

# Оглавление

2	<b>1. Понятие алгоритма</b>	<b>5</b>
3	1.1. Машина с произвольным доступом к памяти . . . . .	5
4	1.2. Алгоритмический язык . . . . .	6
5	1.3. Определение алгоритма . . . . .	7
6	1.4. O-символика . . . . .	8
7	1.5. Числа Фибоначчи . . . . .	10
8	<b>2. Структуры данных</b>	<b>15</b>
9	2.1. Массивы . . . . .	15
10	2.1.1. Расширяемый массив с аддитивной схемой выделения . . . . .	16
11	2.1.2. Расширяемый массив с мультипликативной схемой выделе-	
12	ния. . . . .	17
13	2.2. Списки . . . . .	18
14	2.2.1. Односвязный список . . . . .	18
15	2.2.2. Двусвязный список . . . . .	18
16	2.2.3. Сравнение с массивами . . . . .	19
17	2.3. Абстрактные типы данных . . . . .	19
18	2.3.1. Стек . . . . .	20
19	2.3.2. Очередь . . . . .	21
20	2.4. Двухсторонняя очередь . . . . .	23
21	2.5. Амортизационный анализ . . . . .	23
22	2.5.1. Метод учётных стоимостей . . . . .	23
23	2.5.2. Метод предоплаты . . . . .	28
24	<b>3. Разделяй и властвуй</b>	<b>29</b>
25	3.1. Переход к подзадачам . . . . .	29
26	3.2. Двоичный поиск . . . . .	29
27	3.3. Умножение $n$ -битовых чисел . . . . .	30
28	3.4. Сортировка слиянием . . . . .	30
29	3.5. Основная теорема о рекуррентных соотношениях . . . . .	32
30	3.5.1. Как побороть аддитивные константы? . . . . .	32
31	3.6. Умножение матриц . . . . .	32
32	<b>4. Сортировка</b>	<b>33</b>
33	4.1. Квадратичные сортировки . . . . .	33
34	4.2. Нижняя оценка на сортировку сравнениями . . . . .	34
35	4.3. Сортировка кучей . . . . .	35
36	4.4. Быстрая сортировка . . . . .	36
37	4.4.1. Анализ времени работы . . . . .	36
38	4.4.2. Сортировка массива с повторяющимися элементами . . . . .	37
39	4.5. Линейные сортировки . . . . .	37

40	4.5.1. Сортировка подсчётом . . . . .	37
41	4.5.2. Цифровая сортировка . . . . .	38
42	4.5.3. Блочная сортировка . . . . .	38
43	<b>5. Задачи RMQ и LCA</b>	<b>39</b>
44	5.1. Задача о минимуме на отрезке . . . . .	39
45	5.1.1. Динамическая задача RMQ . . . . .	39
46	5.1.2. Статическая задача RMQ . . . . .	45
47	5.2. Задача о наименьшем общем предке . . . . .	46
48	<b>6. Деревья поиска</b>	<b>51</b>
49	6.1. АДД с быстрым поиском . . . . .	51
50	6.2. Представление корневых деревьев в памяти . . . . .	52
51	6.3. Двоичные деревья поиска . . . . .	53
52	6.3.1. Операции поиска . . . . .	53
53	6.3.2. Добавление элемента . . . . .	54
54	6.3.3. Удаление элемента . . . . .	55
55	6.4. AVL-дерево . . . . .	57
56	6.4.1. Сбалансированность . . . . .	57
57	6.4.2. Вращения . . . . .	57
58	6.5. Splay-дерево . . . . .	59
59	6.5.1. Запросы к splay-дереву . . . . .	59
60	6.5.2. Операция splay . . . . .	61
61	6.5.3. Анализ сложности . . . . .	64
62	6.6. Декартово дерево . . . . .	67
63	6.6.1. Определение и теорема существования . . . . .	67
64	6.6.2. Эффективное построение . . . . .	68
65	6.6.3. Операции . . . . .	69
66	6.6.4. Дуча . . . . .	71
67	6.7. B-деревья . . . . .	71
68	<b>7. Хеширование</b>	<b>73</b>
69	7.1. Прямая адресация . . . . .	73
70	7.2. Хеш-таблица . . . . .	73
71	7.2.1. Хеш-функция . . . . .	74
72	7.2.2. Методы разрешения коллизий . . . . .	74
73	7.2.3. Примеры хеш-функций . . . . .	75
74	7.2.4. Детали реализации . . . . .	77
75	7.3. Гипотеза равномерного хеширования . . . . .	77
76	7.4. Универсальное хеширование . . . . .	79
77	7.5. Совершенное хеширование . . . . .	82
78	<b>8. Числовые алгоритмы</b>	<b>85</b>
79	8.1. Арифметические операции с $n$ -битовыми числами . . . . .	85
80	8.2. Модульная арифметика . . . . .	89
81	<b>9. Альтернативные модели вычисления</b>	<b>93</b>
82	9.1. Работа с большими данными . . . . .	93
83	9.1.1. Модель внешней памяти . . . . .	93
84	9.1.2. Модель cache-oblivious . . . . .	95
85	9.2. Параллельные и распределённые вычисления . . . . .	96

86	9.2.1. Модель PRAM . . . . .	96
87	9.2.2. Модель BSP . . . . .	97
88	9.3. Другие модели . . . . .	101
89	<b>A. Математические факты</b>	<b>103</b>
90	A.1. Суммирование последовательностей . . . . .	103



## 91 Глава 1

### 92 Понятие алгоритма

93 Перед тем, как начать говорить об алгоритмах, нужно сначала разобраться с  
94 самим понятием алгоритма. Все мы имеем какие-то общие представления о том,  
95 что собой представляют алгоритмы обычной в жизни. Наиболее яркий и распро-  
96 странённый пример алгоритма — это кулинарные рецепты, которые описыва-  
97 ют процесс приготовления блюд, разбивая его на *базовые операции*, вроде «посо-  
98 лить», «нарезать», «вскипятить». Для того, чтобы сформулировать понятие алго-  
99 ритма нам нужно договориться о том, какие операции мы будем называть базо-  
100 выми. Другими словами, нам нужно описать то вычислительное устройство, для  
101 которого мы будем разрабатывать алгоритмы.

#### 102 1.1. Машина с произвольным доступом к памяти

103 Математическое описание вычислительного устройства называется *моделью*  
104 *вычисления*. Вы уже могли слышать про такие известные модели, как «машина  
105 Тьюринга», «лямбда-исчисление», «алгоритмы Маркова» и другие. В рамках это-  
106 го курса мы почти всё время будем работать с моделью вычисления *машина с про-*  
107 *извольным доступом к памяти* (*random access memory model, RAM-машина*). Если не  
108 сильно увлекаться формализацией, то можно считать, что речь идёт о компьюте-  
109 ре с одноядерным процессором и прямым доступом к оперативной памяти. Для  
110 нашего удобства и простоты мы даже откажемся от устройства ввода-вывода и  
111 будем предполагать, что входные данные изначально записаны в оперативной  
112 памяти, и решение задачи тоже нужно записать в заранее определённое место в  
113 памяти. Оперативная память в этой модели состоит из ячеек одинакового раз-  
114 мера, пронумерованных целыми числами. Слова «произвольный доступ» в на-  
115 звании этой модели означают, что процессор может обратиться к любой ячейке  
116 оперативной памяти по её номеру, причём это всегда делается за одну операцию  
117 (один такт процессора), т.е. время доступа к ячейке не зависит от её номера.<sup>1</sup>

118 Одна ячейки оперативной памяти современного компьютера обычно имеет  
119 размер 32 или 64 битов. Так как большинство результатов в теории алгоритмов  
120 являются асимптотическими, то нам будет удобно считать, что размер ячейки  
121 не является фиксированным числом, а зависит от некоторого параметра  $n$ , соот-  
122 ветствующего размеру решаемой задачи (обычно в качестве  $n$  будет выступать

---

<sup>1</sup>В современных процессорах всё устроено несколько сложнее. Для оптимизации доступа к оперативной памяти в процессоре имеется несколько уровней кешей. Обращение к оперативной памяти происходит не напрямую, а через кеш: для того, чтобы обратиться к ячейке, которой в кеше нет, её нужно сначала загрузить в кеш. Поэтому доступ к ячейкам, которые в данный момент присутствуют в кеше, происходит значительно быстрее, чем к тем, которых в данный момент в кеше нет.

размер входных данных). Это может показаться странным, т.к. получается, что модель вычисления зависит от задачи, которую мы решаем. На это лучше смотреть следующим образом: если параметр задачи равен  $n$ , то нам гарантируется, что ячейка памяти имеет размер  $w \geq \lceil \log_2(n) \rceil$  битов. Такого количества битов будет достаточно, чтобы представить любое целое число от 0 до  $n$ .<sup>2</sup>

Следующий набор операций в этой модели мы будем называть *базовыми* и считать, что каждая из них выполняется 1 такт процессора.

1. Обращение к ячейке памяти по её номеру (чтение или запись).
2. Логические операции с булевыми значениями.
3. Управляющие операции: условный переход, вызов функции, возврат из функции, прерывание из цикла и т.п.
4. Арифметические операции с  $w$ -битовыми целыми числами: сравнение, сложение, вычитание, умножение, деление и взятие остатка по целому модулю.
5. Битовые операции  $w$ -битовыми целыми числами: инверсия всех битов, побитовые AND, OR и XOR, а также битовые сдвиги влево и вправо на любое число  $k \leq w$ .
6. Операции с вещественными числами с ограниченной точностью: сравнение, сложение, вычитание, умножение, деление, а также приведение (округление) к целому числу.<sup>3</sup>

В предложенной модели мы будем называть *указателем* переменную, хранящую номер ячейки памяти. Кроме того, мы договоримся, что есть специальное значение для указателя (нулевой указатель), которое не соответствует никакой ячейке памяти.

## 1.2. Алгоритмический язык

Для записи алгоритмов нам потребуется выбрать некоторый довольно богатый язык программирования. В данном курсе по причине простоты синтаксиса в качестве такого языка выбрано некоторое подмножество языка Python, которое содержит базовые синтаксические конструкции и работу с массивами

**Важная оговорка.** Из соображений простоты синтаксиса мы будем предполагать, что встроенные списки языка Python на самом деле являются массивами, т.е. следующий код определяет массив из 10 элементов.

```
primes = [2, 3, 5, 7, 11, 13, 17, 19, 23, 29]
```

<sup>2</sup>В литературе такую версию машины с произвольным доступом называют также называют “word RAM” (word random access machine).

<sup>3</sup>Тут нужно быть аккуратным и сразу договориться, что на самом деле мы будем работать не с вещественными числами, а с их приближениями с ограниченной точностью, например, представляя вещественные числа как числа с плавающей точкой или числами с фиксированной точностью. Если разрешить работать с настоящими вещественными числами, то получившаяся модель вычисления будет позволять, например, сохранить данные неограниченного размера в одной ячейке памяти, представив их как дробную часть вещественного числа.

155 Мы будем считать, что `primes` является “классическим” массивом, т.е. элементы  
156 находятся в памяти последовательно друг за другом и, соответственно, по номеру  
157 элемента массива легко вычислить номер ячейки памяти, где хранится этот эле-  
158 мент. На практике интерпретаторы Python этого не гарантируют. Поэтому, если в  
159 реальной программе требуется работа с “настоящими” массивами, то приходит-  
160 ся использовать внешние библиотеки, например, библиотеку `numpy`.

### 161 1.3. Определение алгоритма

162 **Определение 1.3.1.** *Алгоритмом* для решения задачи  $P$  будем называть програм-  
163 му, записанную на алгоритмическом языке, которая получает на вход условие  
164 задачи  $P$  и вычисляет ответ для этой задачи. Мы будем говорить, что алгоритм  
165 *корректен*, если для любого условия задачи  $P$  алгоритм находит верное решение.

166 Так как в описании модели вычисления мы отказались от устройства ввода-  
167 вывода, то все алгоритмы мы будем записывать как функции, которые получают  
168 входные данные в качестве аргументов и возвращают ответ.

169 Для каждого алгоритма, рассматриваемого в этом курсе, мы будем пытаться  
170 отвечать на следующие три вопроса.

- 171 1. Корректен ли этот алгоритм?
- 172 2. Какова сложность этого алгоритма?
- 173 3. Можно ли решить эту задачу быстрее?

174 В тех случаях, когда ответ на третий вопрос будет отрицательным, мы будем го-  
175 ворить, что рассматриваемый алгоритм является «оптимальным».

176 **Пример: проверка наличия элемента в массиве.** Давайте разработаем алго-  
177 ритм, который проверяет, встречается ли заданный элемент в некотором масси-  
178 ве. Договоримся, что алгоритм получает искомым элемент `el` и массив `arr`, как  
179 аргументы функции. Наиболее естественный подход к решению этой задачи —  
180 последовательно проверять все элементы массива `arr` на равенство `el` до перво-  
181 го совпадения. Эта идея реализуется функцией `has_element`.

```
182 # проверка наличия элемента el в массиве arr  
183 def has_element(arr, el):  
184     for x in arr:  
185         if x == el:  
186             return True  
187  
188     return False
```

189 Будем предполагать, что массив `arr` и переменная `el` содержат целые числа, кро-  
190 ме того мы обозначим длину массива `arr` через  $n$ . Попробуем оценить количество  
191 базовых операций, которое совершает данный алгоритм. Количество итераций  
192 в цикле не превосходит  $n$ . На каждой итерации происходит обращение к двум  
193 ячейкам памяти, сравнение, условный переход, увеличение счётчика для перехо-  
194 да к следующей ячейке или возвращение значения. Таким образом, общее число  
195 операций не превосходит  $5n + 1$  (единичка добавляется для возврата значения в  
196 случае, если элемент в массиве отсутствует).

197 Такой точный анализ числа базовых на практике имеет очень мало смысла.  
 198 Дело в том, что при реализации алгоритма на настоящем компьютере реальное  
 199 количество тактов процессора на входе размера  $n$  будет зависеть от множества  
 200 дополнительных факторов, таких как архитектура процессора, язык программи-  
 201 рования, версия и параметры компилятора и/или интерпретатора, используемой  
 202 операционной системы и т.д. Поэтому нас интересует какой-то способ сравни-  
 203 вать эффективность алгоритмов, который бы не зависел от подобных факторов.  
 204 Поэтому вместо точной оценки количества операций мы будем оценивать *ско-*  
 205 *рость роста* количества операций в зависимости от размера входных данных.  
 206 В данном случае в качестве размера входных данных выступает число  $n$ , поэто-  
 207 му количество операций *линейно* возрастает при увеличении  $n$ . В таких случаях  
 208 мы будем говорить, что мы имеем дело с *линейным алгоритмом* и обозначать его  
 209 сложность  $O(n)$ .

210 Перед тем, как начать разбираться с формальными определениями для  $O$ -  
 211 *символики*, нужно отметить ещё один важный момент. При оценке сложности  
 212 функции `has_element` мы оценивали количество операций в самом *худшем слу-*  
 213 *чае*, т.е. тогда, когда в массиве `arr` нет элемента `el`, и функции придётся прове-  
 214 рить весь массив, чтобы убедиться в этом. В дальнейшем мы всегда, если не ука-  
 215 зано обратное, под *сложностью алгоритма* будем понимать количество операций  
 216 в худшем случае.

## 217 1.4. $O$ -символика

218 В следующих определениях функции  $f, g : \mathbb{N} \rightarrow \mathbb{R}_{>0}$  (т.е. функции, )

219 **Определение 1.4.1.** Будем говорить, что функция  $f$  *растёт не быстрее*  $g$  и обо-  
 220 значать  $f(n) = O(g(n))$  или  $f \preceq g$ , если существует такая константа  $c > 0$ , что  
 221  $f(n) \leq c \cdot g(n)$  для всех натуральных  $n$ .

222 *Пример 1.4.1.*  $3n^2 + 4n + 5 = O(n^2)$ . Для того, чтобы это показать, нам нужно  
 223 предъявить подходящую константу  $c$ . В данном случае достаточно выбрать  $c =$   
 224  $12$ , поскольку при всех  $n \geq 1$  выполнено  $3n^2 + 4n + 5 \leq 3n^2 + 4n^2 + 5n^2 = 12n^2$ .

### 225 Свойства

226 1. Для любых константы  $a > 0$  верно  $a \cdot f(n) = O(f(n))$ , т.е. можно отбрасывать  
 227 константные множители.

228 *Пример 1.4.2.*  $10n + 1500 = O(n)$ .

229 2. Для любой константы  $a \geq 0$  верно  $a = O(1)$ .

230 *Пример 1.4.3.*  $10^{100} = O(1)$ .

231 3. Для любых  $a \leq b$  выполняется  $n^a = O(n^b)$ .

232 *Пример 1.4.4.*  $n^2 = O(n^3)$ .

233 4. Все логарифмы с основаниями больше единицы растут одинаково быстро:  
 234 для любых  $a, b > 1$  выполняется  $\log_a n = O(\log_b n)$ . Поэтому, основание ло-  
 235 гарифма в  $O$ -символике обычно не указывают.

236 *Пример 1.4.5.*  $\log_{16} n = O(\log n)$ .



237 5. Любой полилогарифм растёт медленнее любого полинома: для любых  $a, b >$   
 238  $0$  выполняется  $\log^a n = O(n^b)$ .

239 *Пример 1.4.6.*  $\log^{100} n = O(n^{0.01})$ .

240 6. Любой полином растёт медленнее любой экспоненты: для любых  $a > 0$  и  
 241  $b > 1$  выполняется  $n^a = O(b^n)$ .

242 *Пример 1.4.7.*  $n^{100} = O(1.01^n)$ .

243 7. Сумма функций растёт не быстрее, чем самое большое слагаемое: если  $f(n) =$   
 244  $O(g(n))$ , то  $f(n) + g(n) = O(g(n))$ .

245 *Пример 1.4.8.*  $18n^3 + 12n^2 + 5n + 100 + 35 \log_2(n) = O(n^3)$ .

246 *Упражнение 1.4.1.* Приведите формальные доказательства для каждого из этих  
 247 свойств.

248 *Замечание 1.4.1.* Обозначение  $f(n) = O(g(n))$  не совсем корректно с точки зрения  
 249 математики. Дело в том, что знак равенства в таком выражении не обозначает  
 250 равенство в его обычном смысле. Например, из того, что  $2n = O(n)$  и  $3n = O(n)$   
 251 не следует, что  $2n = 3n$ . Для того, чтобы разобраться, почему так происходит,  
 252 нужно задаться вопросом, а какие математические объекты стоят по обеим сто-  
 253 ронам «равенства»  $f(n) = O(g(n))$ ? Следуя определению слева стоит функция  
 254  $f: \mathbb{N} \rightarrow \mathbb{R}_{>0}$ . А какой математический объект скрывается за  $O(g(n))$ ? Под  $O(g(n))$   
 255 понимается класс всех функций, которые растут не быстрее  $g$ . Таким образом,  
 256 математически корректно было бы написать  $f(n) \in O(g(n))$ , т.е. функция  $f$  при-  
 257 надлежит классу (множеству) всех функций, которые растут не быстрее  $g$ . Однако  
 258 в теории алгоритмов (и не только в ней) принято писать  $f(n) = O(g(n))$  и мы то-  
 259 же будем этому следовать.

260 **Преимущества O-символики.** Использование O-символики позволяет охарак-  
 261 теризовать зависимость времени работы алгоритма от размера входных. Получа-  
 262 емые оценки устойчивы к изменению таких факторов, как язык программирова-  
 263 ния или параметры компиляции. Кроме того, сами утверждения про сложность  
 264 алгоритмов носят более простой характер (сложность алгоритма оценивается как  
 265  $O(n^2)$  вместо  $3n^2 + 4n + 5$ ), и значительно упрощается анализ (не нужно скрупу-  
 266 лёзно подсчитывать количество операций, достаточно выделить наиболее тру-  
 267 доёмкие).

268 **Недостатки O-символики.** При использовании O-символики алгоритмы с оди-  
 269 наковыми оценками сложности на практике могут иметь значительно различа-  
 270 ющуюся производительность. Представим себе два алгоритма, имеющие слож-  
 271 ность  $O(n)$ , такие что первый делает порядка  $2n$  операций, а второй —  $2000n$  опе-  
 272 раций. В терминах O-символики эти алгоритмы имеют одинаковую сложность,  
 273 но на практике второй алгоритм значительно проигрывает первому. Более то-  
 274 го, нередко бывает, что, например, для некоторой задачи существует линейный  
 275 алгоритм, но на практике используются простой алгоритм с оценкой  $O(n \log n)$ ,  
 276 который имеет лучшую производительность (это объясняется тем, что констан-  
 277 та скрытая в  $O(n)$  значительно больше, чем  $\log n$  для характерных значений  $n$ ,  
 278 встречающихся в реальных задачах).

279 **Связанные определение.** Введённое выше определение  $O$ -символики позво-  
 280 ляет сравнивать две функции на «меньше или равно». В некоторых случаях нам  
 281 потребуются следующие связанные определения, которые задают все оставшие-  
 282 ся сравнения между функциями.

283 **Определение 1.4.2.** Будем говорить, что функция  $f$  растёт не медленнее  $g$ , и обо-  
 284 значать  $f(n) = \Omega(g(n))$  или  $f \succeq g$ , если существует такая константа  $c > 0$ , что  
 285  $f(n) \geq c \cdot g(n)$  для всех натуральных  $n$ .

286 **Определение 1.4.3.** Будем говорить, что функция  $f$  растёт так же как  $g$ , и обо-  
 287 значать  $f(n) = \Theta(g(n))$  или  $f \asymp g$ , если существуют две константы  $c_1, c_2 > 0$ , такие  
 288 что  $c_1 \cdot g(n) \leq f(n) \leq c_2 \cdot g(n)$  для всех натуральных  $n$ . Другими словами,  $f(n) =$   
 289  $\Theta(g(n))$ , если одновременно выполняется  $f(n) = O(g(n))$  и  $f(n) = \Omega(g(n))$ .

290 Следующие два определения задают строгие отношения, соответствующие «мень-  
 291 ше» и «больше».

292 **Определение 1.4.4.** Будем говорить, что функция  $f$  растёт медленнее  $g$  и обо-  
 293 значать  $f(n) = o(g(n))$ , или  $f \prec g$ , если  $\lim_{n \rightarrow \infty} \frac{f(n)}{g(n)} = 0$ .

294 **Определение 1.4.5.** Будем говорить, что функция  $f$  растёт быстрее  $g$  и обозна-  
 295 чать  $f(n) = \omega(g(n))$ , или  $f \succ g$ , если  $\lim_{n \rightarrow \infty} \frac{g(n)}{f(n)} = 0$ , т.е.  $g(n) = o(f(n))$ .

296 **Применимость такого анализа на практике.**

## 297 1.5. Числа Фибоначчи

298 **Определение 1.5.1.** Числа Фибоначчи — это бесконечная числовая последова-  
 299 тельность целых чисел  $F_0, F_1, F_2, F_3, \dots$ , определённая по следующим правилам:

$$300 \quad F_0 = 0, \quad F_1 = 1, \quad F_n = F_{n-1} + F_{n-2}.$$

301 Таким образом  $F_2 = F_1 + F_0 = 1$ ,  $F_3 = F_2 + F_1 = 2$ ,  $F_4 = F_3 + F_2 = 3$  и т.д.

