона  $[0; 10^{15} - 1]$ . Требуется поддерживать добавление числа х и проверку того, есть ли в таблице число х. Числа, с которыми будет работать таблица, генерируются следующим образом. Пусть имеется четыре целых числа N, X, A, B такие что:

- $1 \le N \le 10^7$
- $1 \le X \le 10^{15}$
- $1 \le A \le 10^3$
- $1 \le B \le 10^{15}$

Требуется N раз выполнить следующую последовательность операций:

- Если X содержится в таблице, то установить A ← (A + AC ) mod 103 , B ← (B + BC ) mod  $10^{15}$  .
- Если X не содержится в таблице, то добавить X в таблицу и установить A ← (A + AD) mod  $10^3$ , B ← (B + BD) mod  $10^{15}$ .
- Установить  $X \leftarrow (X \cdot A + B) \mod 1015$ . Начальные значения X, A и B, а также N, AC, BC, AD и BD даны во входном файле. Выведите значения X, A и B после окончания работы.
- Формат входного файла (input.txt). В первой строке входного файла содержится четыре целых числа N, X, A, B. Во второй строке содержится еще четыре целых числа AC, BC, AD и BD такие что  $0 \le AC$ , AD < 103,  $0 \le BC$ , BD < 1015.
- **Формат выходного файла (output.txt).** Выведите значения X, A и B после окончания работы. 11
- Ограничение по времени. 5 сек.
- Ограничение по памяти. 256 мб.

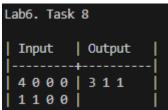
#### Листинг кода.

```
class HashSet:
    def __init__(self, max_size):
        self.size = 0
        self.table = array('q', [EMPTY] * max_size)
```

```
self.aux_list = []
def hash(self, k):
        self.table[index] = k
        self.size += 1
        self.aux list.append(lst)
        lst.append(k)
    self.size += 1
    if cur node == EMPTY:
        return cur node == k
    return k in self.aux list[-cur node]
    return self.size
```

```
def solution(lst):
    n,x,a,b, ac,bc,ad,bd = lst
    size = n * 2 if n < 10**6 else 2017963
    hash_set = HashSet(size) # prime for distribution
    t = 10 ** 15
    for _ in range(n):
        if x in hash_set:
            a = (a + ac) % 1000
            b = (b + bc) % t
        else:
            a = (a + ad) % 1000
            b = (b + bd) % t
            hash_set.add(x)
        x = (x * a + b) % t
    return x, a, b</pre>
```

Чтобы сделать самописный hash-set более эффективным для задания размера его памяти используется простое число, так как распределение хешей в нем более равномерное. Также для оптимизации по памяти используется array, позволяющий хранить либо число, либо ссылку на несколько чисел в выносном массиве.



```
Тест 100 элементов:
Время работы: 0.000153 секунд
Затрачено памяти: 6.7 Килобайт
Тест 10е5 элементов:
Время работы: 0.170469 секунд
Затрачено памяти: 7.2 Мегабайт
Тест 10е6 элементов:
Время работы: 1.392178 секунд
Затрачено памяти: 42.3 Мегабайт
```

	Время выполнения	Затраты памяти
Нижняя граница диапазона значений	0.000153 секунд	6.7 Килобайт

входных данных из текста задачи		
Медиана диапазона значений входных данных из текста задачи	0.170469 секунд	7.2 Мегабайт
Верхняя граница диапазона значений входных данных из текста задачи	1.392178 секунд	42.3 Мегабайт

#### Вывод по задаче:

Алгоритм работает быстро, за линейное время и при этом затрачивает минимум памяти. поскольку не создает связного списка, если элемент с определённым хешем единственен.

#### Вывод

Структуры данных на основе хэширования (Hash-table, Hash-set) позволяют значительно ускорять работу с данными, требующую частого поиска или добавления/удаления объектов по их значению. Они существенно уменьшают время выполнения программы. Взамен делается упор на качество хеш-функции и выделение памяти, которые надо компенсировать при постоянном увеличении размеров.