302 Можно выписать начало этой последовательности:

$$303 \quad 0, 1, 1, 2, 3, 5, 8, 13, 21, 34, 55, 89, 144, \dots$$

304 Насколько быстро члены этой последовательности возрастают? Следующие две  
 305 леммы показывают, что числа Фибоначчи возрастают экспоненциально.

306 **Лемма 1.5.1.** Для любого  $n \geq 6$ ,  $F_n \geq 2^{n/2}$ .

307 *Доказательство.* Докажем индукцией по  $n$ .

308 1. **База.**  $F_6 = 8 \geq 2^3 = 2^{6/2}$ ,  $F_7 = 13 \geq 12 = 2^3 \cdot 1.5 > 2^3 \cdot \sqrt{2} = 2^{7/2}$ .

309 2. **Предположение.** Предположим, что  $F_n \geq 2^{n/2}$  для всех  $n < k$  для некоторого  
 310  $k \geq 2$ .

311 3. **Переход.** Докажем, что  $F_k \geq 2^{k/2}$ . По предположению  $F_{k-1} \geq 2^{(k-1)/2}$  и  $F_{k-2} \geq$   
 312  $2^{(k-2)/2}$ , следовательно

$$313 \quad F_k = F_{k-1} + F_{k-2} \geq 2^{(k-1)/2} + 2^{(k-2)/2} \geq 2 \cdot 2^{(k-2)/2} = 2^{k/2}.$$

314

□

315 **Лемма 1.5.2.** Для любого  $n \geq 0$ ,  $F_n \leq 2^n$ .

316 *Доказательство.* Докажем индукцией по  $n$ .

317 1. **База.**  $F_0 = 0 \leq 2^0$ ,  $F_1 = 1 \leq 2^1$ .

318 2. **Предположение.** Предположим, что  $F_n \leq 2^n$  для всех  $n < k$  для некоторого  
319  $k \geq 2$ .

320 3. **Переход.** Докажем, что  $F_k \leq 2^k$ . По предположению индукции  $F_{k-1} \leq 2^{k-1}$   
321 и  $F_{k-2} \leq 2^{k-2}$ , следовательно,

$$322 \quad F_k = F_{k-1} + F_{k-2} \leq 2^{k-1} + 2^{k-2} \leq 2 \cdot 2^{k-1} = 2^k.$$

323

□

324 **Следствие 1.5.1.** Для любого  $n \geq 6$  выполняется

$$325 \quad 2^{n/2} \leq F_n \leq 2^n.$$

326 Числа Фибоначчи очень хорошо изучены, в том числе известна и точная оцен-  
327 ка  $F_n = \Theta(\varphi^n)$ , где  $\varphi = \frac{1+\sqrt{5}}{2}$  — золотое сечение. Эта оценка следует из формулы  
328 Бине для вычисления чисел Фибоначчи:

$$329 \quad F_n = \frac{\left(\frac{1+\sqrt{5}}{2}\right)^n - \left(\frac{1-\sqrt{5}}{2}\right)^n}{\sqrt{5}} = \frac{\varphi^n - (-\varphi)^{-n}}{\varphi - (-\varphi)^{-1}} = \frac{\varphi^n - (-\varphi)^{-n}}{2\varphi - 1}.$$

330 Доказывать эту формулу не входит в наши цели. Более того, если мы захотим ис-  
331 пользовать эту формулу для вычисления чисел Фибоначчи, то нам потребуются  
332 работать с числами с плавающей точкой, а следовательно придётся как-то оце-  
333 нивать погрешность. В результате окажется, что использовать данную формулу  
334 можно только для относительно небольших  $n$ .

335 Вместо это давайте напишем простую функцию, которая вычисляет  $F_n$  по опре-  
336 делению. Так как числа Фибоначчи определяются через рекуррентно, то есте-  
337 ственно вычислять числа их при помощи рекурсии.

338 *# рекурсивная реализация вычисления чисел Фибоначчи*

```
339 def fib_rec(n):
340     if n == 0:
341         return 0
342
343     if n == 1:
344         return 1
345
346     return fib_rec(n-1) + fib_rec(n-2)
```

347 Структура этой функции дословно повторяет определение чисел Фибоначчи, по-  
348 этому более-менее понятно, что алгоритм корректен. Тем не менее, давайте при-  
349 ведём *формальное доказательство корректности*. В данном случае корректность  
350 алгоритма означает, что для целого положительного  $n$  функция  $\text{fib\_rec}(n)$  воз-  
351 вращает  $F_n$ .

352 **Лемма 1.5.3.** Для любого целого  $n \geq 0$  функция  $\text{fib\_rec}(n)$  возвращает  $F_n$ .

353 *Доказательство.* Докажем индукцией по  $n$ .

354 1. **База.** Для  $n = 0$  и  $n = 1$  функция возвращает соответственно  $F_0$  и  $F_1$ .

355 2. **Предположение.** Предположим, что для всех  $n$  меньших некоторого  $k \geq 2$   
356 функция  $\text{fib\_rec}(n)$  возвращает  $F_n$ .

357 3. **Переход.** Докажем, что  $\text{fib\_rec}(k)$  возвращает  $F_k$ . По предположению  $\text{fib\_rec}(k-1)$   
358 и  $\text{fib\_rec}(k-2)$  возвращают соответственно  $F_{k-1}$  и  $F_{k-2}$ . Следовательно  $\text{fib\_rec}(k)$   
359 возвращает  $F_{k-1} + F_{k-2} = F_k$ .

360 □

361 Разобравшись с корректностью этого алгоритма мы можем попробовать оце-  
362 нить время его работы. Пусть  $T(n)$  — это количество операций, которые которые  
363 требуются для вычисления  $\text{fib\_rec}(n)$ . Тогда  $T(0) = 1$  (одно сравнение),  $T(1) = 2$   
364 (два сравнения), а  $T(n) = T(n-1) + T(n-2) + 3$  (два сравнения и сложение). По-  
365 лученное соотношение очень похоже на определение чисел Фибоначчи. Покажем  
366 по индукции, что  $T(n) \geq F_n$  для любого  $n$ .

367 1. **База.**  $T(0) \geq F_0$ ,  $T(1) \geq F_1$ .

368 2. **Предположение.** Предположим, что это верно для всех  $n$  меньше некото-  
369 рого  $k \geq 2$ .

370 3. **Переход.** Докажем, что  $T(k) \geq F_k$ . По предположению  $T(k-1) \geq F_{k-1}$  и  
371  $T(k-2) \geq F_{k-2}$ , следовательно

$$372 \quad T(k) = T(k-1) + T(k-2) + 3 \geq F_{k-1} + F_{k-2} + 3 > F_k.$$

373 Таким образом получается, что количество операций, которые нужны функции  
374  $\text{fib\_rec}$  для вычисления  $F_n$  не меньше  $F_n$ , т.е.  $T(n) = \Omega(2^{n/2})$ . Это время работы  
375 совершенно неприемлемо, т.к. для вычисления всего лишь  $F_{100}$  нам потребуется  
376 более  $2^{50} > 10^{15}$  операций (более 10 дней на процессоре с частотой 1GHz).

377 *Замечание 1.5.1.* Мы здесь оценили время работы предложенного алгоритма *сни-*  
378 *зу.* Можно аналогичным способом по индукции доказать, что  $T(n) \leq 2F_{n+1}$  и та-  
379 *ким* образом получить верхнюю оценку  $T(n) = O(2^n)$ , но для данного алгоритма  
380 это не очень интересно, т.к. мы уже знаем, что он очень медленный.

381 Для того, чтобы разобраться, почему этот естественный алгоритм работает  
382 так медленно, давайте нарисуем схему рекурсивных вызовов функции  $\text{fib\_rec}$ .

383 На этой схеме видно, что для вычисления  $F_n$  мы сначала рекурсивно вычис-  
384 ляем  $F_{n-1}$  и  $F_{n-2}$ . При этом для вычисления  $F_{n-1}$  мы снова вычисляем  $F_{n-2}$ . Здесь  
385 и кроется проблема этого алгоритма: мы не запоминаем, что уже однажды вы-  
386 числили  $F_{n-2}$ , и вычисляем его заново. На схеме видно, что  $F_{n-3}$  вычисляется уже  
387 трижды,  $F_{n-4}$  — пять раз, и т.д. Из этого анализа возникает естественная идея,  
388 что мы должны запоминать те числа Фибоначчи, которые мы уже вычислили.  
389 Это приводит нас к следующему алгоритму вычисления чисел Фибоначчи с со-  
390 хранением результатов в массиве.

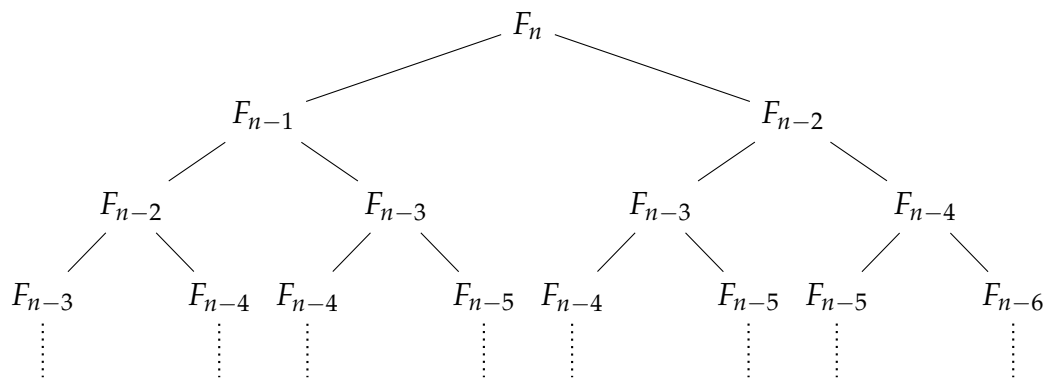


Рис. 1.5.1. Схема рекурсивных вызовов fib\_rec(n).

```

391 # вычисления числа Фибоначчи с сохранением результатов в массиве
392 def fib_array(n):
393     # изначально в массиве хранятся F0 и F1
394     a = [0, 1]
395     for i in range(2, n+1):
396         # добавляем в массив Fi = Fi-1 + Fi-2
397         a.append(a[i-1] + a[i-2])
398
399     # возвращаем Fn
400     return a[n]
  
```

Если более детально проанализировать алгоритм fib\_array, то можно заметить, что на каждой итерации мы пользуемся только двумя самыми последними элементами массива. Это приводит нас к мысли, что достаточно сохранить только два последних элемента. В этом случае вместо массива мы заведём две дополнительные переменные.

```

406 # вычисления числа Фибоначчи с двумя дополнительными переменными
407 def fib_vars(n):
408     if n == 0:
409         return 0
410
411     # изначально переменные хранят F0 и F1
412     a, b = [0, 1]
413     for i in range(2, n+1):
414         a, b = b, b + a
415
416     return b
  
```

Полученный алгоритм так же работает за  $O(n)$ , но в отличие от fib\_array использует константное количество ячеек памяти.

*Замечание 1.5.2.* Вычисление можно ещё ускорить, если записать вычисление числа Фибоначчи на языке матриц и воспользоваться алгоритмом быстрого возведения в степень. В результате можно получить алгоритм со сложностью  $O(\log n)$  использующий  $O(1)$  ячеек памяти.



## 423 Глава 2

# 424 Структуры данных

425 Для эффективной работы с данными используются различные способы орга-  
426 низации их хранения в памяти — *структуры данных*. В данной главе мы погово-  
427 рим про самые базовые структуры данных, а так же научимся оценивать амор-  
428 тизированную сложность *запросов* к ним<sup>1</sup>.

### 429 2.1. Массивы

430 Наиболее простой структурой данных для RAM-машины является массив. Мас-  
431 сив представляется собой *непрерывную* область памяти, разбитую на ячейки оди-  
432 накового размера. Зная адрес начала массива можно легко вычислить адрес про-  
433 извольной ячейки. Языки программирования обычно уже имеют готовые функ-  
434 ции для работы с массивами, но нам было бы полезно разобраться, как они устро-  
435 ены.

436 Если нам изначально известно, сколько элементов будет храниться в массиве,  
437 или же у нас есть некоторая оценка сверху на количество элементов, то мы можем  
438 сразу же *выделить* массив достаточного размера (т.е. договориться, что некото-  
439 рая непрерывная область памяти будет использоваться для хранения элементов  
440 массива). Для добавления элемента в конец заранее выделенного массива доста-  
441 точно просто скопировать его в ещё незаполненную ячейку и увеличить счётчик,  
442 который хранит количество заполненных ячеек. Таким образом, сложность за-  
443 проса на добавление элемента —  $O(1)$  операций. Запрос на удаление последнего  
444 элемента реализовать ещё проще, т.к. нам достаточно просто уменьшить счётчик  
445 числа заполненных ячеек.

446 *Замечание 2.1.1.* Здесь и далее при оценке сложности запросов к структурам дан-  
447 ных мы всегда будем предполагать, что размер элементов структуры данных фик-  
448 сирован и не зависит от размера входа. Другими словами, мы будем предпола-  
449 гать, что сложность копирования любого элемента структуры данных констант-  
450 на, т.е. стоит нам  $O(1)$  операций. Если это не так, то в большинстве случаев по-  
451 лученные оценки сложности легко скорректировать, домножив их на сложность  
452 копирования одного элемента. Например, если в массиве хранятся строки дли-  
453 ны  $\sqrt{n}$ , то оценки сложности всех запросов, которые приводят к копированию  
454 элементов, нужно умножить на  $\sqrt{n}$ : добавление элемента в конец будет иметь  
455 сложность  $O(\sqrt{n})$ , а удаление последнего элемента —  $O(1)$ , т.к. не требует копи-  
456 рования элементов.

---

<sup>1</sup>Для того, чтобы избежать путаницы мы будем говорить об *операциях* процессора и *запросах* к структурам данных.

457 Если мы захотим удалить или добавить элемент в начале или где-то в сере-  
 458 дине массива, то нам потребуется “сдвинуть” все последующие элементы, поэто-  
 459 му эти запросы будут иметь сложность  $O(n)$ , где  $n$  — это количество элементов в  
 460 массиве.

461 Однако, мы далеко не всегда заранее знаем, сколько элементов будет храниться  
 462 в массиве. В этом случае можно начать с небольшого массива и расширять его по-  
 463 степенно, добавляя новые ячейки тогда, когда они потребуются. Предположим,  
 464 что все ячейки массива заполнены, а нам нужно добавить к массиву ещё один эле-  
 465 мент? Вполне возможно, что ячейки памяти, которые расположены сразу же по-  
 466 сле конца массива, свободны. Тогда мы могли бы расширить массив за счёт при-  
 467 соединения к нему необходимого количества ячеек памяти (например, это может  
 468 сделать функция `realloc` в языке C). Но что делать, если ячейки памяти после  
 469 конца массива уже заняты или наш язык программирования не поддерживает  
 470 расширение массивов? В этом случае вместо расширения существующего масси-  
 471 ва нам придётся создать новый массив большего размера и скопировать туда все  
 472 данные из исходного массива (массив обязательно должен лежать в непрерыв-  
 473 ной области памяти, поэтому мы не можем расположить дополнительные ячейки  
 474 где-то отдельно). Давайте рассмотрим две различные реализации этой идеи.

### 475 2.1.1. Расширяемый массив с аддитивной схемой выделения

476 Предположим, что нам нужно добавить один новый элемент в конец масси-  
 477 ва размера  $s$ , который уже полностью заполнен, т.е. хранит  $s$  элементов. В этом  
 478 случае мы можем выделить новый массив размера  $s + c$ , где  $c$  — некоторая целая  
 479 положительная константа, скопировать туда все элементы из исходного массива  
 480 и добавить новый элемент в первую из добавленных пустых ячеек (см. рис. 2.1.1).  
 Это называют *расширяемым массивом с аддитивной схемой выделения*. Какую слож-

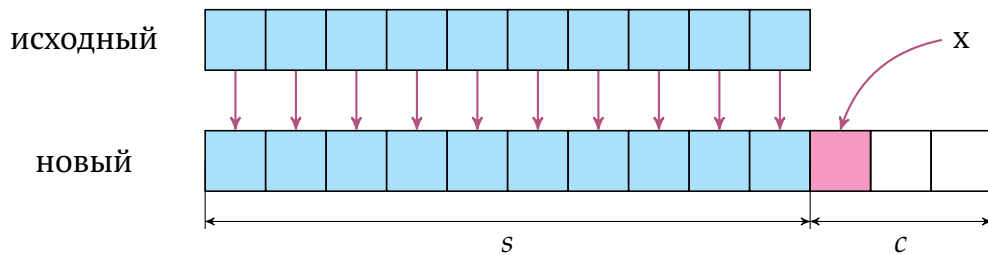


Рис. 2.1.1. Добавление элемента  $x$  в заполненный расширяемый массив с аддитивной схемой выделения при  $c = 3$ .

481 ность будет иметь запрос на добавление элемента в при такой реализации рас-  
 482 ширения? Пусть в массиве  $n$  элементов. Если свободных ячеек нет, то при добав-  
 483 лении нового элемента нам потребуется скопировать  $n$  элементов из исходного  
 484 массива и один новый элемент. Таким образом запрос на добавление элемента  
 485 будет иметь сложность  $\Theta(n)$ .

487 Полезно также оценить суммарные накладные расходы на расширение мас-  
 488 сива. Пусть мы начали с пустого массива, и в процессе работы он увеличился  
 489 до размера  $n$ . Сколько дополнительных операций потребовалось? Для простоты  
 490 будем учитывать только операции копирования элементов. При первом нетри-  
 491 виальном расширении массива было скопировано  $c$  элементов, при втором —  $2c$   
 492 элементов, и так далее. В сумме получается, что на все расширения массива по-



493 требовалось

$$494 \quad c + 2c + 3c + \dots + (n - c) = \frac{n(n/c - 1)}{2} = \Theta(n^2)$$

495 дополнительных копирований (см. лемму A.1.1). Получается, что при такой ре-  
496 ализации расширяемого массива любой алгоритм, который использует массив  
497 размера  $n$ , имеет сложность  $\Omega(n^2)$ .

### 498 2.1.2. Расширяемый массив с мультипликативной схемой выделения.

499 Вместо увеличения размера массива на константу, можно увеличивать размер  
500 массива в константу раз. Эта идея называется *расширяемым массивом с мульти-*  
501 *пликативной схемой выделения*. Выберем некоторое  $\alpha > 1$  (наиболее типичным  
502 значением является  $\alpha = 2$ ). Если в текущем массиве размера  $s$  все ячейки запол-  
нены, то выделим новый массив размера  $\lceil \alpha s \rceil$  (см. рис. 2.1.2).

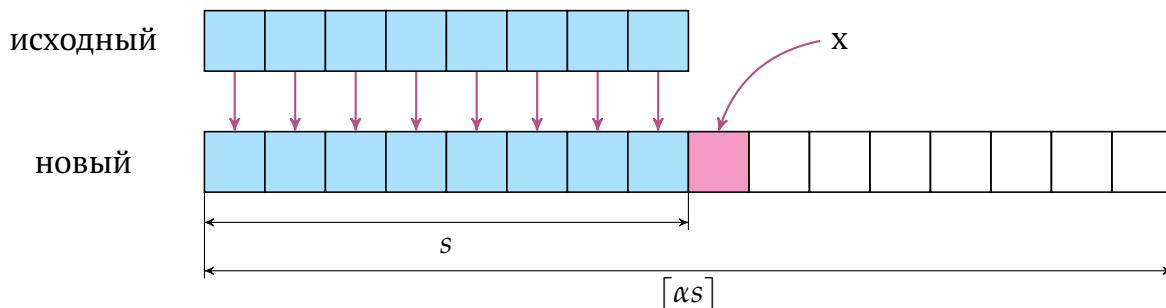


Рис. 2.1.2. Добавление элемента  $x$  в заполненный расширяемый массив с мультипликативной схемой выделения при  $\alpha = 2$ .

503 Давайте оценим также суммарные накладные расходы на расширение массива.  
504 Пусть мы начали с массива некоторого фиксированного размера, и в процессе  
505 работы размер увеличился до  $n$ . При последнем расширении массива было ско-  
506 пировано  $\lfloor n/\alpha \rfloor$  элементов, при предпоследнем —  $\lfloor n/\alpha^2 \rfloor$ , и так далее. В сумме на  
507 расширения массива потребовалось  $\lfloor n/\alpha \rfloor + \lfloor n/\alpha^2 \rfloor + \dots + c$  дополнительных ко-  
508 пирований, где  $c$  — начальный размер массива. Эту сумму можно оценить сверху  
509 суммой бесконечной убывающей геометрической прогрессии (см. лемму A.1.3)

$$511 \quad \lfloor n/\alpha \rfloor + \lfloor n/\alpha^2 \rfloor + \dots + c \leq \sum_{k=1}^{\infty} \lfloor n/\alpha^k \rfloor < \sum_{k=1}^{\infty} n/\alpha^k = n \cdot \frac{1}{1 - 1/\alpha} = \Theta(n).$$

512 Таким образом при мультипликативной схеме выделения накладные расходы  
513 линейные. Это можно переформулировать следующим образом: в среднем на  
514 каждый элемент массива требуется  $O(1)$  дополнительных операций на расши-  
515 рение массива. Стоит отметить, что не смотря на это, запрос на добавление в  
516 массив будет иметь сложность  $\Theta(n)$  в худшем (сложность в худшем оценивает  
517 количество операций в самом “плохом” случае).

518 *Упражнение 2.1.1.* Как реализовать запрос на удаление элемента из массива так,  
519 чтобы в среднем на каждый элемента массива дополнительно требовалось  $O(1)$   
520 операций, и при этом свободное место использовалось бы экономно (т.е., напри-  
521 мер, если из большого массива удалить почти все элементы, то зарезервирован-  
522 ное свободное место должно пропорционально уменьшиться)?

## 2.2. Списки

Другая базовая структура данных — это список. В отличие от массива, где все данные хранятся в непрерывной области памяти, элементы списка могут быть разбросаны по памяти произвольным образом. Выделяют два основных вида списков: *односвязный* и *двусвязный*.

### 2.2.1. Односвязный список

Каждый элемент *односвязного списка* хранится в отдельной структуре с двумя полями — *узле списка*. Кроме самого элемента каждый узел хранит *указатель* на следующий узел — номер ячейки памяти, где находятся данные следующего узла списка (см. рис. 2.2.1). Односвязный список задаётся парой указателей: на первый узел списка, *голову списка*, и на последний узел списка, *хвост списка*.

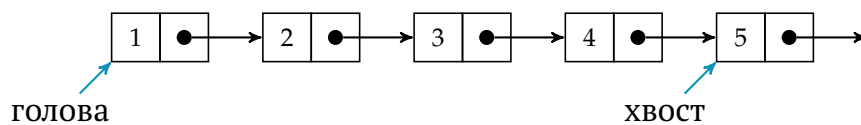


Рис. 2.2.1. Односвязный список.

Такая организация хранения данных позволяет перебирать элементы от головы к хвосту списка, причём переход к следующему элементу выполняется за  $O(1)$ , удалять и добавлять элементы, находящиеся в начале списка, за  $O(1)$ , добавлять элементы в конец списка за  $O(1)$ , а так же удалять и добавлять элементы в произвольном месте списка за  $O(1)$  по ссылке на предыдущий элемент см. рис. 2.2.2. При этом обращение по индексу и удаление из конца списка будут работать за  $\Theta(n)$  (для удаления последнего элемента нужно найти предпоследний узел, а для этого придётся пройти весь список от головы до хвоста).

Кроме того, списки позволяют эффективно осуществлять некоторые более сложные операции. Например, можно «вырезать» часть списка и выделить её в отдельный список или наоборот «вклеить» один список внутрь другого, причём обе операции будут работать за  $O(1)$ , т.к. они требуют изменения всего лишь пары указателей.

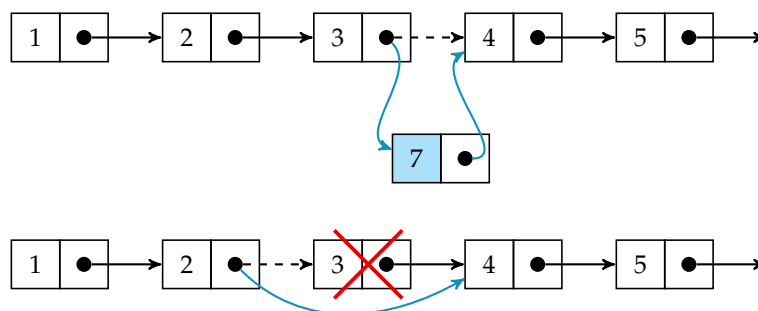


Рис. 2.2.2. Добавление и удаление элемента односвязного списка.

### 2.2.2. Двусвязный список

Если есть необходимость перебирать элементы списка в обоих направлениях, то применяется *двусвязный список*. Элементы двусвязного списка хранятся в уз-

550 лах вместе с указателями на следующий и предыдущий узлы списка (см. рис. 2.2.3).  
 551 За счёт хранения двух указателей двусвязный список позволяет добавлять и уда-  
 552 лять элементы в произвольном месте списка за  $O(1)$ . При этом доступ к элементу  
 553 по индексу так же осуществляется за  $\Theta(n)$ .

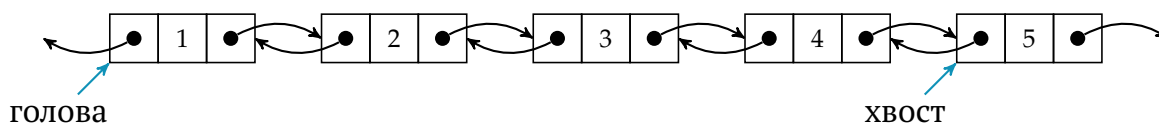


Рис. 2.2.3. Двусвязный список.

554 По функциональности двусвязный список ни в чём не уступает односвязному,  
 555 в т.ч. так же позволяет эффективно реализовать операции «вырезания» и «вкле-  
 556 ивания». Тем не менее, в тех приложениях, где память очень ограничена (напри-  
 557 мер, внутри кода драйверов или во встроенных системах), предпочитают исполь-  
 558 зовать односвязные списки, если их возможностей достаточно для решения по-  
 559 ставленной задачи.

560 *Упражнение 2.2.1.* Придумайте, как можно хранить двусвязный список так, чтобы  
 561 в каждом узле хранился только один указатель, но при этом мы могли бы пере-  
 562 бирать элементы списка в обоих направлениях. (Подсказка: нужно использовать  
 563 битовые операции с целыми числами.)

### 564 2.2.3. Сравнение с массивами

565 В следующей таблице представлены сложности основных запросов к массивам и спискам.

Запрос	Массив	Односвязный список	Двусвязный список
Обращение по номеру	$O(1)$	$O(n)$	$O(n)$
Добавление в начало	$O(n)$	$O(1)$	$O(1)$
Удаление из начала	$O(n)$	$O(1)$	$O(1)$
Добавление в конец	$O(n)$	$O(1)$	$O(1)$
Удаление из конца	$O(1)$	$O(n)$	$O(1)$
Добавление в середину	$O(n)$	$O(1)$	$O(1)$
Удаление из середины	$O(n)$	$O(1)$	$O(1)$

Таблица 2.1. Сравнение сложности запросов к массивам и спискам.

566

## 567 2.3. Абстрактные типы данных

568 **Определение 2.3.1.** *Абстрактный тип данных (АТД)* — это математическое опи-  
 569 сание структуры данных в терминах реализуемых ей запросов.

570 Понятия абстрактного типа данных и структуры данных соотносятся так же,  
 571 как в объектно-ориентированном программировании соотносятся понятия ин-  
 572 терфейса и класса. Когда мы определяем структуру данных, то мы подробно опи-  
 573 сываем то, как эта структура устроена, и какие запросы к этой структуре можно  
 574 реализовать эффективно. При определении абстрактного типа данных мы опи-  
 575 сываем только те запросы, которые этот тип данных должен реализовывать, но

576 ничего не говорим про его устройство в памяти и реализацию запросов. Таким  
 577 образом, у одного абстрактного типа данных может быть несколько реализаций,  
 578 основанных на разных структурах данных. В этом разделе мы рассмотрим два  
 579 базовых абстрактных типа данных и поговорим о том, как их можно реализовать  
 580 на основе массивов и списков.

### 581 2.3.1. Стек

582 Абстрактный тип данных *стек* определяется следующими запросами:

- 583 • `push(x)` — добавить элемент в стек,
- 584 • `pop()` — вытолкнуть элемент стека.

585 Порядок, в котором элементы выталкиваются из стека, является обратным по-  
 586 рядку добавления, т.е. первым выталкивается тот элемент, который был добав-  
 587 лен последним (см. рис. 2.3.1). Про стек говорят, что он реализует концепцию  
 588 «последним пришёл — первым ушёл» («last in, first out», LIFO). В обычной жиз-  
 589 ни такой порядок можно встретить, например, если положить несколько книг в  
 590 стопку на стол: сверху будет лежать та книга, которую мы положили последней.  
 591 В соответствии с этой аналогией про элемент, который добавлен последним, го-  
 592 ворят, что он находится *на вершине* стека, а про тот, что был добавлен первым —  
 что он находится *на дне* стека.

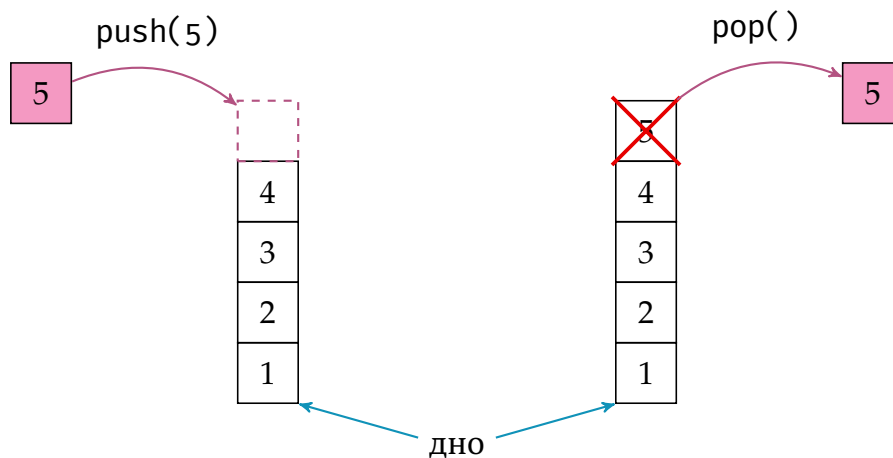


Рис. 2.3.1. Операции со стеком: добавление и выталкивание элемента.

593

### 594 Реализация на списках

595 Стек эффективно реализуется на односвязном списке, а следовательно и на  
 596 двухсвязном: элементы добавляются и выталкиваются из начала списка.

### 597 Реализация на массиве

598 При эффективной реализации стека на массиве элементы добавляются в ко-  
 599 нец массива и выталкиваются из конца массива. Если массив заполняется полно-  
 600 стью, то происходит его расширение по мультипликативной схеме.

### 2.3.2. Очередь

Абстрактный тип данных *очередь* определяется следующими запросами:

- `enqueue(x)` — добавить элемент  $x$  в очередь,
- `dequeue()` — извлечь элемент из очереди.

Как и в обычной жизни, порядок, в котором элементы извлекаются из очереди, должен соответствовать порядку, в котором они очередь добавляются. Про очередь говорят, что она реализует концепцию «первым пришёл — первым ушёл» («first in, first out», FIFO).

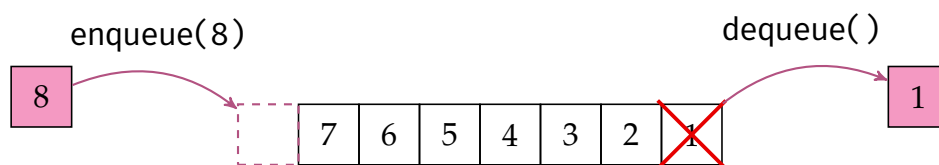


Рис. 2.3.2. Операции с очередью.

### Реализация на списках

Очередь можно эффективно реализовать на двусвязном списке, т.к. он позволяет добавлять элементы в начало и удалять из конца за  $O(1)$ . Если немного подумать, то становится ясно, что очередь можно эффективно реализовать и на односвязном списке — для этого нужно добавлять элементы в конец списка, а удалять из начала.

### Реализация на массиве

Для эффективной реализации очереди на массиве требуется использовать идею *кольцевого буфера*. Для этого в массиве поддерживаются два указателя: на начало очереди и на конец очереди (см. рис. 2.3.3). Элементы добавляются по указателю на конец очереди, выталкиваются по указателю на начало. При этом, если указатель достигает конца массива, то он «перескакивает по циклу» на начало массива. Если после добавления элемента в очередь указатели на начало и конец встретились, т.е. весь массив заполнен элементами очереди, то массив расширяется по мультипликативной схеме, и элементы очереди копируются в новый массив от начала до конца.

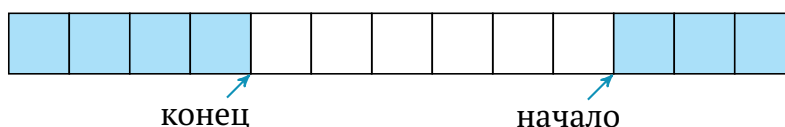


Рис. 2.3.3. Реализация очереди на массиве с использованием кольцевого буфера.

## 625 Реализация на двух стеках

626 В некоторых случаях можно описать реализацию одного абстрактного типа  
 627 данных через другой абстрактный тип данных. В частности в некоторых алгорит-  
 628 мах используется следующий способ реализации очереди на двух стеках. Пусть  
 629 далее  $s_1$  и  $s_2$  — это два стека, на основе которых мы должны реализовать запро-  
 630 сы enqueue и dequeue. Будем предполагать, что стеки реализуют запросы push,  
 631 pop и также запрос empty, который возвращает True, если стек пуст. Реализация  
 632 запроса enqueue будет затрагивать только первый стек.

```
633 def enqueue(x):
634     s1.push(x)
```

635 При запросе dequeue проверяется, есть ли элементы во втором стеке. Если второй  
 636 стек пуст, то мы перекаладываем элементы первого стека во второй, и после этого  
 637 выталкиваем верхний элемент из второго стека (см. рис. 2.3.4).

```
638 def dequeue(x):
639     if s2.empty():
640         while not s1.empty():
641             s2.push(s1.pop())
642
643     return s2.pop()
```

644 При такой реализации очереди запрос dequeue будет иметь сложность  $\Theta(n)$   
 645 в худшем (при условии, что стек реализован эффективно). Однако, можно пока-  
 646 зать, что суммарные накладные расходы на все запросы dequeue так же будут  
 647 линейными. Для этого нужно заметить, что для каждого элемента очереди будет  
 648 сделано не более двух запросов push и не более двух запросов pop. Следовательно,  
 649 сложность любой последовательности запросов к такой очереди для  $n$  элементов  
 650 ограничена  $O(n)$ .

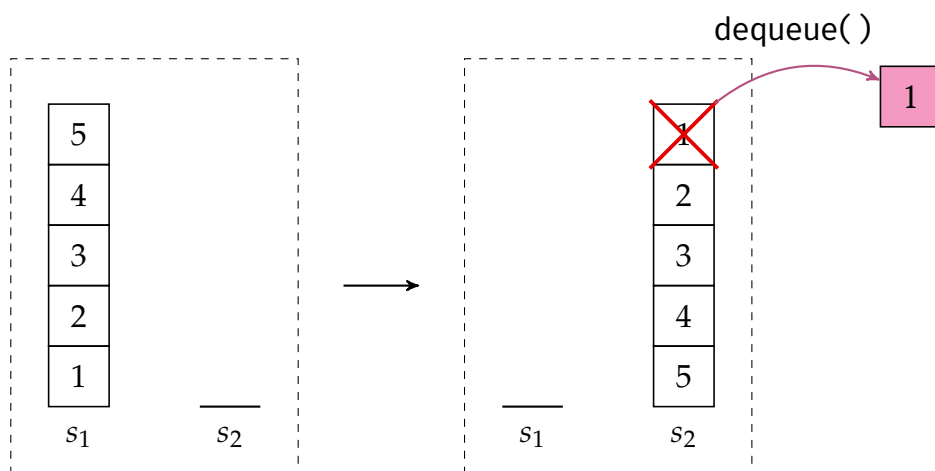


Рис. 2.3.4. Если при запросе dequeue() стек  $s_2$  пуст, то содержимое стека  $s_1$  поэлементно перекаладывается в стек  $s_2$ .

651 *Упражнение 2.3.1.* Объясните, почему не получится реализовать стек при помощи  
 652 нескольких (константного количества) очередей.

## 653 2.4. Двухсторонняя очередь

654 *Двухсторонняя очередь* позволяет добавлять и выталкивать элементы с обеих  
655 сторон, т.е. определяется следующими запросами:

- 656 • `enqueue_back(x)` — добавляет элемент  $x$  в конец очереди,
- 657 • `enqueue_front(x)` — добавляет элемент  $x$  в начало очереди,
- 658 • `dequeue_front()` — вытаскивает элемент из начала очереди.
- 659 • `dequeue_back()` — вытаскивает элемент из конца очереди.

### 660 Реализация на списках

661 Двухстороннюю очередь можно эффективно реализовать на двусвязном спис-  
662 ке, но не на односвязном, т.к. он не позволяет эффективно удалять из конца спис-  
663 ка.

### 664 Реализация на массиве

665 Эффективная реализация двухсторонней очереди на массиве использует ту же  
666 идею кольцевого буфера, только теперь указатели могут двигаться в обоих на-  
667 правлениях. Соответственно, если указатель на начало массива нужно сдвинуть  
668 влево, то он «перепрыгивает по циклу» на последний элемент массива.

## 669 2.5. Амортизационный анализ

670 *Амортизационный анализ* (англ. amortized analysis) — это способ подсчёта сред-  
671 ней сложности запроса к структуре данных. Он позволяет оценить *амортизиро-*  
672 *ванную стоимость*<sup>2</sup> запроса — средняя сложность (в худшем случае), где усредне-  
673 ние берётся по всем запросам к структуре данных. Мы рассмотрим два метода  
674 получения таких оценок.

### 675 2.5.1. Метод учётных стоимостей

676 Опишем технику оценки амортизированной стоимости запросов к структу-  
677 рам данных. Пусть к некоторой структуре данных, содержащей не более  $n$  эле-  
678 ментов, выполняются  $m$  запросов, и пусть  $c_1, \dots, c_m$  — *истинные* стоимости этих  
679 запросов (т.е. реальное количество операций, потраченных на их выполнение).  
680 Если (индивидуальную) сложность каждого запроса нельзя оценить лучше, чем  
681  $O(f(n))$ , то суммарная стоимость  $m$  таких запросов можно оценить так

$$682 \sum_{i=1}^m c_i = O(m \cdot f(n)).$$

---

<sup>2</sup>В амортизационном анализе принято использовать термин «стоимость» как синоним «сложности»

683 Пусть  $S = \{S_0, S_1, \dots, S_m\}$  — это множество состояний структуры данных, где  
 684  $S_0$  — это исходное состояние, а  $S_i$  — состояние после запроса с номером  $i$ . Для лю-  
 685 бой функции  $\Phi : S \rightarrow \mathbb{R}_+$  можно определить *учётную стоимость запроса  $i$  отно-*  
 686 *сительно  $\Phi$* , обозначаемую  $\tilde{c}_i$ , следующим образом:

$$687 \quad \tilde{c}_i = c_i + \Phi(S_i) - \Phi(S_{i-1}).$$

688 Для удобства введём обозначение  $\Phi_i = \Phi(S_i)$ . В таких обозначениях

$$689 \quad \tilde{c}_i = c_i + \Phi_i - \Phi_{i-1}.$$

690 Просуммируем это соотношение по всем  $i$ :

$$691 \quad \sum_{i=1}^m \tilde{c}_i = \sum_{i=1}^m c_i + \sum_{i=1}^m \Phi_i - \sum_{i=1}^m \Phi_{i-1} = \sum_{i=1}^m c_i + \underbrace{\Phi_m - \Phi_0}_{\Delta\Phi}.$$

### 692 Основная идея метода

693 Пусть про  $\Phi$  дополнительно известно, что  $\Delta\Phi \geq 0$  (можно потребовать, напри-  
 694 мер, чтобы  $\Phi(S_0) = 0$ ). Тогда мы можем оценить сумму истинных стоимостей как

$$695 \quad \sum_{i=1}^m c_i = \sum_{i=1}^m \tilde{c}_i - \Delta\Phi \leq \sum_{i=1}^m \tilde{c}_i.$$

696 Заметим, что у нас есть большая свобода в выборе функции  $\Phi$ , т.к. от неё по-  
 697 чти ничего не требуется. Если нам удастся подобрать  $\Phi$  такую, что учётная стои-  
 698 мость<sup>3</sup> любого запроса будет ограничена  $O(g(n))$ , т.е.  $\tilde{c}_i = O(g(n))$  для любого  $i$ ,  
 699 и при этом  $g(n) = o(f(n))$ , то в результате мы получим лучшую оценку на общее  
 700 число операций

$$701 \quad \sum_{i=1}^m c_i \leq \sum_{i=1}^m \tilde{c}_i = O(m \cdot g(n)) = o(m \cdot f(n)).$$

702 В таком случае мы будем говорить<sup>4</sup>, что *амортизированная стоимость запроса* не  
 703 более  $O(g(n))$ .

### 704 Пример: амортизация стоимости печи в пекарне

705 Слово «амортизация» в контексте анализа алгоритмов используется в том же  
 706 смысле, в котором оно используется в экономике. Попробуем разобрать идею ме-  
 707 тода учётных стоимостей на следующей несложной экономической задаче. Пред-  
 708 ставим, что вам в наследство досталась пекарня, и теперь вы разрабатываете для  
 709 неё бизнес-план. Для простоты будем считать, что вы планируете выпекать толь-  
 710 ко один вид булочек. Вы оценили затраты на производство одной булочки в 40  
 711 рублей — сюда входят затраты на ингредиенты, оплату труда пекаря, аренду и  
 712 прочее. Осталось учесть только одну «мелочь»: печь имеет ограниченный ресурс

<sup>3</sup>Учётная стоимость используется как удобная абстракция, которая не имеет «физического» смысла — в зависимости от разных  $\Phi$  мы будем получать разные учётные стоимости, которые ничего не говорят про истинную стоимость запросов.

<sup>4</sup>Тут важно не забывать, что мы оценили амортизированную стоимость запроса относительно последовательности из  $m$  запросов. Сама по себе эта оценка ничего не говорит про сложность одного запроса.



и её требуется заменить после выпекания 10 тысяч булочек. Стоимость новой печи 50 тысяч рублей. Таким образом, первые 9999 тысяч булочек действительно будут стоить вам по 40 рублей за каждую. А вот булочка с номером 10000 обойдётся вам уже в 50040 рублей. Вместо того, чтобы требовать с «несчастливого» покупателя булочки с номером 10000 дополнительные 50 тысяч рублей, вы решаете разделить эту сумму между всеми покупателями. Для этого вы добавляете 5 рублей к стоимости каждой булочки (т.е. считаете, что затраты на производство равны 45 рублям), и в результате накапливаете 50 тысяч непосредственно к тому моменту, когда печь нужно будет заменить.

В наших обозначениях, истинная стоимость производства  $i$ -ой булочки  $c_i = 40$  для всех  $i < 10000$  и  $c_{10000} = 50040$ . В качестве потенциала выберем функцию

$$\Phi(i) = 5i \pmod{50000},$$

которая отражает количество отложенных денег после производства  $i$  булочки. Если после производства булочки мы накопили заветные 50 тысяч, то мы их сразу же тратим на обновление печи. Тогда для всех  $i < 10000$  учётная стоимость производства булочки равна 45 рублям:

$$\tilde{c}_i = 40 + \Phi(i) - \Phi(i - 1) = 40 + 5i - 5(i - 1) = 45.$$

В то же время учётная стоимость производства булочки с номером 10000 так же получится равной 45 рублям:

$$\tilde{c}_{10000} = 50040 + \Phi(10000) - \Phi(9999) = 50040 + 0 - 5 \cdot 9999 = 45.$$

Таким образом, не смотря на то, что истинная стоимость производства некоторых булочек очень высокая (т.к. включает в себя стоимость замены печи), за счёт амортизации мы ограничили затраты на каждую булочку всего 45 рублями. Аналогично мы поступаем в применении к анализу алгоритмов — мы разделяем истинную стоимость «дорогих» операций между всеми операциями и оцениваем получившиеся учётные стоимости.

**Замечание.** Для этого простого примера нам удалось придумать естественные интерпретации для функции  $\Phi$  и учётных стоимостей. В дальнейшем мы столкнёмся со значительно более сложными потенциалами, которым будет сложно придать какой-то физический смысл.

### Пример: анализ расширяющегося массива

Применим метод учётных стоимостей для оценки сложности запроса добавления элемента в расширяющийся массив. Мы уже знаем, что вне зависимости от схемы выделения памяти в худшем случае сложность добавления нового элемента в массив  $O(n)$ . Однако мы также показали, что при мультипликативной схеме выделения памяти в среднем на добавление каждого нового элемента мы потратим  $O(1)$ . Давайте докажем это при помощи метода учётных стоимостей.

Будем предполагать, что мы начинаем с пустого массива и последовательно добавляем в него  $n$  элементов, кроме того коэффициент расширения  $\alpha = 2$ . Для получения оценки нам достаточно оценить количество копирований как самой трудоёмкой части запроса.

754 Для доказательства нам нужно предъявить функцию потенциала, относитель-  
 755 но которой учётная стоимость запросов будет константной. Чтобы подобрать та-  
 756 кую функцию, давайте посмотрим, сколько копирований происходит при добав-  
 757 лении нового элемента в конец массива в зависимости от его номера. Если это  
 758 первый элемент, то требуется одно копирование: выделяется массив единичного  
 759 размера и в него копируется первый элемент. Для добавления второго элемента  
 760 требуется выделить массив размера два и скопировать туда первый и второй эле-  
 761 менты. И так далее. Заметим, что для добавления четвёртого элемента нам требу-  
 762 ется только одно копирование, т.к. к этому моменту мы уже выделили массив раз-  
 763 мера 4, одна из ячеек которого свободна. Аналогично, только одно копирование  
 требуется для пятого, шестого, седьмого и восьмого элементов (см. рис. 2.5.1).

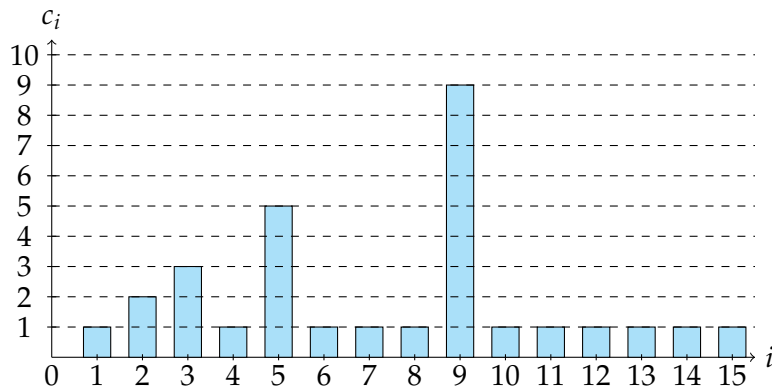


Рис. 2.5.1. Стоимость добавления элемента в расширяющийся массив.

764 Искомая функция потенциала в должна вести себя в некотором смысле про-  
 765 тивоположно тому, как ведут себя стоимости операций  $c_i$ : в те моменты, когда  $c_i$   
 766 большое, функция потенциала должна резко убывать, компенсируя таким обра-  
 767 зом  $c_i$ , а в тех случаях, когда для добавления требуется только одна операция, то  
 768 потенциал должен возрасти не сильно. Для того, чтобы определить такую функ-  
 769 цию, нам нужно как-то описать состояние расширяющегося массива. Давайте за  $l$   
 770 обозначим размер массива, а за  $s$  — количество заполненных ячеек (см. рис. 2.5.2).

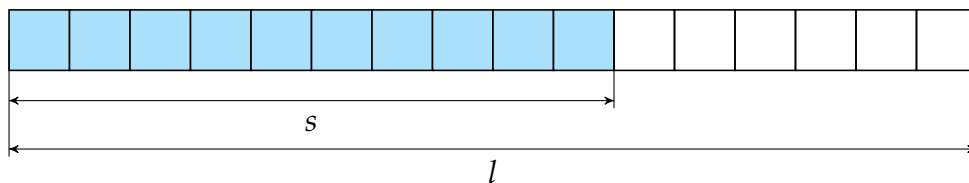
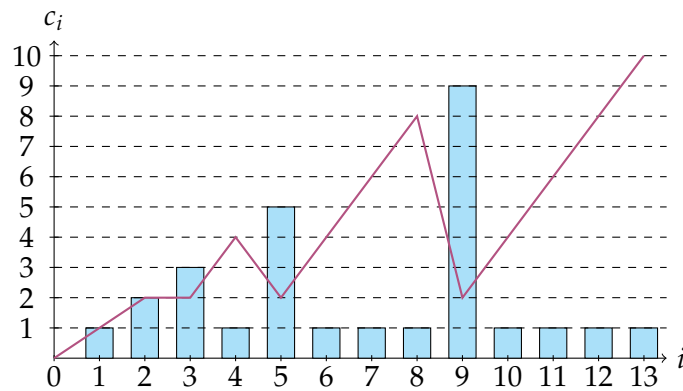


Рис. 2.5.2. Массив размера  $l$  с  $s$  элементами.

772 Заметим, что  $s$  с каждым добавленным элементом возрастает на 1, а  $l$  либо  
 773 не изменяется, либо увеличивается в 2 раза, если массив был полностью запол-  
 774 нен, т.е., когда для добавления нового элемента требуются дополнительные ко-  
 775 пирования. На основе этих наблюдений можно догадаться, что нам может подой-  
 776 ти линейная функция потенциала, которая положительно зависит от  $s$  и отрица-  
 777 тельно от  $l$  — тогда эта функция будет медленно возрастать, пока в массиве есть  
 778 свободные ячейки, и резко падать, если массив заполнился и произошло расши-  
 779 рение. Осталось подобрать подходящие коэффициенты. Исходя из того, что при  
 780 расширении массива уменьшение функции потенциала должно компенсировать

Рис. 2.5.3. Потенциал  $\Phi$  и стоимость добавления элемента массива.

781 дополнительные затраты на копирования, предлагается использовать функцию  
782  $\Phi(s, l) = 2s - l$  (см. рис. 2.5.3).

783 Нам осталось показать, что учётные стоимости запросов на добавление дей-  
784 ствительно ограничены константой относительно предложенной функции по-  
785 тенциала  $\Phi(s, l) = 2s - l$ . По определению учётная стоимость равна

$$786 \quad \tilde{c}_i = c_i + \Phi_i - \Phi_{i-1}.$$

787 Пусть  $l$  и  $s$  обозначают соответственно размер массива и количество элементов в  
788 нём до добавления элемента с номером  $i$ . Нам нужно рассмотреть два случая.

789 1. В массиве есть пустые ячейки, т.е.  $s < l$ . Тогда для добавления нам потребу-  
790 ется только одно копирование, т.е.  $c_i = 1$ , следовательно

$$791 \quad \tilde{c}_i = c_i + \Phi(s + 1, l) - \Phi(s, l) = 1 + (2(s + 1) - l) - (2s - l) = 3.$$

792 2. Массив полностью заполнен, т.е.  $s = l$ . В этом случае нужно скопировать все  
793  $s$  элементов в новый массив и потом скопировать туда же новый элемент,  
794 т.е.  $c_i = s + 1$ . Поэтому

$$795 \quad \tilde{c}_i = c_i + \Phi(s + 1, 2l) - \Phi(s, l) = (s + 1) + (2(s + 1) - 2l) - (2s - l) = 3$$

796 (тут мы пользуемся тем, что  $s = l$ ).

797 Мы показали, что относительно такой функции потенциала учётные стоимо-  
798 сти запросов ограничены константой. Осталось заметить, что потенциал пустого  
799 массива  $\Phi_0 = \Phi(0, 0) = 0$ . Это позволяет нам заключить, что запрос на добав-  
800 ление элемента в конец массива при использовании мультипликативной схемы  
801 расширения имеет амортизированную оценку сложности  $O(1)$ .

802 *Упражнение 2.5.1.* Подберите функцию потенциала для произвольного коэффи-  
803 циента расширения  $\alpha > 1$ .

804 *Упражнение 2.5.2.* Рассмотрим двоичный счётчик, который реализует запрос на  
805 увеличение значения на единицу. Если в счётчике хранится число  $n$ , то для увели-  
806 чения его на единицу в худшем случае необходимо изменить значения  $1 + \lfloor \log_2 n \rfloor$   
807 битов, т.е. сложность этого запроса равна  $O(\log_2 n)$ . Покажите, что амортизиро-  
808 ванную стоимость запроса на увеличение не превосходит  $O(1)$ .

### 809 2.5.2. Метод предоплаты

810 *Метод предоплаты* значительно менее формален, чем метод учётных стоимо-  
811 стей, и за счёт этого позволяет упростить некоторые рассуждения. Метод основан  
812 на следующей идее: для того, чтобы показать, что суммарное количество опера-  
813 ций невелико, можно каждому элементу структуры данных «выдать» некоторое  
814 количество монеток, которыми этот элемент сможет «оплачивать» операций с  
815 ним. Если показать, что каждому элементу хватит выданных монеток, и что в  
816 сумме на  $m$  запросов к структуре данных, хранящей  $n$  элементов, нам потребует-  
817 ся не более  $O(m \cdot g(n))$  монет, то можно заключить, что амортизированную сто-  
818 имость запросов ограничена  $O(g(n))$ .

#### 819 **Пример: анализ очереди на двух стеках**

820 В разделе 2.3.2 мы рассматривали реализацию очереди на двух стеках. Вос-  
821 пользуемся методом предоплаты для того, чтобы показать, что запросы к такой  
822 очереди имеют амортизированную стоимость  $O(1)$ . Для этого каждому элемен-  
823 ту, который попадает в очередь, выдаётся четыре монетки на оплату операций со  
824 стеками. Легко увидеть, что четырёх монеток достаточно: каждый элемент дол-  
825 жен оплатить две операции push и две операции pop. В сумме на оплату любых  $m$   
826 запросов нам потребуется не более  $4m$  монеток. Следовательно амортизирован-  
827 ная стоимость запросов к такой очереди  $O(1)$ .

## 828 Глава 3

# 829 Разделяй и властвуй

### 830 3.1. Переход к подзадачам

831 В этом курсе мы будем изучать не только алгоритмы для конкретных задач,  
832 но и общие методы построения алгоритмов. Начнём с одного из самых естествен-  
833 ных подходов, который получил название «разделяй и властвуй» («divide and conquer»).

834 Идея этого метода заключается в следующих трёх шагах:

- 835 • разделить вход алгоритма на кусочки меньшего размера,
- 836 • решить *подзадачи* для этих кусочков рекурсивно,
- 837 • вычислить решение для исходной задачи на основании решения подзадач.

838 Далее мы проиллюстрируем работу этого метода на нескольких примерах и в за-  
839 вершение докажем общую теорему, которая поможет нам оценивать сложность  
840 алгоритмов, построенных на этой идее.

### 841 3.2. Двоичный поиск

842 Рассмотрим задачу поиска заданного элемента в упорядоченном массиве. Те,  
843 кто когда-либо пытался найти слово в словаре (имеется в виду настоящая бумаж-  
844 ная книга), уже интуитивно знают как нужно действовать. Не стоит пытаться най-  
845 ти слово последовательно просматривая все страницы словаря. Нужно действо-  
846 вать умнее и открыть словарь примерно на середине. Посмотреть, какие слова  
847 находятся на открытой странице, и в зависимости от этого продолжить поиск в  
848 первой или второй половине книги.

849 Используя эту же идею мы построим эффективный алгоритм для этой задачи.  
850 Давайте возьмём средний элемент массива и сравним его с искомым. Если иско-  
851 мый элемент меньше, то он может быть только в левой половине. Если искомый  
852 элемент больше, то он, соответственно, может быть только во второй половине  
853 массива. А если нам повезло и искомый элемент равен среднему элементу, то мы  
854 можем завершить работу.

```
855 # ищем самое левое вхождение x в array[left:right]  
856 def binary_search(array, x, left, right):  
857     #  
858     if right == left:  
859         return None
```

860

```
861 if right == left + 1:  
862     if array[left] == x:  
863         return left  
864     else:  
865         return None  
866  
867     mid = left + (right - left) // 2  
868  
869     if array[mid] < x:  
870         return binary_search(array, x, mid + 1, right)  
871     elif array[mid] == x and array[mid - 1] < x:  
872         return mid  
873     else:  
874         return binary_search(array, x, left, mid)
```

### 875 3.3. Умножение $n$ -битовых чисел

### 876 3.4. Сортировка слиянием

877 Алгоритм был изобретён Джоном фон Нейманом в 1945 году.

```
878 # сортировка слиянием
879 def merge_sort(array):
880     n = len(array)
881
882     # пустой и одноэлементный массивы сортировать не нужно
883     if n < 2:
884         return array
885
886     # рекурсивно сортируем левую и правую половины массива
887     array1 = merge_sort(array[:n/2])
888     array2 = merge_sort(array[n/2:])
889
890     # возвращаем результат слияние двух упорядоченных массивов
891     return merge(array1, array2)
892
893 # слияние двух упорядоченных массивов
894 def merge(array1, array2):
895     n, m = len(array1), len(array2)
896
897     # переменные цикла
898     i, j = 0, 0
899
900     # массив для результата слияния
901     result = []
902
903     # перебираем элементы массивов
904     while i < n and j < m:
905         # на каждой итерации выбираем минимальный из текущих
906         # и записываем его в массив с результатом
907         if array1[i] < array2[j]:
908             result.append(array1[i])
909             i = i + 1
910         else:
911             result.append(array2[j])
912             j = j + 1
913
914     # добавляем в конец оставшиеся элементы
915     while i < n:
916         result.append(array1[i])
917         i = i + 1
918
919     while j < m:
920         result.append(array2[j])
921         j = j + 1
922
923     return result
```

### 924 3.5. Основная теорема о рекуррентных соотношениях

925 **Теорема 3.5.1.** Пусть  $T(n)$  удовлетворяет следующему рекуррентному соотноше-  
926 нию:

$$927 T(1) = O(1), \quad T(n) = a \cdot T(\lceil n/b \rceil) + O(n^d)$$

#### 928 3.5.1. Как побороть аддитивные константы?

929 Рекуррентное соотношение в теореме 3.5.1 задаётся в очень удобной форме.  
930 На практике соотношение может выглядеть менее красивым. Рассмотрим два слу-  
931 чая.

932 1. Пусть дано следующее соотношение:  $T(n) = 2 \cdot T(\lceil n/2 \rceil - 4) + O(n)$ . В этом  
933 случае мы можем сказать, что  $T(n)$  — возрастающая функция, и поэтому

$$934 T(n) = 2 \cdot T(\lceil n/2 \rceil - 4) + O(n) \leq 2 \cdot T(\lceil n/2 \rceil) + O(n),$$

935 что позволяет применить теорему о рекуррентных соотношениях и полу-  
936 чить оценку  $T(n) = O(n \log n)$ .

937 2. Пусть теперь задано соотношение:  $T(n) = 2 \cdot T(\lceil n/2 \rceil + 4) + O(n)$ . Предду-  
938 щие рассуждения здесь уже не сработают, т.к.  $\lceil n/2 \rceil + 4 > \lceil n/2 \rceil$ . Мы могли  
939 бы повторить доказательство теоремы 3.5.1 для того, чтобы показать, что  
940 она верна для  $T(n) = a \cdot T(\lceil n/b \rceil + c) + O(n^d)$  с произвольной константой  $c$ .  
941 Это же можно объяснить при помощи следующего трюка с заменой пере-  
942 менной. Введём переменную  $m = n - 8$ , т.е.  $n = m + 8$ , и заменим перемен-  
943 ную в заданном рекуррентном соотношении. Получим

$$944 T(m + 8) = 2 \cdot T(\lceil (m + 8)/2 \rceil + 4) + O(m + 8) = 2 \cdot T(\lceil m/2 \rceil + 8) + O(m).$$

945 Теперь рассмотрим функцию  $T'(m) = T(m + 8)$ . Для неё выполняется следу-  
946 ющее рекуррентное соотношение:

$$947 T'(m) = 2 \cdot T'(\lceil m/2 \rceil) + O(m),$$

948 для которого по теореме 3.5.1 получается оценка  $O(m \log m)$ . Нам осталось  
949 снова пройти по цепочке замен и получить требуемую оценку

$$950 T(n) = T(m + 8) = T'(m) = O(m \log m) = O((n - 8) \log(n - 8)) = O(n \log n).$$

951 **Упражнение 3.5.1.** Покажите, что теорема 3.5.1 верна для рекуррентных соотно-  
952 шений вида  $T(n) = a \cdot T(\lceil n/b \rceil + c) + O(n^d)$  с произвольной константой  $c$ .

### 953 3.6. Умножение матриц



## 954 Глава 4

# 955 Сортировка

### 956 4.1. Квадратичные сортировки

```
957 def swap(array, i, j):  
958     array[i], array[j] = array[j], array[i]
```

```
959 def selection_sort(array):  
960     n = len(array)  
961  
962     for i in range(n):  
963         min_pos = i  
964         for j in range(i + 1, n):  
965             if array[min_pos] > array[j]:  
966                 min_pos = j  
967  
968         swap(array, i, min_pos)  
969  
970     return array
```

```
971 def insertion_sort(array):  
972     n = len(array)  
973  
974     for i in range(n):  
975         for j in range(i, 0, -1):  
976             if array[j - 1] > array[j]:  
977                 swap(array, j - 1, j)  
978  
979     return array
```

```
980 def bubble_sort(array):  
981     n = len(array)  
982  
983     for i in range(n - 1):  
984         for j in range(n - i - 1):  
985             if array[j] > array[j + 1]:  
986                 swap(array, j, j + 1)  
987  
988     return array
```

## 989 4.2. Нижняя оценка на сортировку сравнениями

990 Мы уже знаем, что массив можно отсортировать за  $O(n \log n)$  при помощи сор-  
 991 тировки слиянием. А можно ли быстрее? Оказывается, что если при сортировке  
 992 использовать *только сравнения* элементов, то сортировать быстрее не получится.  
 993 Будем говорить, что сортировка *основана на сравнениях*, если она не использует  
 994 никаких других операций с ключами, кроме сравнений (т.е., например, не ис-  
 995 пользуется арифметических операций с ключами).

996 **Теорема 4.2.1.** *Любой алгоритм сортировки массива длины  $n$  требует  $\Omega(n \log n)$*   
 997 *сравнений.*

998 Для доказательства этой теоремы мы сначала докажем её для более слабой  
 999 модели вычислений.

1000 **Определение 4.2.1.** *Дерева сравнений* — это модель вычислений, для которой  
 1001 каждая программа представляется в виде двоичного дерева, в котором внутрен-  
 1002 ние вершины соответствуют сравнениям элементов, а листья — перестановкам  
 1003 элементов.

1004 Для объяснения того, как работает эта модель, давайте рассмотрим следую-  
 щую программу для сортировки трёх элементов. На вход программе подаётся

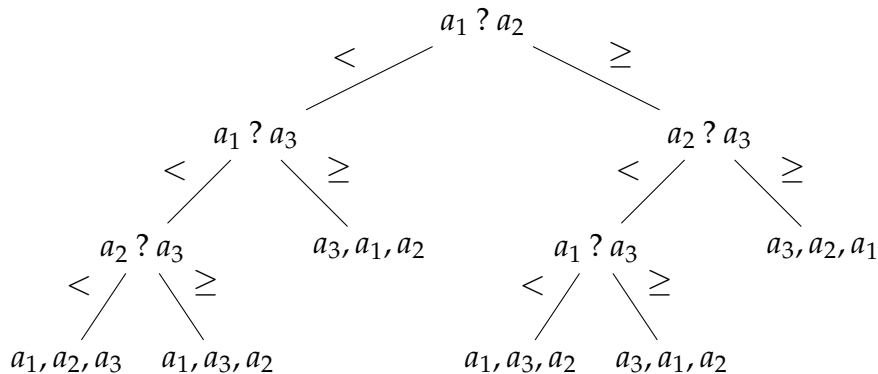


Рис. 4.2.1. Дерево сравнений для сортировки трёх элементов.

1005 набор из трёх чисел:  $a_1, a_2, a_3$ . В каждый момент времени программа находится  
 1006 в одной из вершин дерева. Один шаг программы состоит в выполнении сравне-  
 1007 ния двух элементов, которые указаны в текущей вершине, и переходе в одного  
 1008 из сыновей текущей вершины, в зависимости от результата сравнения. Выпол-  
 1009 нение программы начинается в корне дерева и завершается в одном из листьев,  
 1010 результат выполнения — перестановка элементов, записанная в соответствую-  
 1011 щем листе.

1013 В данной программе на первом шаге выполняется сравнение  $a_1$  и  $a_2$ . В зависи-  
 1014 мости от результата сравнения исполнение переходит в одного из сыновей: если  
 1015 оказывается, что  $a_1 < a_2$ , то исполнение переходит в левого сына, а иначе — в пра-  
 1016 вого. Исполнение продолжается, пока не достигнет листа. Проверьте, что данная  
 1017 программа действительно корректно сортирует любой набор из трёх элементов.

1018 Теперь мы докажем, что в глубина любого дерева сравнений, упорядочиваю-  
 1019 щее  $n$  элементов, имеет глубину  $\Omega(n \log n)$ .

1020 **Лемма 4.2.1.** *Любое дерево сравнений корректно сортирующее  $n$  элементов имеет*  
 1021 *глубину  $\Omega(n \log n)$ .*

*Доказательство.* Если дерево сравнений сортирует  $n$  элементов корректно, то среди его листьев должны встречаться все возможные  $n!$  перестановок из  $n$  элементов. Таким образом, количество листьев у любого такого дерева не меньше  $n!$ . Теперь нам потребуется воспользоваться следующим свойством двоичных деревьев: в дереве глубины  $h$  не более  $2^h$  листьев. Действительно, если дерево высоты  $h$  — полное, т.е. в него нельзя добавить вершин, не увеличив высоты, то в нём ровно  $2^h$  листьев, следовательно в произвольном дереве глубины  $h$  листьев не больше  $2^h$ . Отсюда получаем следующее соотношение на глубину дерева  $h$ :

$$2^h \geq n! = 1 \cdot 2 \cdot \dots \cdot n.$$

Логарифмируя обе части неравенство получаем

$$h \geq \log(n!) = \log(1) + \log(2) + \dots + \log(n).$$

Оставим только половину слагаемых справа.

$$h \geq \log(\lceil n/2 \rceil) + \dots + \log(n) \geq \frac{n}{2} \cdot \log(\lceil n/2 \rceil) = \Omega(n \log n).$$

1022

□

### 1023 4.3. Сортировка кучей

1024 **4.4. Быстрая сортировка**

```
1025 def quick_sort(array, left, right):
1026     if right <= left + 1:
1027         return
1028
1029     pivot = random.randint(left, right - 1)
1030
1031     mid = partition(array, pivot, left, right)
1032
1033     quick_sort(array, left, mid)
1034     quick_sort(array, mid + 1, right)
1035
1036     return array
```

```
1037 def partition(array, pivot, left, right):
1038     x = array[pivot]
1039
1040     swap(array, pivot, right - 1)
1041
1042     m = left
1043
1044     for p in range(left, right - 1):
1045         if array[p] < x:
1046             swap(array, p, m)
1047             m += 1
1048
1049     swap(array, m, right - 1)
1050     return m
```

1051 **4.4.1. Анализ времени работы**

1052 **4.4.2. Сортировка массива с повторяющимися элементами**

```
1053 def quick_sort3(array, left, right):
1054     if right <= left + 1:
1055         return
1056
1057     pivot = random.randint(left, right - 1)
1058
1059     (m1, m2) = partition3(array, pivot, left, right)
1060
1061     quick_sort3(array, left, m1)
1062     quick_sort3(array, m2, right)
1063
1064     return array
```

```
1065 def partition3(array, pivot, left, right):
1066     x = array[pivot]
1067
1068     swap(array, pivot, right - 1)
1069
1070     m = left
1071     p = left
1072     q = right - 1
1073
1074     while p < q:
1075         if array[p] == x:
1076             q -= 1
1077             swap(array, p, q)
1078         else:
1079             if array[p] < x:
1080                 swap(array, p, m)
1081                 m += 1
1082             p += 1
1083
1084     for i in range(q, right):
1085         swap(array, m + i - q, i)
1086
1087     return m, m + right - q
```

1088 **4.5. Линейные сортировки**1089 **4.5.1. Сортировка подсчётом**

```
1090 def counting_sort(array, max_value):
1091     counters = [0 for k in range(max_value + 1)]
1092
1093     for k in array:
```

```
1094     counters[k] += 1
1095
1096     result = []
1097     for k in range(max_value + 1):
1098         for i in range(counters[i]):
1099             result.append(k)
1100
1101     return result
```

```
1102 def counting_sort_stable(array, max_value):
1103     counters = [0 for k in range(max_value + 1)]
1104
1105     for (k, v) in array:
1106         counters[k] += 1
1107
1108     idx = 0
1109     indices = []
1110     for k in range(max_value + 1):
1111         indices.append(idx)
1112         idx += counters[k]
1113
1114     result = [None for k in range(len(array))]
1115     for (k, v) in array:
1116         result[indices[k]] = (k, v)
1117         indices[k] += 1
1118
1119     return result
```

1120 **4.5.2. Цифровая сортировка**

1121 **4.5.3. Блочная сортировка**

## 1122 Глава 5

### 1123 Задачи RMQ и LCA

1124 В данном разделе мы изучим две задачи, которые неожиданным образом ока-  
1125 зываются связанными: *задача о минимуме на отрезке* и *задача о наименьшем об-*  
1126 *щем предке*.

#### 1127 5.1. Задача о минимуме на отрезке

1128 **Определение 5.1.1.** *Задача о минимуме на отрезке (range minimum query, RMQ):*  
1129 *по последовательности  $A = (a_1, \dots, a_n)$  и паре индексов  $(i, j)$ ,  $i \leq j$ , найти такой*  
1130 *индекс  $k \in [i, j]$ , что  $a_k = \min\{a_i, \dots, a_j\}$ .*

1131 В такой постановке эта задача имеет сложность  $\Theta(n)$ : её можно решить за  $O(n)$   
1132 *прямым вычислением*, т.е. перебрав элементы  $a_i, \dots, a_j$ , и при этом нельзя решить  
1133 быстрее, т.к. для пары индексов  $(1, n)$  нам потребуется перебрать все элементы  
1134 последовательности. Задача становится более интересной, если к одной и той же  
1135 последовательности совершается множество запросов. В этом случае мы можем  
1136 вычислить какую-то вспомогательную информацию или использовать структу-  
1137 ры данных, которые позволят отвечать на такие запросы значительно быстрее.

##### 1138 5.1.1. Динамическая задача RMQ

1139 Разрешим дополнительно кроме запроса минимума на отрезке также изме-  
1140 нять элементы последовательности  $A = (a_1, \dots, a_n)$  (но не их количество) при  
1141 помощи запроса  $\text{change}(i, x)$ , который заменяет элемент  $a_i$  на  $x$ . В таких случа-  
1142 ях принято говорить о *динамической постановке задачи RMQ*. Давайте подумаем,  
1143 какой сложности для таких запросов следует ожидать в лучшем случае. Предпо-  
1144 ложим, что на вычисление вспомогательной информации и на построение струк-  
1145 тур данных мы хотим потратить не более  $O(n)$  (на меньшее надеяться не стоит,  
1146 т.к. нам требуется прочитать все элементы последовательности). Можно доказать  
1147 следующую оценку.

1148 **Теорема 5.1.1.** *Если сложность вычисления вспомогательной информации и по-*  
1149 *строения структур данных не более  $O(n)$ , то максимальная из сложностей запро-*  
1150 *сов  $\text{minimum}$  и  $\text{change}$  в динамической постановке задачи RMQ не меньше  $\Omega(\log n)$ .*

1151 *Доказательство.* Будем доказывать от обратного. Предположим, что существу-  
1152 ет алгоритм, позволяющий выполнять запросы  $\text{minimum}(i, j)$  и  $\text{change}(i, x)$  на  
1153 массиве длины  $n$  за  $o(\log n)$  после предобработки массива функцией  $\text{rmq-preprocess}$ ,  
1154 сложность которой не более  $O(n)$ . Покажем, что в таком случае мы можем от-  
1155 сортировать произвольных массив длины  $n$  сравнениями за  $o(n \log n)$ , что про-  
1156 тиворечит нижней оценке  $\Omega(n \log n)$  на сортировку сравнениями. Пусть нам дан

1157 массив  $A$  длины  $n$ , и пусть  $M$  — это некоторое значение, которое больше любого  
 1158 элемента массива  $A$ . Рассмотрим следующую функцию, которая упорядочивает  
 1159 массив  $A$  используя алгоритм для RMQ.

```

1160 # сортирует массив используя алгоритм для RMQ
1161 def sort-using-rmq(A):
1162     n = len(A)
1163
1164     # предобработка для RMQ, не более O(n)
1165     rmq-preprocess(A)
1166
1167     # массив для результата
1168     result = []
1169
1170     for i in range(n):
1171         # вычисляем минимум для всего массива, o(log n)
1172         k = minimum(0, n - 1)
1173         # сохраняем минимум
1174         result.append(A[k])
1175         # заменяем элемент k на значение M, o(log n)
1176         change(k, M)
1177
1178     return result
  
```

1179 Сложность этой функции складывается из сложности предобработки  $O(n)$ , слож-  
 1180 ности  $n$  запросов `minimum` и сложности  $n$  запросов `change`. По нашему предпо-  
 1181 ложению сложность любого запроса не более  $o(\log n)$ , следовательно суммарная  
 1182 сложность ограничена  $o(n \log n)$ .  $\square$

1183 Таким образом, самая оптимистичная оценка — это  $O(\log n)$  на запрос, при  
 1184 условии, что на вычисление вспомогательной информации и на построение струк-  
 1185 тур данных тратится не более  $O(n)$  операций. Такую оценку позволяет достичь  
 1186 специальная структура данных, называемая *деревом отрезков*.

1187 **Определение 5.1.2.** *Дерево отрезков* для массива длины  $n$  — это двоичное дере-  
 1188 во, с каждой вершиной которого ассоциирован отрезок массива по следующему  
 1189 правилу:

- 1190 • с корнем  $r$  ассоциирован отрезок  $S(r) = [1, n]$ ,
- 1191 • для каждой вершины  $v$  с отрезком  $S(v) = [i, j]$ ,
  - 1192 – если  $i = j$ , то  $v$  — лист,
  - 1193 – в противном случае  $v$  имеет два потомка  $u$  и  $w$ , причём

$$1194 \quad S(u) = [i, m] \quad \text{и} \quad S(w) = [m + 1, j], \quad \text{где} \quad m = \lfloor (i + j) / 2 \rfloor.$$

1195 Легко проверить, что у дерева отрезков для массива длины  $n$  ровно  $n$  листьев —  
 1196 по одному на каждый элемент массива. Кроме того, так как отрезки во внутрен-  
 1197 них вершинах при переходе к потомкам уменьшаются вдвое, то высота такого  
 1198 дерева отрезков ограничена  $\lceil \log n \rceil$ . Будем называть отрезки, которые соответ-  
 1199 ствуют вершинам дерева, *каноническими* и докажем следующую лемму.



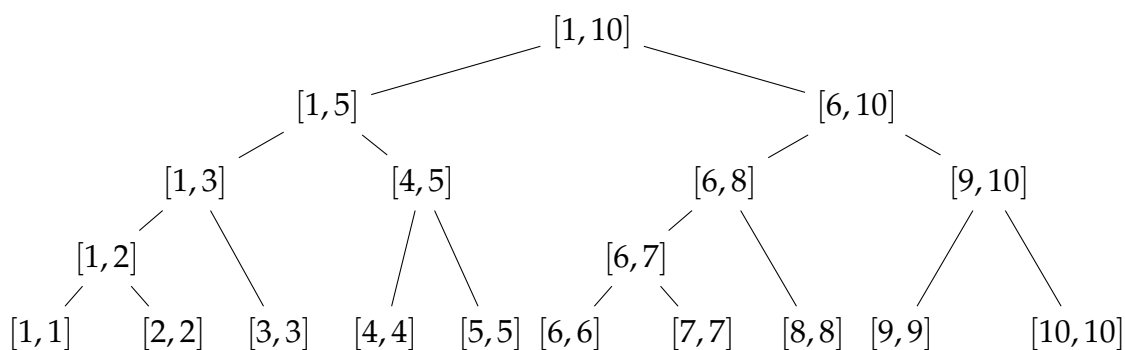


Рис. 5.1.1. Дерево отрезков для массива длины 10.

1200 **Лемма 5.1.1.** *Любой отрезок  $s \subseteq [1, n]$  может быть представлен в виде объединения*  
 1201 *непересекающихся канонических отрезков  $s_1, \dots, s_k$ , где  $k = O(\log n)$ .*

1202 *Доказательство.* Будем задавать отрезки парами целых чисел и в каждой вер-  
 1203 *шине двоичного дерева дополнительно хранить поле `segment` — ассоциирован-*  
 1204 *ный отрезок. Нам будет удобно воспользоваться следующей вспомогательной функ-*  
 1205 *цией для пересечения двух отрезков.*

```
1206 # вычисляет отрезок — пересечение двух отрезков
1207 def segment_intersect((l1, r1), (l2, r2)):
1208     # левая граница — это максимум из левых границ
1209     l = max(l1, l2)
1210
1211     # правая граница — это минимум из правых границ
1212     r = min(r1, r2)
1213
1214     # если окажется, что l > r, то пересечение пустое
1215     return (l, r)
```

1216 Рассмотрим следующую рекурсивную процедуру, которая по отрезку  $[l, r]$  возвра-  
 1217 *щает список канонических отрезков, задающий разбиение  $[l, r]$ .*

```
1218 # разбивает отрезок s на множество канонических отрезков
1219 def decompose((l, r), v = root):
1220     # обрабатываем пустой отрезок
1221     if r < l:
1222         return []
1223
1224     # возвращаем канонический отрезок
1225     if (l, r) == v.segment:
1226         return [v.segment]
1227
1228     # разбиваем s
1229     s_left = segment_intersect((l, r), v.left.segment)
1230     s_right = segment_intersect((l, r), v.right.segment)
1231
1232     # выполняем рекурсивные вызовы и объединяем результаты
1233     return decompose(s_left, v.left) + decompose(s_right, v.right)
```

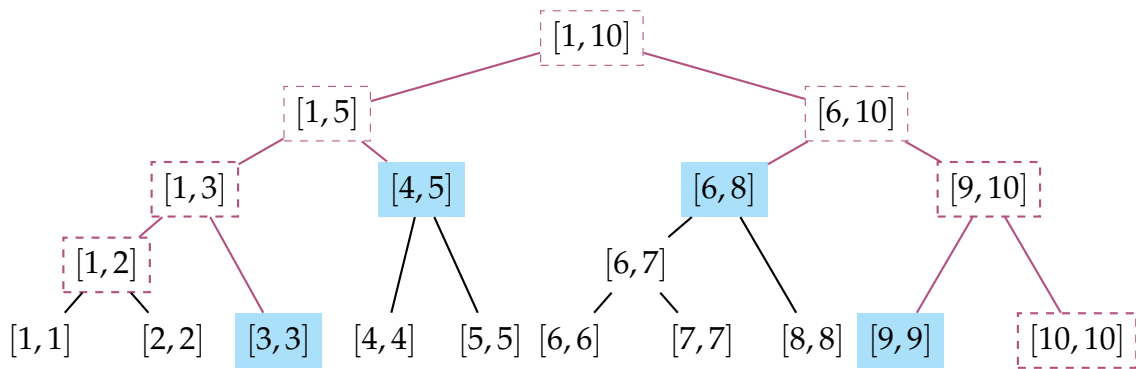


Рис. 5.1.2. Разбиение отрезка  $[3, 9]$  функцией `decompose` для массива длины 10. В результате отрезок разбивается на канонические отрезки  $[3, 3]$ ,  $[4, 5]$ ,  $[6, 8]$ ,  $[9, 9]$ . Остальные просмотренные вершины выделены пунктиром.

1234 Сначала докажем корректность, а именно покажем, что данная процедура дей-  
 1235 ствительно задаёт разбиение отрезка на канонические. Заметим, что при каждом  
 1236 рекурсивном вызове `decompose(s, v)` поддерживается следующий инвариант:  
 1237 отрезок  $s$  всегда является подотрезком  $v.segment$ . Будем доказывать коррект-  
 1238 ность этой процедуры по индукции по длине отрезка, ассоциированного с вер-  
 1239 шиной  $v$ . База индукции: процедура `decompose` корректна для листьев. Теперь  
 1240 предположим, что процедура корректна для всех вершин с отрезками длины ме-  
 1241 нее  $k$ , и покажем корректность для вершины с отрезком длины  $k$ . Если мы попали  
 1242 в один из **if**-ов, то процедура корректна. Если же произошли рекурсивные вызо-  
 1243 вы, то по предположению индукции оба вызова вернут корректное разбиение для  
 1244 левой и правой частей отрезка  $s$ , что в сумме даст корректное разбиение для  $s$ .

1245 Осталось доказать, что в разбиении получится не более  $O(\log n)$  отрезков. Сна-  
 1246 чала рассмотрим частный случай, когда один из концов отрезка  $s$  совпадает с со-  
 1247 ответствующим концом отрезка  $v.segment$ . Для определённости давайте пред-  
 1248 положим, что у  $s$  и  $v.segment$  общее начало. Тогда могут быть три варианта вы-  
 1249 полнения `decompose(s, v)`:

- 1250 • если  $s == v.segment$ , то вызов завершится без рекурсивных вызовов,
- 1251 • если  $s\_left == v.left.segment$ , то внутри `decompose(s\_left, v.left)` не  
 1252 будет рекурсивных вызовов, а у  $s\_right$  и  $v.right$  будет общее начало,
- 1253 • если  $s\_left != v.left.segment$ , то отрезок  $s\_right$  — пустой, следователь-  
 1254 но внутри `decompose(s\_right, v.right)` не будет рекурсивных вызовов, а  
 1255 у  $s\_left$  и  $v.left.segment$  будет общее начало.

1256 Другими словами, при каждом вызове `decompose` рекурсия будет продолжаться  
 1257 не более, чем в одной ветке. Поэтому всего будет не более  $2\lceil \log n \rceil$  вызовов  
 1258 `decompose`, и следовательно, `decompose(s, v)` вернёт не более  $O(\log n)$  отрезков.

1259 Теперь предположим, что у  $s$  и  $v.segment$  нет общих концов. В этом случае  
 1260 есть два варианта выполнения `decompose(s, v)`:

- 1261 • если один из отрезков  $s\_left$  и  $s\_right$  пустой, то рекурсия продолжится  
 1262 только в одном из рекурсивных вызовов `decompose`,
- 1263 • если оба отрезка  $s\_left$  и  $s\_right$  непустые, т.е. в этой вершине произо-  
 1264 шло нетривиальное разбиение отрезка  $s$  на две части, то теперь в обоих ре-  
 1265 курсивных вызовах `decompose` один из концов отрезка совпадает с концом  
 1266 отрезка, ассоциированного с вершиной.

1267 Пока мы находимся в условиях первого варианта, рекурсия продолжается толь-  
 1268 ко в одном из рекурсивных вызовов. Как только происходит нетривиальное раз-  
 1269 биение отрезка, то в обоих рекурсивных вызовах мы попадаем в условия част-  
 1270 ного случая, который мы рассмотрели выше, и про который мы доказали оценку  
 1271  $O(\log n)$  на количество отрезков. Поэтому в сумме в общем случае `decompose` вер-  
 1272 нёт не более  $O(\log n) + O(\log n) = O(\log n)$  канонический отрезков.  $\square$

1273 *Замечание 5.1.1.* Может показаться, что и время работы функции `decompose` огра-  
 1274 ничено  $O(\log n)$ . Действительно, мы доказали, что функция `decompose` посещает  
 1275  $O(\log n)$  вершин дерева. Однако, эта функция возвращает список, и на объедине-  
 1276 ние списков этих списков в Python дополнительно потратится  $O(\log^2 n)$ .

1277 *Упражнение 5.1.1.* Перепишите функцию `decompose` так, чтобы её сложность была  
 1278 не более  $O(\log n)$ .

1279 Применим дерево отрезков для решения поставленной задачи. Пусть дан неко-  
 1280 торый массив  $A$  длины  $n$ . Построим дерево отрезков для этого массива и в каждой  
 1281 вершине запишем индекс минимального элемента на соответствующем отрезке.  
 1282 Отметим, что построение дерева и заполнение индексов можно организовать за  
 1283  $O(n)$ : будем заполнять индексы в вершинах от листьев к корню, индекс в листе —  
 1284 это номер соответствующего элемента массива, а индексом во внутренней вер-  
 1285 шине выбирается один из индексов её сыновей — тот, что соответствует меньше-  
 му элементу массива.

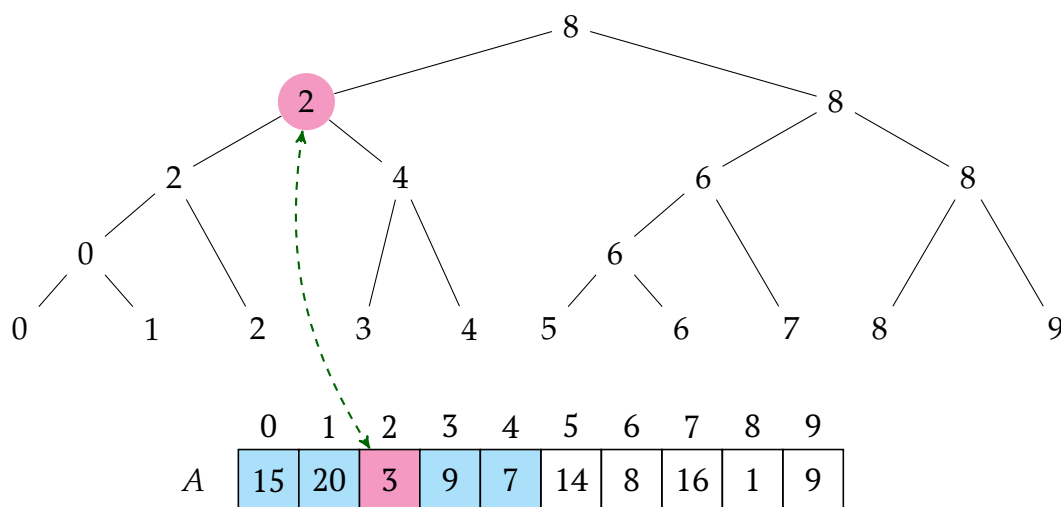


Рис. 5.1.3. Дерево отрезков с индексами минимумов для массива длины 10. Выде-  
 ленный элемент хранит индекс минимального элемента среди первых пяти эле-  
 ментов массива

1286 Пришло время реализовать процедуры `minimum` и `change`. Будем предпола-  
 1287 гать, что в каждой вершине хранится значение `value`, которое равно минималь-  
 1288 ному элементу на соответствующем отрезке. Процедура `minimum` реализуется ана-  
 1289 логично процедуре `decompose`, посещает те же узлы дерева отрезков, и, следова-  
 1290 тельно, имеет сложность  $O(\log n)$ .  
 1291

```

1292 # действуем аналогично процедуре decompose
1293 def minimum(l, r), v = root):
1294     # обрабатываем пустой отрезок
1295     if r < l:
1296         return None
  
```

```

1297
1298 # если отрезок совпал с отрезком вершины
1299 if (l,r) == v.segment:
1300     # возвращаем минимум на этом отрезке
1301     return A[v.index]
1302
1303 # разбиваем s
1304 s_left = segment_intersect((l,r), v.left.segment)
1305 s_right = segment_intersect((l,r), v.right.segment)
1306
1307 # выполняем рекурсивные вызовы
1308 left = minimum(s_left)
1309 right = minimum(s_right)
1310
1311 # объединяем результаты
1312 if left == None:
1313     return right
1314 if right == None:
1315     return left
1316
1317 # если оба вызова вернули не None
1318 if A[left] < A[right]:
1319     return left
1320 else
1321     return right

```

1322 Теперь опишем процедуру изменения элемента массива. При изменении эле-  
1323 мента нам нужно обновить минимумы во всех вершинах-предках соответствующе-  
1324 го листа, количество таких вершин не больше высоты дерева, т.е.  $O(\log n)$ .

```

1325 # изменяем элемент массива
1326 def change(i, x, v = root):
1327     # если мы пришли в соответствующий лист
1328     if v.segment == (i,i):
1329         v.value = x
1330         return
1331
1332     # рекурсивно спускаемся до соответствующего листа
1333     if i <= v.left.segment[1]:
1334         change(i, x, v.left)
1335     else
1336         change(i, x, v.right)
1337
1338     # обновляем значение в вершине
1339     v.value = min(v.left.value, v.right.value)

```

1340 Таким образом мы научились решать задачу RMQ в динамической постановке  
1341 за  $O(n)$  на построение дерева отрезков и  $O(\log n)$  на каждый запрос. Заметим,  
1342 что наша конструкция не завязана на какие-либо свойства операции минимума,  
1343 кроме ассоциативности. Поэтому аналогичный подход может быть использован,

1344 например, для максимума, суммы, произведения и других ассоциативных опе-  
1345 раций.

### 1346 5.1.2. Статическая задача RMQ

1347 Теперь предположим, что массив для задачи RMQ задан и никогда не меня-  
1348 ется. В этом случае принято говорить о *статической постановке задачи RMQ*. Мы  
1349 уже умеем отвечать на запрос RMQ за  $O(\log n)$  при помощи дерева отрезков. В  
1350 статическое постановке мы можем получить дополнительное ускорение за счёт  
1351 *предобработки* — можно вычислить какие-то дополнительные данные о массиве,  
1352 когда уже известно, что дальше он изменяться не будет, и использовать их  
1353 для более быстрого ответа на запрос. В дальнейшем, говоря о сложности реше-  
1354 ния задачи RMQ, мы будем оперировать парами оценок  $(f(n), g(n))$ , где  $f(n)$  —  
1355 это сложность предобработки, а  $g(n)$  — сложность ответа на запрос RMQ. В этих  
1356 обозначениях сложность прямого вычисления  $(O(1), O(n))$ .

1357 Перед тем как перейти к решению статической задачи RMQ, рассмотрим более  
1358 простую задачу RSQ.

1359 **Определение 5.1.3.** *Задача о сумме на отрезке (range sum query, RSQ):* по после-  
1360 довательности  $A = (a_1, \dots, a_n)$  и паре индексов  $(i, j)$ ,  $i \leq j$ , вычислить сумму  
1361  $S_{i,j} = a_i + \dots + a_j$ .

1362 В статической постановке эту задачу легко решить за  $(O(n), O(1))$ . Давайте вы-  
1363 числим вспомогательные *частичные суммы*

$$1364 \quad S_0 = 0, \quad \forall k \in [n], \quad S_k = a_1 + \dots + a_k.$$

1365 Это легко сделать за  $O(n)$ . Теперь для ответа на запрос  $RSQ(i, j)$  нам достаточно  
1366 вычислить разность соответствующих частичных сумм:

$$1367 \quad RSQ(i, j) = S_j - S_{i-1}.$$

1368 Для RMQ это не работает, т.к. операция минимума не имеет обратной.

1369 **Полное предвычисление.** Начнём с наиболее простого подхода: раз уж мы от-  
1370 дельно учитываем сложность предобработки и сложность запроса, то давайте сде-  
1371 лаем все вычисления на этапе предобработки. Для этого вычислим RMQ для всех  
1372 возможных пар индексов  $(i, j)$  и запишем ответы в таблицу размера  $n \times n$ . Такое  
1373 решение будет иметь сложность  $(O(n^2), O(1))$ . (Если действовать „в лоб“ и запол-  
1374 нять эту таблицу прямым вычислением для каждого запроса, то получится и во-  
1375 все  $(O(n^3), O(1))$ , но нетрудно заметить, что такую таблицу можно заполнить за  
1376  $O(n^2)$ , если действовать немного умнее.)

1377 **Метод разреженной таблицы.** Назовём отрезок *стандартным*, если его дли-  
1378 на равна некоторой степени двойки. На этапе предобработки вычислим ответы  
1379 на запросы RMQ только для стандартных отрезков. Это потребует  $O(n \log n)$  вре-  
1380 мени и памяти (опять же такая оценка достигается, только если заполнять соот-  
1381 ветствующую табличку от меньших отрезков к большим). Как воспользоваться  
1382 этими данными для ответа на запрос RMQ для произвольного отрезка? Восполь-  
1383 зуемся следующим наблюдением.

1384 **Лемма 5.1.2.** Любой отрезок является либо стандартным, либо может быть пред-  
1385 ставлен как объединение двух пересекающихся стандартных отрезков.

1386 *Доказательство.* Пусть отрезок  $[l, r]$  не является стандартным. Выберем  $k$  такое,  
1387 что

$$1388 \quad 2^{k+1} > r - l + 1 > 2^k.$$

1389 Такое  $k$  обязательно существует, иначе отрезок был бы стандартным. Тогда  $[l, r]$   
1390 можно представить как объединение двух стандартных отрезков длины  $2^k$ :

$$1391 \quad [l, r] = [l, l + 2^k - 1] \cup [r - 2^k + 1, r].$$

1392 (Можно проверить, что  $l + 2^k - 1 \geq r - 2^k + 1$ , т.е.  $l - r + 1 \geq 2^{k+1} - 1$ ) □

1393 Теперь для ответа на запрос нам нужно либо вернуть значение из таблицы, ес-  
1394 ли отрезок оказался стандартным, либо представить отрезок в виде объединения  
1395 двух стандартных и взять минимум из минимумов на этих отрезках.

1396 *Замечание 5.1.2.* В методе разреженной таблицы мы пользуемся свойством *идем-*  
1397 *потентности* операции минимума, т.е. тем, что  $\min\{a, a\} = a$ . Это позволяет по-  
1398 лучать ответ как минимум из минимумов на пересекающихся отрезках. Поэтому  
1399 этот подход не будет работать для операций, которые этим свойством не обла-  
1400 дают. Например, метод разреженной таблицы не подходит для задачи RSQ, т.к.  
1401 сумма на отрезке не равна сумме на двух пересекающихся отрезках, которые его  
1402 покрывают. К счастью, для RSQ у нас есть более эффективное решение.

1403 *Замечание 5.1.3.* При оценке сложности этого метода предполагалось, что мы уме-  
1404 ем находить ближайшую степень двойки (подходящее число  $k$  из леммы 5.1.2)  
1405 за  $O(1)$ . Этого можно достичь, заполнив таблицу размера  $n$ , но на практике это  
1406 обычно реализуется при помощи битовых операций, время работы которых в  
1407 этом случае (число  $n$  ограничено размером памяти) можно считать константным.

## 1408 5.2. Задача о наименьшем общем предке

1409 Теперь пришло время поговорить о второй задаче этого раздела.

1410 **Определение 5.2.1.** *Задача о наименьшем/нижайшем общем предке (least/lowest*  
1411 *common ancestor, LCA):* по паре вершин  $(u, v)$  корневого дерева  $T$  найти их самого  
1412 нижнего общего предка. Другими словами, требуется найти наиболее удалённую  
1413 от корня вершину, лежащую одновременно и на пути от  $u$  к корню и на пути от  $v$   
1414 к корню. Поэтому, если  $v$  является предком  $u$ , то  $LCA(u, v) = v$ .

1415 Для решения этой задачи мы применим *сведёние* задачи LCA к задаче RMQ, т.е.  
1416 научимся по условию задачи LCA строить условие для задачи RMQ, а потом из от-  
1417 вета восстанавливать ответ для исходной задачи. Существует множество различ-  
1418 ных видов сведений, некоторые из которых мы изучим, когда будем говорить о  
1419 NP-трудных задачах. Следующее определение довольно неформально, но вполне  
1420 достаточно для наших целей.

1421 **Определение 5.2.2.**  $(O(n), O(1))$ -сведением задачи  $A$  к задаче  $B$  будем называть  
1422 тройку алгоритмов  $(F, G, H)$  таких, что

- 1423 • алгоритм  $F$  по условию для задачи  $A$  размера  $n$  строит условие для задачи  $B$
- 1424 размера  $O(n)$  за время  $O(n)$ ,

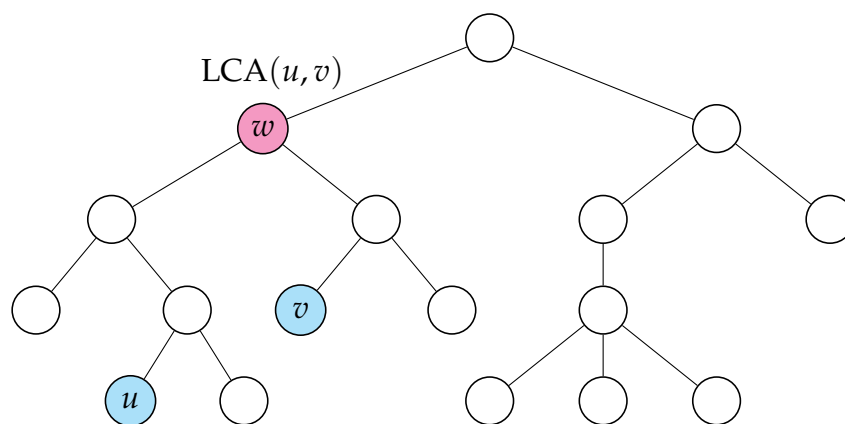


Рис. 5.2.1. Вершина  $w$  является ближайшим общим предком  $u$  и  $v$ .

- 1425 • алгоритм  $G$  по запросу для задачи  $A$  строит соответствующий запрос для
- 1426 задачи  $B$  за  $O(1)$ ,
- 1427 • алгоритм  $H$  по ответу для запроса к задаче  $B$  восстанавливает ответ к исход-
- 1428 ному запросу для задачи  $A$  за время  $O(1)$ .

1429 **Утверждение 5.2.1.** Если существует  $(O(n), O(1))$ -сведение задачи LCA к задаче

1430 RMQ, то задача LCA решается за время  $(O(n \log n), O(1))$ .

1431 *Доказательство.* Для решения задачи LCA на дереве  $T$  с  $n$  вершинами на стадии

1432 предобработки применим алгоритм  $F$ , который по  $T$  построит массив  $A$  размера

1433  $O(n)$  за линейное время. Для массива  $A$  построим разреженную таблицу за время

1434  $O(n \log n)$ . На стадии запросов для каждого запроса LCA будем применять алго-

1435 ритм  $G$ , вычислим полученный RMQ-запрос методом разреженной таблицы, и

1436 в завершение применим алгоритм  $H$ , который вернёт нам ответ для исходного

1437 LCA-запроса. В результате на предобработку мы потратим  $O(n \log n)$ , а каждый

1438 запрос будет обрабатываться за  $O(1)$ . □

1439 **Сведение LCA к RMQ.** Построим сведение LCA к RMQ. Пусть нам дано корневое

1440 дерево  $T$  с  $n$  вершинами. Выполним *эйлеров обход* дерева  $T$ , т.е. рекурсивно обой-

1441 дём дерево, проходя по каждому ребру ровно два раза (см. рис. 5.2.2), и выпишем

1442 для вершины на пути её номер и глубину в дереве.

1443 *# построение эйлерова обхода*

1444 **def** euler-traversal( $v = \text{root}$ ,  $\text{depth} = 0$ ,  $\text{result} = \text{None}$ ):

1445 *# создаём список для хранения результата*

1446 **if**  $\text{result} == \text{None}$ :

1447  $\text{result} = []$

1448

1449 *# добавим в обход текущую вершину*

1450  $\text{result.append}(v, \text{depth})$

1451

1452 *# рекурсивно вызовем для каждого ребёнка и объединим результаты*

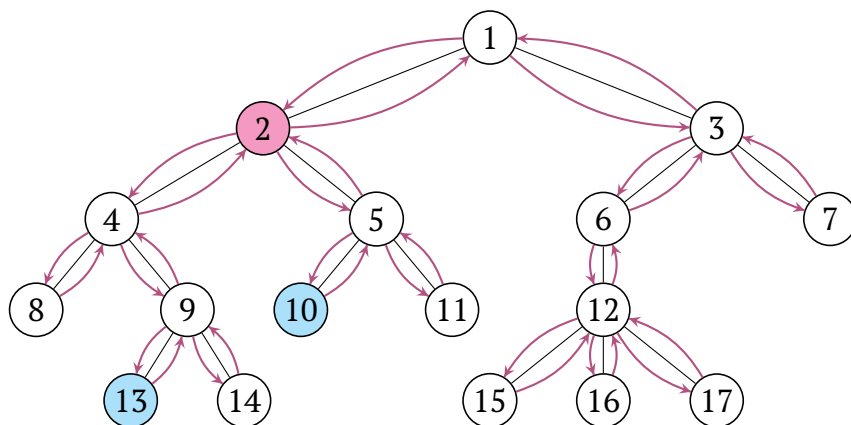
1453 **for**  $c$  **in**  $v.\text{children}$ :

1454  $\text{euler-traversal}(c, \text{depth} + 1, \text{result})$

1455  $\text{result.append}(v, \text{depth})$

1456

1457 **return**  $\text{result}$



вершина	1	2	4	8	4	9	13	9	14	9	4	2	5	10	5	11	...
глубина	0	1	2	3	2	3	4	3	4	3	2	1	2	3	2	3	...

Рис. 5.2.2. Эйлеров обход дерева  $T$  и пример вычисления  $LCA(13, 10)$  через RMQ.

1458 В результате обхода мы получим массив пар длины  $2n - 1$ . Заполним вспомога-  
 1459 тельный массив  $first$  длины  $n$ , где для каждой вершины  $v$  в ячейке с номером  $v$   
 1460 запоним индекс её первого вхождения в эйлеров обход. Теперь по эйлерову об-  
 1461 ходу можно построить задачу RMQ, в которой глубина вершины будет ключом, а  
 1462 номер вершины — дополнительной информацией. Докажем следующую теоре-  
 1463 му.

1464 **Теорема 5.2.1.**  $LCA(v, u)$  равно вершине, на которой достигается  $RMQ(first[v], first[u])$ .

1465 Давайте для начала проверим это утверждение на примере. Сначала найдём  
 1466 на рис. 5.2.2 первое вхождение вершин 13 и 5. Потом на отрезке между этими дву-  
 1467 мя ячейками найдём вершину с минимальной глубиной — это будет вершина 2,  
 1468 и действительно она является ближайшим общим предком вершин 13 и 5. Теперь  
 1469 перейдём к доказательству теоремы.

1470 *Доказательство.* Давайте переформулируем условие теоремы в терминах дере-  
 1471 ва: ближайший общий предок  $u$  и  $v$  — это вершина с минимальной глубиной, ко-  
 1472 торая встречается в эйлеровом обходе между первым вхождением  $u$  и первым  
 1473 вхождением  $v$ . Покажем сначала, что ближайший общий предок обязательно встре-  
 1474 тится между первым вхождением  $u$  и первым вхождением  $v$ . Действительно, в  
 1475 дереве между парой вершин есть только один простой путь, и ближайший об-  
 1476 щий предок всегда лежит на этом пути, поэтому по пути из  $u$  в  $v$  мы обязательно  
 1477 посетим ближайшего общего предка. Теперь покажем, что на этом пути не бу-  
 1478 дет других вершин с меньшей глубиной. Предположим, что такая вершина есть.  
 1479 Следовательно, по пути из  $u$  в  $v$  эйлеров обход прошёл через ближайшего общего  
 1480 предка и поднялся выше, т.е. поднялся по ребру, по которому до этого когда-то  
 1481 спустился. Но ведь для того, чтобы дойти до  $v$ , необходимо будет снова спустить-  
 1482 ся по этому ребру, а такого не бывает, т.к. каждое ребро проходится ровно два  
 1483 раза.  $\square$

1484 Для завершения построения сведения нам осталось описать, как будут устро-  
 1485 ены алгоритмы  $F, G, H$ . Алгоритм  $F$  по дереву  $T$  строит его эйлеров обход и запол-  
 1486 няет массив  $first$ . Алгоритм  $G$  преобразует запрос  $LCA(v, u)$  в запрос  $RMQ(first[v], first[u])$ .



1487 Алгоритм  $H$  возвращает номер вершины, записанный в ячейке массива, на ко-  
1488 торой достигается минимум на соответствующем отрезке. Нетрудно проверить,  
1489 что эти алгоритмы имеют требуемую сложность.

1490 *Замечание 5.2.1.* Получившаяся в результате сведения задача RMQ обладает важ-  
1491 ными дополнительными свойствами, а именно любые два соседних значения от-  
1492 личаются на единицу. Задачи RMQ с таким свойством называют  $\pm 1$ -RMQ. Для  
1493 статической задачи  $\pm 1$ -RMQ существует алгоритм Фарака-Колтона — Бендера [1],  
1494 который решает её за  $(O(n), O(1))$ . Таким образом и задачу LCA можно решить за  
1495  $(O(n), O(1))$  при помощи построенного нами сведения. Но оказывается, что мож-  
1496 но пойти дальше, и построить сведение задачи RMQ к LCA (по массиву для RMQ  
1497 можно построить декартово дерево и на нём запускать алгоритм для LCA). В ре-  
1498 зультате композиции этих двух сведений можно построить сведение задачи RMQ  
1499 к  $\pm 1$ -RMQ, а следовательно решить задачу RMQ в общем случае за  $(O(n), O(1))$ .  
1500 Однако на практике этот алгоритм не используется, т.к. асимптотическое уско-  
1501 рение нивелируется большими константами и сложностью реализации.

1502 *Упражнение 5.2.1.* Постройте  $(O(n), O(1))$ -сведение задачи RMQ к задаче LCA.



## 1503 Глава 6

### 1504 Деревья поиска

#### 1505 6.1. АД с быстрым поиском

1506 **Определение 6.1.1** (множество, set). Абстрактный тип данных *множество* поз-  
1507 воляет хранить набор однотипных элементов и реализует следующие запросы.

- 1508 • `insert(k)` — добавить `k` в множество.
- 1509 • `delete(k)` — удалить `k` из множество.
- 1510 • `find(k)` — найти `k` в множестве.

1511 Если разрешается хранить несколько одинаковых значений, то такой АД назы-  
1512 вается *мультимножеством*.

1513 **Определение 6.1.2** (отображение, словарь, map, dictionary). Абстрактный тип  
1514 данных *отображение* позволяет хранить набор пар “ключ-значение” и реализует  
1515 следующие запросы.

- 1516 • `insert(k, v)` — добавить значение `v` с ключом `k`.
- 1517 • `delete(k)` — удалить запись с ключом `k`.
- 1518 • `find(k)` — найти запись с ключом `k`.

1519 Если разрешается хранить несколько пар с одинаковыми ключами, то такой АД  
1520 называется *мультиотображение*.

1521 Заметим, что имея реализацию множества, мы можем реализовать и отобра-  
1522 жение: будем хранить в множестве пары и сравнивать их только по первому по-  
1523 лю. Поэтому в дальнейшем мы будем говорить только о реализации множества,  
1524 понимая при этом, что все полученные результаты справедливы и для отображе-  
1525 ния.

1526 Нас будет интересовать эффективная реализация множества в случае, если на  
1527 элементах множества задан порядок (т.е. элементы множества можно сравни-  
1528 вать). Для эффективной реализации множества применяются деревья поиска.

1529 *Замечание 6.1.1.* Если множество или отображение не изменяется после началь-  
1530 ного заполнения (т.е. сначала выполняются все запросы `insert()`, а потом посту-  
1531 пают только запросы `find()`), то можно реализовать эти структуры на упорядо-  
1532 ченном массиве: отсортируем массив после заполнения и будем искать элементы  
1533 двоичным поиском.

## 1534 6.2. Представление корневых деревьев в памяти

1535 Перед тем, как говорить о деревьях поиска, полезно обсудить, как можно хра-  
1536 нить корневые деревья в памяти компьютера. Существует два основных подхода.

1537 • **Хранение в массиве.** Такой способ мы применяли для хранения двоичной  
1538 кучи. Этот способ хорошо подходит для хранения двоичных деревьев, кото-  
1539 рые являются полными или почти полными. В этом случае корень хранится  
1540 в ячейке с номером 1, его сыновья в ячейках 2 и 3, и в общем случае для вер-  
1541шины в ячейке  $i$  номера ячеек сыновей и предка вычисляются по формулам:

$$1542 \quad \text{left}(i) = 2i, \quad \text{right}(i) = 2i + 1, \quad \text{parent}(i) = \lfloor i/2 \rfloor.$$

1543 Такой подход можно обобщить на деревья любой ограниченной арности.  
1544 Преимущество данного подхода в отсутствии необходимости хранить ссыл-  
1545 ки и отсутствии фрагментации памяти. Однако, если дерево сильно отлича-  
1546 ется от полного, то в данном представлении будет слишком много пустых  
1547 ячеек, т.е. память будет использовать неэкономно.

1548 • **Хранение в виде набора записей.** При таком способе хранения каждая  
1549 вершина задаётся записью со ссылками на потомков, родителя и другими  
1550 данными. Само дерево задаётся ссылкой на корень. При таком подходе не  
1551 важно, является ли дерево полным, — количество используемой памяти про-  
1552 порционально количеству вершин в дереве. Более того, данный подход поз-  
1553 воляет хранить деревья произвольной арности — для этого в вершинах мож-  
1554 но хранить массив или список потомков. К минусам данного подхода сто-  
1555 ит отнести возможную фрагментацию памяти (зависит от способа хране-  
1556 ния записей). Кроме того, могут потребоваться дополнительные затраты на  
1557 управление памятью.

1558 **Представление «левый ребёнок — правый сосед».** Как мы уже отметили, пред-  
1559 ставление корневого дерева в виде набора записей позволяет хранить деревья  
1560 произвольной арности. Однако, в некоторых случаях желательно, чтобы каждая  
1561 вершина имела фиксированный размер. В таких случаях можно использовать пред-  
1562 ставление «левый ребёнок — правый сосед». В этом случае каждая вершина хра-  
1563 нит две ссылки: на левого ребёнка и на правого соседа. Такое представление, с  
1564 одной стороны, сохраняет фиксированным размер вершины, а с другой стороны  
1565 позволяет эффективно перемещаться по дереву, как если бы мы хранили потом-  
ков вершины в виде списка.

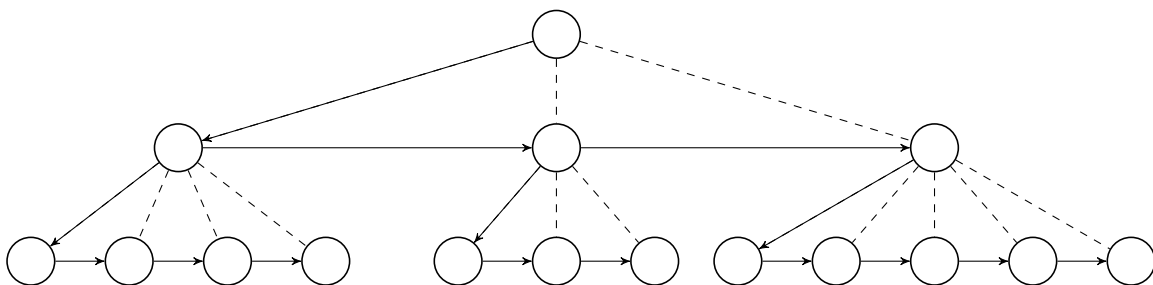


Рис. 6.2.1. Представление «левый ребёнок — правый сосед».

1567 **6.3. Двоичные деревья поиска**

1568 **Определение 6.3.1.** Двоичное дерево поиска — это структура данных, которая пред-  
 1569 ставляет корневое двоичное дерево, состоящее из однотипных узлов. Для значе-  
 1570 ний, которые хранятся в узлах, поддерживается следующий инвариант: элемент  
 1571 в узле не меньше всех элементов в его левом поддереве и не больше всех элемен-  
 1572 тов в его правом поддереве.

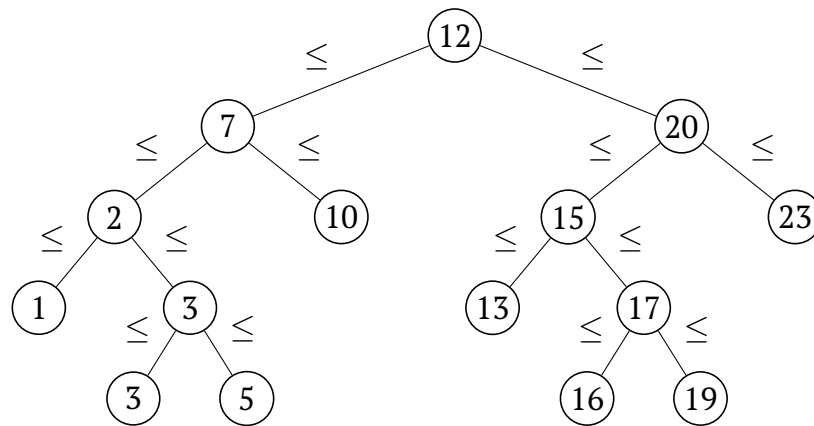


Рис. 6.3.1. Двоичное дерево поиска.

1573 Для реализации в каждом узле будем хранить следующие поля:

- 1574 • `key` — элемент множества,
- 1575 • `left` — левый потомок,
- 1576 • `right` — правый потомок,
- 1577 • `parent` — родитель.

1578 Само же дерево будет задаваться указателем на корень.

1579 **6.3.1. Операции поиска**

1580 Начнём с рекурсивного обхода дерева, при котором элементы будут переби-  
 1581 раться в порядке возрастания (это т.н. *inorder* обхода дерева).

```
1582 # inorder обход дерева
1583 def inorder-traverse(v = root):
1584     if v == None:
1585         return
1586     inorder-traverse(v.left)
1587     print v.key
1588     inorder-traverse(v.right)
```

1589 **Утверждение 6.3.1.** Дерево  $T$  является деревом поиска тогда и только тогда, когда  
 1590 его *inorder* обход печатает ключи в неубывающем порядке.

1591 Тогда поиск элемента можно описать следующим псевдокодом.

```

1592 # Поиск элемента в дереве поиска
1593 def find(k, v = root):
1594     if v == None or k == v.key:
1595         return v
1596     if k < v.key:
1597         return find(k, v.left)
1598     else
1599         return find(k, v.right)

```

1600 Структура дерева поиска позволяет так же легко найти, например, минимальный  
1601 или максимальный элемент.

```

1602 # Поиск минимального элемента
1603 def min_element(v = root):
1604     if v.left == None:
1605         return v
1606     return min_element(v.left)
1607
1608 # Поиск максимального элемента
1609 def max_element(v = root):
1610     if v.right == None:
1611         return v
1612     return max_element(v.right)

```

1613 Заметим, что предыдущие функции возвращают не элемент множества, а узел  
1614 дерева. Имея ссылку на элемент дерева, можно перейти к следующему или преды-  
1615 дущему элементу.

```

1616 # Поиск следующего элемента
1617 def succ(v):
1618     if v.right != None:
1619         return min_element(v.right)
1620     while v.parent != None and v.parent.left != v:
1621         v = v.parent
1622     return v.parent
1623
1624 # Поиск предыдущего элемента
1625 def pred(v):
1626     if v.left != None:
1627         return max_element(v.left)
1628     while v.parent != None and v.parent.right != v:
1629         v = v.parent
1630     return v.parent

```

### 1631 6.3.2. Добавление элемента

1632 Будем всегда добавлять новый элемент в качестве листа дерева. Поэтому нам  
1633 нужно запустить поиск по дереву и найти подходящего родителя, т.е. такой узел,  
1634 у которого нет одного из сыновей и его значение находится в правильном соот-  
1635 ношении с добавляемым элементом.

```

1636 # Добавление нового элемента
1637 # NB: если одинаковые элементы запрещены,
1638 # то наличие элемента нужно проверить отдельно вызовом find
1639 def insert(k):
1640     newnode = node(k)
1641     if root != None:
1642         p = root
1643         while True:
1644             if k < p.key:
1645                 if p.left == None:
1646                     p.left = newnode
1647                     newnode.parent = p
1648                     return newnode
1649                 p = p.left
1650             else # k >= p.key
1651                 if p.right == None:
1652                     p.right = newnode
1653                     newnode.parent = p
1654                     return newnode
1655                 p = p.right
1656     else
1657         root = newnode
1658     return newnode

```

### 1659 6.3.3. Удаление элемента

1660 Заметим, что никаких проблем с удалением листа нет — у листа нет потом-  
 1661 ков, которые нужно куда-то перевешивать. Теперь предположим, что у удаляе-  
 1662 мой вершины нет одного из потомков, например левого. Тогда после удаления  
 этой вершины её место займёт её правый потомок.

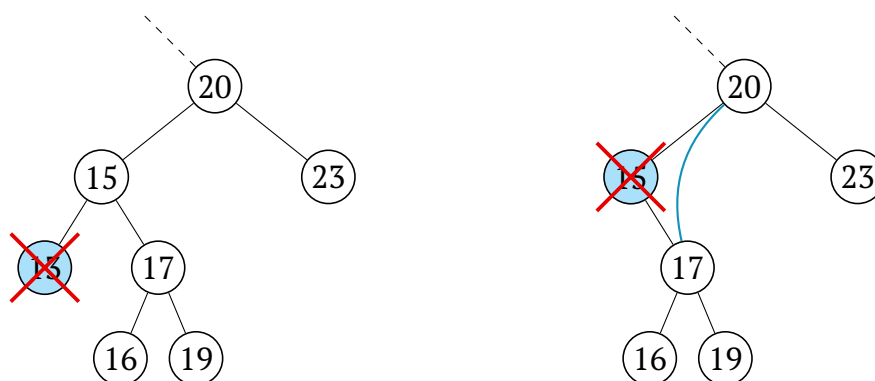


Рис. 6.3.2. Удаление листа и удаление вершины с одним потомком.

1663  
 1664 Осталось понять, что делать в случае удаления вершины, у которой есть оба  
 1665 потомка. Пусть мы хотим удалить некоторую вершину  $v$ , у которой есть оба по-  
 1666 томка. Тогда в поддереве правого потомка обязательно будет вершина  $u = \text{succ}(v)$ ,  
 1667 причём у вершины  $u$  гарантированно не будет левого потомка (следующая за  $v$   
 1668 вершина — это минимальная вершина в поддереве её правого потомка). Обме-  
 1669 няем вершины  $u$  и  $v$  местами и удалим  $v$  одним из описанных выше способов

1670 (теперь у вершины  $v$  нет левого потомка). Нетрудно заметить, что в результате  
1671 дерево останется деревом поиска.

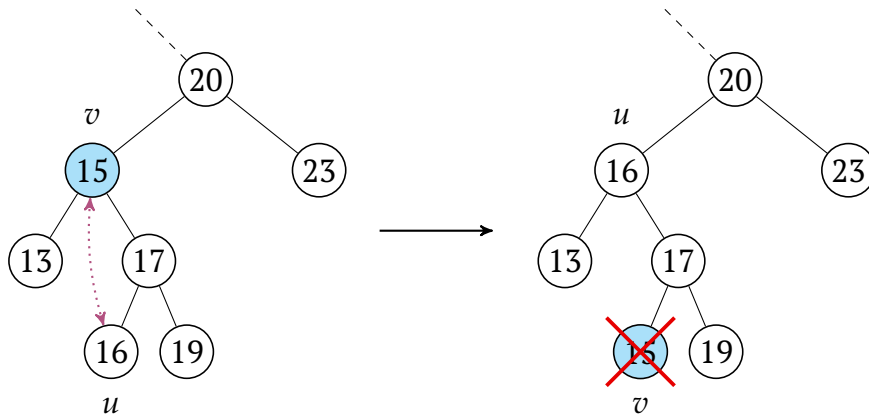


Рис. 6.3.3. Удаление вершины с двумя потомками.

1671

1672 **Упражнение 6.3.1.** Написать соответствующий код для удаление вершины.

1673 **Утверждение 6.3.2.** Все рассмотренные операции с деревом работают за  $O(h)$ , где  
1674  $h$  — высота дерева.

1675 Действительно, все рассмотренные выше операции делают не более одного  
1676 спуска и не более одного подъёма по дереву. При этом в худшем случае высота  
1677 дерева с  $n$  элементами может достигать  $n - 1$  (например, если последовательно  
1678 добавлять в дерево поиска целые числа от 1 до  $n$ , то мы получим „бамбук“ длины  
1679  $n - 1$ ), т.е. операции с таким деревом будут работать за  $O(n)$ .

1680 Для того, чтобы использование дерева поиска было разумным, нам нужно по-  
1681 заботиться о том, чтобы дерево имело высоту поменьше, тогда и операции будут  
1682 выполняться быстрее. Минимальная высота двоичного дерева достигается, если  
1683 оно *полное*, то есть все уровни дерева, кроме последнего, заполнены. Высота пол-  
1684 ного двоичного дерева с  $n$  вершинами равна  $\lceil \log(n + 1) \rceil - 1$ , т.е. время работы  
1685 всех рассмотренных операций будет не более  $O(\log n)$ .

1686 Однако, поддерживать дерево в таком состоянии очень затратно (можно при-  
1687 думать примеры, когда при удалении или добавлении вершины нам придётся пе-  
1688 ревесить множество других вершин, чтобы снова сделать дерево полным). Вме-  
1689 сто полноты двоичного дерева можно накладывать более слабые требования —  
1690 требования *сбалансированности*, которые позволяют получить такую же асимп-  
1691 тотическую оценку  $O(\log n)$  на время работы операций, но при этом можно под-  
1692 держать дерево в сбалансированном состоянии, тратя на это не более  $O(\log n)$   
1693 на каждую операцию. Существует несколько реализаций *сбалансированных де-*  
1694 *реьев поиска*, наиболее известные из которых — *красно-чёрные деревья* и *AVL-*  
1695 *дерево*. Описание красно-чёрных деревьев можно посмотреть в [2] — на практике  
1696 их использую чаще, но эта структура данных значительно более сложна в описа-  
1697 нии, нежели AVL-дерево.

1698 Можно заметить, что оценка  $O(\log n)$  на стоимость *всех* операции для дерева  
1699 поиска — это лучшее, чего можно добиться. Действительно, пусть мы изобрели  
1700 какое-то особенное дерево, для которого все операции работают за  $o(\log n)$ . То-  
1701 гда мы можем построить дерево для произвольного набора элементов  $a_1, \dots, a_n$   
1702 за  $o(n \log n)$  ( $n$  раз добавить вершину за  $o(\log n)$ ). Если мы теперь обойдём это де-  
1703 рево (это можно сделать за  $O(n)$  операций), то таким образом получим элементы



1704  $a_1, \dots, a_n$  в порядке возрастания, т.е. мы научились сортировать произвольные  
 1705 элементы за  $O(n \log n)$ , что противоречит нижней оценке на сортировку сравне-  
 1706 ниями.

## 1707 6.4. AVL-дерево

### 1708 6.4.1. Сбалансированность

1709 *AVL-дерево* [3] придумано советскими учёными Г.М. Адельсоном-Вельским и  
 1710 Е.М. Ландисом. Будем называть двоичное дерево *сбалансированным*, если у лю-  
 1711 бой его вершины высота её левого поддерева отличается от высоты её правого  
 1712 поддерева не более чем на 1.

1713 **Утверждение 6.4.1.** *В сбалансированном дереве высоты  $h$  не более  $2^h - 1$  вершины.*

1714 *Доказательство.* Равенство достигается в случае полного двоичного дерева.  $\square$

1715 **Утверждение 6.4.2.** *В сбалансированном дереве высоты  $h$  не менее  $1.6^h$  вершин.*

1716 *Доказательство.* Пусть  $M_h$  — это минимальное количество вершин в сбаланси-  
 1717 рованном дереве высоты  $h$ . Тогда  $M_h = 1 + M_{h-1} + M_{h-2}$ . Заметим, что это ре-  
 1718 куррентное соотношение очень похоже на соотношения для чисел Фибоначчи:  
 1719  $F_n = F_{n-1} + F_{n-2}$ . Легко проверить, что  $M_h = F_{h+3} - 1$ . Используя оценку для чи-  
 1720 сел Фибоначчи  $F_n \geq ((\sqrt{5} + 1)/2)^n > 1.6^n$ , получаем  $M_h \geq F_{h+3} - 1 \geq 1.6^{h+3} - 1 >$   
 1721  $1.6^h$ .  $\square$

1722 **Следствие 6.4.1.** *Сбалансированное дерево с  $n$  вершинами имеет высоту  $\Theta(\log n)$ .*

1723 *Доказательство.* Прологарифмируем  $2^h > n > 1.6^h$ .  $\square$

1724 Мы показали, что свойства сбалансированности достаточно для того, чтобы  
 1725 высота дерева была логарифмической, а следовательно, операции поиска, добав-  
 1726 ления и удаления вершины работали бы за  $O(\log n)$ . Осталось разобраться, как  
 1727 поддерживать дерево сбалансированным после добавления или удаления верши-  
 1728 ны.

### 1729 6.4.2. Вращения

1730 Для поддержания дерева в сбалансированном состоянии применяются локаль-  
 1731 ные коррекции дерева, называемые *вращениями*. Вращения применяются, если  
 1732 после добавления или удаления вершины, дерево перестало быть сбаланси-  
 1733 рованным. Вращение позволяет восстановить нарушенный баланс и при этом де-  
 1734 рево остаётся деревом поиска.

1735 **Утверждение 6.4.3.** *При добавлении (удалении) вершины в сбалансированное дере-  
 1736 во, баланс мог нарушиться только у тех вершин, которые лежат на пути от корня  
 1737 к добавленной (удалённой) вершине. При этом разница высот поддеревьев в таких  
 1738 вершинах не может быть больше 2.*

1739 Пусть  $x$  — первая вершина на пути от добавленной/удалённой вершины к кор-  
 1740 ню, в которой нарушился баланс. Для восстановления баланса к ней нужно при-  
 1741 менить одно из двух вращений.

1742 **Малые вращения.** Применяются в случае, когда левое поддереву левого по-  
 1743 томка имеет высоту на 2 больше, чем правое поддерево  $x$ , и в симметричном слу-  
 1744 чае. Пусть  $y$  — левый потомок  $x$ . При *малом правом вращении*  $y$  занимает место  $x$ ,  
 1745 а  $x$  в свою очередь становится её правым потомком и забирает одно из поддер-  
 1746 вьев, см. рис. 6.4.1 (высота поддерева  $\beta$  может совпадать с высотой  $\alpha$  или быть на  
 1747 1 меньше). При этом дерево остаётся деревом поиска. Для того, чтобы проверить  
 1748 это, выпишем отношения вершин этого дерева до и после вращения и в обоих  
 случаях мы получим одинаковый порядок  $\alpha \leq y \leq \beta \leq x \leq \gamma$ .

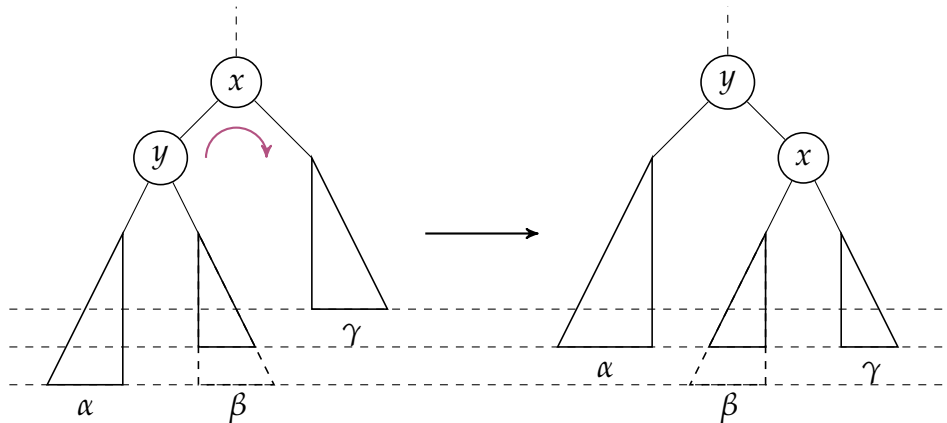


Рис. 6.4.1. Малое правое вращение AVL-дерева.

1749

1750 **Большое вращение.** Большие вращения применяются в случаях, когда малые  
 1751 вращения не могут исправить баланс, т.е. когда правое поддерево левого потомка  
 1752  $x$  имеет высоту на 1 больше, чем высота правого поддерева  $x$ , и в симметричном  
 1753 случае. Большое вращение складывается из двух малых вращений, см. рис. 6.4.2  
 1754 (высота одного из поддеревьев  $\beta$  и  $\gamma$  должна совпадать с высотой  $\alpha$ , а  $y$  второго  
 1755 может быть такой же или на 1 меньше). Можно проверить, что в этом случае од-  
 1756 ним малым вращением не исправить баланс в вершине  $x$ . Большое вращение —  
 1757 это композиция двух малых вращений, каждое из которых сохраняет дерево де-  
 1758 ревом поиска, поэтому в результате большого вращения дерево так же остаётся  
 деревом поиска.

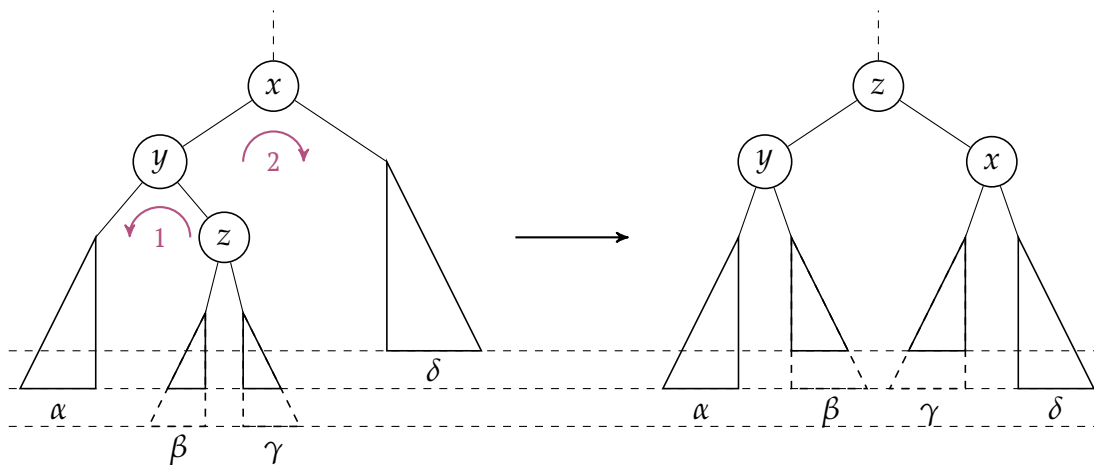


Рис. 6.4.2. Большое правое вращение AVL-дерева.

1759

1760 **Утверждение 6.4.4.** Малое и большое вращения покрывают все возможные случаи  
 1761 нарушения баланса, когда высота поддеревьев отличается на 2.

1762 **Восстановление сбалансированности.** Теперь опишем как происходит вос-  
 1763 становление сбалансированности при помощи вращений. Предположим, что мы  
 1764 добавили или удалили какую-то вершину из дерева. Будем подниматься от этой  
 1765 вершины к корню, пока не встретим какую-нибудь вершину, в которой баланс  
 1766 нарушен. Пусть  $x$  — такая первая вершина. Так как мы добавили или удалили  
 1767 только одну вершину, то разница высот поддеревьев  $x$  не может быть больше 2,  
 1768 а значит, мы можем применить вращение, чтобы восстановить баланс. При этом  
 1769 вращения не могут нарушить баланс ниже по дереву, поэтому если после восста-  
 1770 новления баланса в  $x$  в дереве и остались несбалансированные вершины, то они  
 1771 находятся выше по дереву на пути к корню. Будем продолжать подниматься и  
 1772 исправлять баланс во всех встреченных несбалансированных вершинах.

1773 *Замечание 6.4.1.* В данном рассуждении пропущен важный момент: мы не объ-  
 1774 яснили, почему не может случиться так, что после нескольких вращений баланс  
 1775 в одной из вершин-предков нарушится более чем на 2. Для того, чтобы это объ-  
 1776 яснить, нужно сделать следующие наблюдения.

- 1777 • Вращения могут только уменьшить высоту какого-то поддерева. Поэтому  
 1778 интересующая нас ситуация, когда у какой-то вершины выше баланс нару-  
 1779 шился более чем на 2, могла произойти только при удалении.
- 1780 • При удалении вершины баланс нарушается только в одной вершине (при  
 1781 добавлении может нарушиться в нескольких). И в дальнейшем, каждое вра-  
 1782 щение, которое исправляет это дисбаланс может нарушить баланс только в  
 1783 одной вершине выше, т.к. может только уменьшить высоту.

1784 Из этого можно сделать вывод, что такая ситуация невозможна.

1785 **Оценка времени работы.** Количество вращений не больше, чем вершин на пу-  
 1786 ти от листа к корню, т.е.  $O(\log n)$ . Каждое вращение можно выполнить за  $O(1)$ ,  
 1787 т.к. оно затрагивает константное количество вершин. В сумме получается, что  
 1788 на восстановление баланса после изменяющей дерево операции нам потребует-  
 1789 ся  $O(\log n)$ .

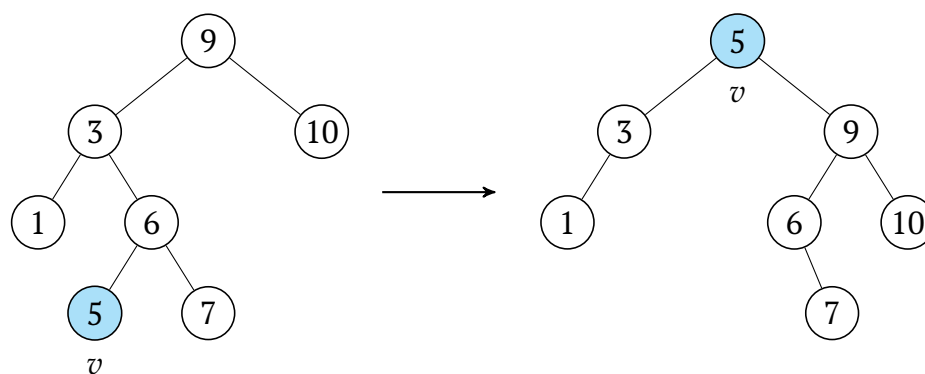
1790 *Замечание 6.4.2.* Для эффективной реализации AVL-дерева в каждой вершине  
 1791 следует хранить и поддерживать высоту её поддерева.

## 1792 6.5. Splay-дерево

1793 *Splay-дерево* — это самобалансирующееся двоичное дерево поиска, изобре-  
 1794 тённое Слеатором и Тарьяном [4, 1] в 1983. В отличие от AVL-дерева, которое  
 1795 обеспечивает оценки *в худшем*, splay-дерево имеет *амортизированную* оценку  $O(\log n)$   
 1796 на время работы основных операций. При этом реализовать splay-дерево немно-  
 1797 го проще, чем AVL-дерево, т.к. не нужно хранить и поддерживать высоту под-  
 1798 дерева в каждой вершине, а также не нужно рассматривать разные случаи при  
 1799 удалении вершины.

### 1800 6.5.1. Запросы к splay-дереву

1801 Ключевую роль в splay-дереве играет операция  $\text{splay}(v)$ , которая перестра-  
 1802 ивает дерево при помощи вращений так, чтобы вершина  $v$  оказалась в корне.

Рис. 6.5.1. Работа операции `splay(v)`.

1803 Более подробно устройство этой операции мы рассмотрим позже, а пока с её по-  
 1804 мощью научимся реализовывать другие операции.

1805 Операция `splay` вызывается для каждой вершины, которая была найдена, до-  
 1806 бавлена или удалена.

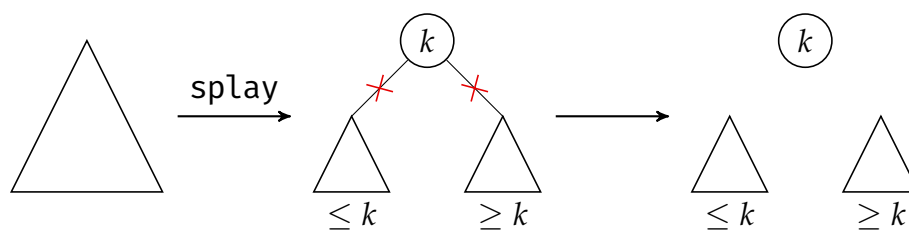
```

1807 # Поиск элемента в splay-дереве
1808 def splay-find(k, v = root):
1809     if k == v.key:
1810         splay(v)
1811         return v
1812     if k < v.key:
1813         if v.left == None:
1814             splay(v)
1815             return None
1816         return splay-find(k, v.left)
1817     else
1818         if v.right == None:
1819             splay(v)
1820             return None
1821         return splay-find(k, v.right)
1822
1823 # Добавление вершины в splay-дереве
1824 def splay-insert(k):
1825     # вызов insert для обычного дерева поиска
1826     v = insert(k)
1827     splay(v)
1828     return v
  
```

1829 Для реализации удаления давайте введём две вспомогательные операции —  
 1830 `split` и `merge`. Операция `split(k)` дерево на три части: дерево с элементами  $\leq k$ ,  
 1831 вершина с ключом  $k$ , дерево с элементами  $\geq k$ .

```

1832 # Разбиение дерева относительно ключа k
1833 # предполагается, что ключ k в дереве есть
1834 def splay-split(k):
1835     v = splay-find(k)
1836     T1 = v.left
1837     T2 = v.right
  
```

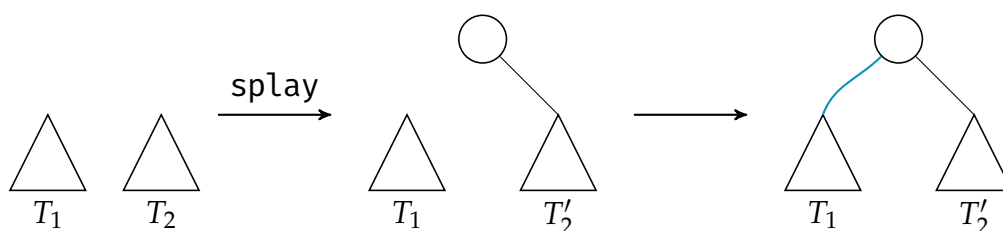
Рис. 6.5.2. Работа операции splay-split( $k$ ).

```

1838 v.left = None
1839 v.right = None
1840 return (T1, v, T2)

```

1841 Операция merge объединяет два дерева в одно (требуется, чтобы все ключи первого дерева были меньше ключей второго дерева).

Рис. 6.5.3. Работа операции splay-merge( $T_1$ ,  $T_2$ ).

```

1842
1843 # Слияние двух деревьев
1844 # ключи T1 меньше ключей T2
1845 def splay-merge(T1, T2):
1846     v = min_element(T2)
1847     splay(v)
1848     v.left = T1
1849     return v

```

1850 Теперь можно описать операцию удаления вершины.

```

1851 # Удаление вершины с ключом k
1852 def splay-remove(k):
1853     (T1, v, T2) = splay-split(k)
1854     return splay-merge(T1, T2)

```

1855 Как мы увидим далее, сложность операции splay ограничена  $O(h)$ , где  $h$  —  
1856 высота дерева, соответственно, верно следующее утверждение.

1857 **Утверждение 6.5.1.** Все рассмотренные операции имеют сложность не более  $O(h)$ .

## 1858 6.5.2. Операция splay

1859 Мы научились реализовывать все необходимые операции со splay-деревом,  
1860 кроме самой операции splay, которая используется в каждой из рассмотренных  
1861 операций как подпроцедура. Теперь мы рассмотрим устройство операции splay  
1862 и докажем амортизированную оценку  $O(\log n)$  на её сложность. Начнём со следу-  
1863 ющих двух утверждений.

1864 **Лемма 6.5.1.** *Истинная стоимость  $\text{splay}(v)$  пропорциональна глубине вершины*  
1865  *$v$ .*

1866 Эта лемма сама по себе не является особенно примечательной, т.к. реализация  
1867 операции  $\text{splay}$  „в лоб“ при помощи малых вращений для AVL дерева имела бы  
1868 такую же оценку (каждое малое вращение уменьшает глубину вершины на 1).

1869 Вторая лемма значительно более интересная.

1870 **Лемма 6.5.2.** *Существует целочисленная функция потенциала  $\Phi$ , определённая на*  
1871 *двоичных деревьях и принимающая неотрицательные значения, такая, что*

1872 (a) *для любого дерева с  $n$  вершинами  $\Phi(T) = O(n \log n)$ ,*

1873 (b) *учётная стоимость любой операции  $\text{splay}$  относительно  $\Phi$  не более  $O(\log n)$ ,*

1874 (c) *добавление листа увеличивает потенциал не более чем на  $\log n$ ,*

1875 (d) *сумма потенциалов деревьев, получившихся из  $T$  удалением корня, меньше ис-*  
1876 *ходного потенциала  $T$  не более чем на  $\log n$ ,*

1877 (e) *при соединении одного дерева к корню другого потенциал получившегося де-*  
1878 *рева не более чем на  $\log n$  превосходит сумму потенциалов исходных деревьев.*

1879 Вместе лемма 6.5.1 и пункты (a) и (b) леммы 6.5.2 позволяют сделать следую-  
1880 щее неожиданное заключение. Предположим, что мы начали с некоторого  $\text{splay}$ -  
1881 дерева с  $n$  вершинами и  $m \geq n$  раз выполнили следующую операцию: нашли са-  
1882 мую глубокую вершину и вызвали для неё  $\text{splay}$ . По лемме 6.5.2 суммарная (ис-  
1883 тинная) стоимость этих операций не будет превосходить  $O(m \log n) + O(n \log n) =$   
1884  $O(m \log n)$ . С другой стороны мы ведь каждый раз вызывали  $\text{splay}$  для самой глу-  
1885 бокой вершины, а значит, по лемме 6.5.1 количество операций было пропорцио-  
1886 нально высоте дерева. Получается, что и высота дерева в среднем была  $O(\log n)$ ,  
1887 т.е. за счёт вызовов  $\text{splay}$  дерево самобалансируется и „в среднем“ имеет неболь-  
1888 шую высоту.

1889 Теперь покажем, как можно обобщить это умозаключение на  $m$  произволь-  
1890 ных (из рассмотренных выше) операций со  $\text{splay}$ -деревом. Все рассмотренные  
1891 операции устроены более-менее одинаково: мы сначала спускаемся по дереву  
1892 до некоторой вершины, а потом вызываем для неё  $\text{splay}$ . Рассмотрим подробнее  
1893 операцию  $\text{splay-find}$ : в поиске ключа мы спускаемся от корня вниз пока не най-  
1894 дём подходящую вершину или не убедимся, что её нет, а потом вызываем  $\text{splay}$   
1895 от вершины, на которой поиск завершился. Заметим что если мы спустились на  
1896 глубину  $d$ , то  $\text{splay}$  совершит  $\Theta(d)$  операций, чтобы поднять вершину в корень (по  
1897 лемме 6.5.1). Таким образом, сколько бы операций спуска по дереву мы не сдела-  
1898 ли, примерно столько же операций подъёма по дереву произойдёт внутри  $\text{splay}$   
1899 (тут важно, что если мы не нашли вершину, то мы всё-равно вызываем  $\text{splay}$  по-  
1900 следней вершины в поиске). Получается, что стоимость  $\text{splay-find}$  лишь в кон-  
1901 станту раз больше стоимости вызова  $\text{splay}$  внутри неё. Соответственно, если над  
1902 деревом размера  $n$  выполнить  $m \geq n$  операций, некоторые из которых  $\text{splay}$ , а  
1903 некоторые —  $\text{splay-find}$ , то суммарная сложность получится не более  $O(m \log n)$ ,  
1904 т.е. в среднем  $O(\log n)$  на каждую операцию.

1905 Можно ли перенести это рассуждение и на другие операции? Рассмотрим, на-  
1906 пример,  $\text{splay-insert}$ . Аналогичные рассуждения показывают, что истинная сто-  
1907 имость  $\text{splay-insert}$  лишь в константу раз больше, чем истинная стоимость  $\text{splay}$   
1908 внутри неё. Однако тут мы не учли, что  $\text{splay-insert}$ , в отличие от  $\text{splay-find}$ ,

1909 изменяет дерево, а значит, изменяет и его потенциал, и таким образом учётная  
 1910 стоимость `splay-insert` может оказаться значительно больше, чем учётная сто-  
 1911 имость `splay`. Тут нам нужно воспользоваться пунктом (с) леммы 6.5.2 о том, что  
 1912 добавление листа изменяет потенциал не более чем на  $\log n$ , а значит, учётная  
 1913 стоимость `splay-insert` тоже ограничена  $O(\log n)$ .

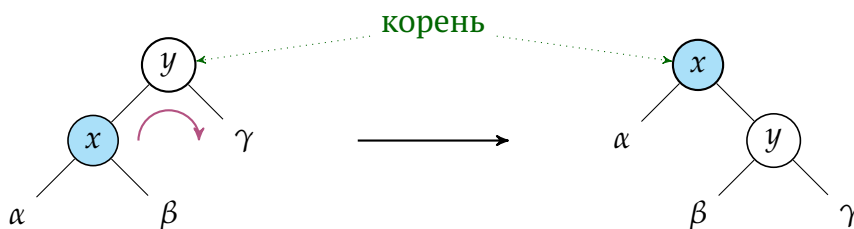
1914 Аналогичные рассуждения позволяют показать, что учётная стоимость опе-  
 1915 рации `splay-remove` не превышает  $O(\log n)$ . Удаление вершины состоит из двух  
 1916 спусков по дереву, из двух вызовов `splay` для соответствующих вершин, удале-  
 1917 нии корня и подсоединении одного дерева к корню другого. Осталось сказать, что  
 1918 сложность спусков не более чем в константу раз превышает сложность соответ-  
 1919 ствующих операций `splay`, а операции удаления корня (пункт (d) леммы 6.5.2) и  
 1920 подсоединения одного дерева к корню другого (пункт (e) леммы 6.5.2) суммарно  
 1921 изменяют потенциал не более чем на  $\log n$ .

1922 Таким образом, если над `splay`-деревом выполнить  $m \geq n$  произвольных опе-  
 1923 раций, в процессе которых в дереве никогда не будет более  $n$  элементов, то сум-  
 1924 марная сложность будет не более  $O(m \log n)$ , т.е. в среднем  $O(\log n)$  на каждую  
 1925 операцию.

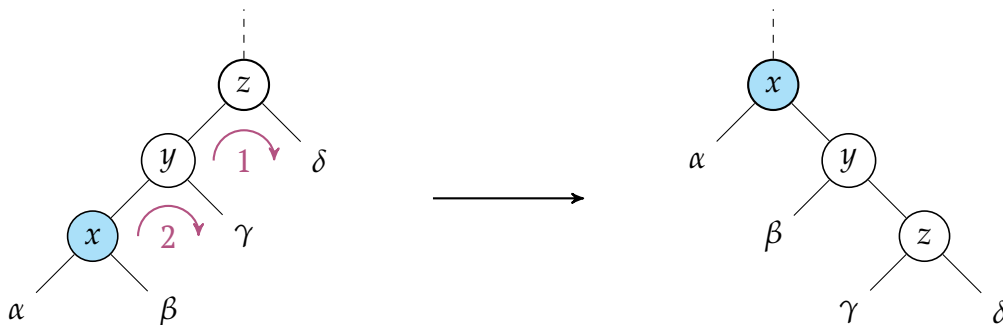
1926 Осталось разобраться в устройстве функции и доказать леммы 6.5.1 и 6.5.2.

1927 **Вращения.** Как уже было сказано, операция `splay` основана на вращениях. Опи-  
 1928 шем её работу для вершины  $x$ . Если  $x$  является корнем, то ничего делать не нуж-  
 1929 но. В противном случае применяется один из трёх шагов-вращений (или один из  
 1930 симметричных им).

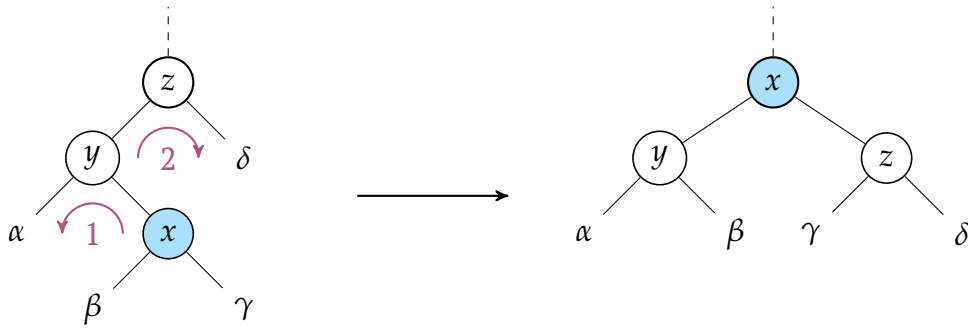
- 1931 • **zig-шаг.** Применяется только в том случае, если предок вершины  $x$  является  
 1932 корнем дерева.



- 1934 • **zigzig-шаг.** Применяется, если  $x$  является левым сыном и левым внуком од-  
 1935 новременно.



- 1937 • **zigzag-шаг.** Применяется, если  $x$  является правым сыном и левым внуком  
 1938 одновременно.



1939

1940 *Замечание 6.5.1.* При реализации zigzig и zigzag-шагов **важен порядок враще-**  
 1941 **ний.** Если бы в zigzig-шаге можно было бы поменять вращения местами, то не  
 1942 было бы необходимости рассматривать zigzig и zigzag-шаги отдельно, ведь они  
 1943 соответствовали бы двум zig-шагам. Но в этом случае получится другая структу-  
 1944 ра данных, про которую лемму 6.5.2 уже доказать не получается.

### 1945 6.5.3. Анализ сложности

1946 Из описания операции splay следует лемма 6.5.1 (количество вращений про-  
 1947 порционально глубине, а каждое вращение требует не более  $O(1)$  операций). До-  
 1948 кажем лемму 6.5.2.

1949 *Доказательство леммы 6.5.2.* Для каждой вершины  $x$  определим *вес*  $w(x)$  равным  
 1950 количеству вершин в поддереве с корнем  $x$ . Теперь мы можем определить функ-  
 1951 цию *потенциала для вершины*  $\Phi(x)$  равной  $\lfloor \log w(x) \rfloor$ . И наконец, определим по-  
 1952 тенциал для всего дерева как сумму потенциалов его вершин.

1953 Потенциал каждой вершины не превосходит  $O(\log n)$ , соответственно, потен-  
 1954 циал дерева не превосходит  $O(n \log n)$ . Таким образом  $\Phi$  удовлетворяет требова-  
 1955 нию пункта (а).

1956 Рассмотрим один вызов функции  $\text{splay}(x)$  для дерева с корнем  $r$ . После вы-  
 1957 зова функции  $\text{splay}(x)$  потенциал вершины  $x$  изменится на  $\Phi_0(r) - \Phi_0(x)$ , где  
 1958  $\Phi_0$  обозначает функцию потенциала до вызова  $\text{splay}$ . Мы покажем, что учётная  
 1959 стоимость этой операции относительно  $\Phi$  не превосходит  $3(\Phi_0(r) - \Phi_0(x)) + 1$ .  
 1960 Поскольку потенциал любой вершины не превосходит  $\log n$ , то и учётная стои-  
 1961 мость  $\text{splay}$  будет ограничена  $3 \log n + 1 = O(\log n)$ .

1962 Разобьём операцию  $\text{splay}$  на шаги в соответствии с её определением и оценим  
 1963 учётную стоимость каждого шага по-отдельности. Пусть  $\Phi$  обозначает функцию  
 1964 потенциала до совершения шага, а  $\Phi'$  — после. Тогда утверждается, что

- 1965 1. учётная стоимость zig-шага не превосходит  $3(\Phi'(x) - \Phi(x)) + 1$ .
- 1966 2. учётная стоимость zigzig-шага не превосходит  $3(\Phi'(x) - \Phi(x))$ .
- 1967 3. учётная стоимость zigzag-шага не превосходит  $3(\Phi'(x) - \Phi(x))$ .

1968 После вызова  $\text{splay}$  вершина  $x$  станет корнем дерева, поэтому суммируя все эти  
 1969 оценки по все шагам мы получаем, что учётная стоимость  $\text{splay}(x)$  не превос-  
 1970 ходит  $3(\Phi_0(r) - \Phi_0(x)) + 1$  (мы пользуемся тем, что в  $\text{splay}$  не более одного zig  
 1971 шага и тем, что потенциал корня не изменяется).

1972 Начнём с оценки zig-шага. Истинная стоимость любого вращения  $O(1)$ . Для  
 1973 простоты будем считать, что она равна 1. Тогда учётная стоимость zig-шага скла-  
 1974 дывается из его истинной стоимости и приращения потенциала дерева, т.е. равна

$$1975 \tilde{c} = 1 + \Phi'(x) + \Phi'(y) - \Phi(x) - \Phi(y).$$



1976 Здесь  $\Phi'(x) = \Phi(y)$  — это потенциал корня, поэтому  $\Phi'(y) \leq \Phi'(x)$ . Подставляем  
1977 и получаем требуемую оценку

$$1978 \quad \tilde{c} = 1 + \Phi'(y) - \Phi(x) \leq 1 + \underbrace{\Phi'(x) - \Phi(x)}_{\geq 0} \leq 1 + 3(\Phi'(x) - \Phi(x)).$$

1979 Теперь докажем оценку для zigzig-шага — для zigzag-шага оценка доказыва-  
1980 ется аналогично. Выпишем учётную стоимость шага:  $\tilde{c} = 1 + \Delta\Phi$ , где

$$1981 \quad \Delta\Phi = \Phi'(x) + \Phi'(y) + \Phi'(z) - \Phi(x) - \Phi(y) - \Phi(z).$$

1982 Воспользуемся тем, что потенциал корня не изменяется, т.е.  $\Phi'(x) = \Phi(z)$ , а так-  
1983 же тем, что

$$1984 \quad \Phi'(z) \leq \Phi'(y) \leq \Phi'(x) \quad \text{и} \quad \Phi(y) \geq \Phi(x). \quad (6.5.1)$$

1985 Получаем, что

$$1986 \quad \Delta\Phi = \Phi'(y) + \Phi'(z) - \Phi(x) - \Phi(y) \leq 2(\Phi'(x) - \Phi(x)).$$

1987 Рассмотрим два случая. Пусть сначала  $\Phi'(x) > \Phi(x)$ , тогда

$$1988 \quad \tilde{c} \leq 1 + 2(\Phi'(x) - \Phi(x)) \leq 3(\Phi'(x) - \Phi(x)) \quad (6.5.2)$$

1989 (здесь мы пользуемся целочисленностью потенциала), и утверждение доказано.

1990 Теперь предположим, что  $\Phi'(x) = \Phi(x)$  (меньше быть не может, т.к. после это-  
1991 го шага  $x$  стал корнем поддерева). В этом случае неравенство (6.5.2) уже невер-  
1992 но. Однако это не значит, что доказываемое утверждение неверно, т.к. неравен-  
1993 ство (6.5.2) получилось применением неравенств (6.5.1), и в этот момент мы мог-  
1994 ли что-то потерять. Действительно, покажем, что в этом случае

$$1995 \quad \Delta\Phi < 0 = 2(\Phi'(x) - \Phi(x)) = 3(\Phi'(x) - \Phi(x)),$$

1996 а следовательно,

$$1997 \quad \tilde{c} = 1 + \Delta\Phi \leq 0 = 3(\Phi'(x) - \Phi(x)),$$

1998 Проведём доказательство от обратного. Предположим, что  $\Phi'(x) = \Phi(x)$  и  $\Delta\Phi =$   
1999  $0$ , т.е. что все неравенства (6.5.1) на самом деле являются равенствами. Тогда по-  
2000 лучается, что все шесть потенциалов  $\Phi'(x)$ ,  $\Phi'(y)$ ,  $\Phi'(z)$ ,  $\Phi(x)$ ,  $\Phi(y)$ ,  $\Phi(z)$  равны меж-  
2001 ду собой и равны какому-то числу  $k$ . Покажем, что так не бывает. Пусть  $w$  и  $w'$  обо-  
2002 значают веса вершин до и после шага. Заметим, что после шага веса поддерева  
2003  $\alpha, \beta, \gamma, \delta$  не изменились — могли измениться только веса вершин  $x, y, z$ . Тогда

$$\begin{aligned} 2004 \quad w'(x) &= w'(\alpha) + w'(\beta) + w'(\gamma) + w'(\delta) + 3 \\ &= (1 + w(\alpha) + w(\beta)) + (1 + w'(\gamma) + w'(\delta)) + 1 \\ &= w(x) + w'(z) + 1 \geq 2^k + 2^k + 1 \geq 2^{k+1}, \end{aligned} \quad (6.5.3)$$

2005 а значит,  $\Phi'(x) \geq k + 1$ , что противоречит нашему предположению.

2006 Осталось доказать пункты (с)–(е).

2007 (с) Добавление листа увеличивает вес всех вершин на пути от корня к листу. У  
2008 скольких из них при этом может увеличиться потенциал? Только у тех, вес  
2009 которых стал равен  $2^k$  для некоторого целого  $k$ . Так как вес вершин на пути  
2010 от листа к корню строго возрастает в промежутке  $0$  до  $n$ , то таких вершин не  
2011 может быть более  $\log n$ .

2012 (d) Удаление корня уменьшает суммарный потенциал деревьев на потенциал  
2013 корня, т.е. не более чем на  $\log n$ .

2014 (e) Подсоединение одного дерева к корню другого увеличивает потенциал кор-  
2015 ня, а он не может превышать  $\log n$ .  $\square$

2016 *Упражнение 6.5.1.* Докажите оценку для zigzag-шага.

2017 *Упражнение 6.5.2.* Проверьте, что если изменить порядок вращений в zigzig-шаге,  
2018 то желаемое утверждение не получится доказать.

2019 Если немного обобщить функцию потенциала в этой лемме, то можно дока-  
2020 зать интересное свойство splay-деревьев.

2021 **Лемма 6.5.3.** Пусть в splay-дереве хранятся  $n$  ключей  $k_1, k_2, \dots, k_n$ . Рассмотрим на-  
2022 бор из  $m \geq n$  запросов splay-find таких, что каждый ключ запрашивается хотя  
2023 бы один раз. Для каждого  $i \in [n]$  обозначим через  $q_i \geq 1$  количество запросов ключа  
2024  $k_i$ . Тогда суммарное количество операций на выполнение всех запросов будет огра-  
2025 ничено  $O(m + \sum_i q_i \log(m/q_i))$ .

*Доказательство.* Давайте посмотрим, как изменится доказательство леммы 6.5.2, если каждой вершине  $x$  назначить некоторый вес  $\rho(x) \geq 1$  и определить  $w(x)$  как сумму весов всех вершин в поддереве с корнем  $x$ . Тогда потенциал  $\Phi(x)$  вершины  $x$  будет ограничен  $O(\log(\sum_{i=1}^n \rho(i)))$ . При таком определении  $w$  неравенство (6.5.3), остаётся верным:

$$\begin{aligned} w'(x) &= w'(\alpha) + w'(\beta) + w'(\gamma) + w'(\delta) + \rho(x) + \rho(y) + \rho(z) \\ &= (\rho(x) + w(\alpha) + w(\beta)) + (\rho(z) + w'(\gamma) + w'(\delta)) + \rho(y) \\ &\geq w(x) + w'(z) + 1 \geq 2^k + 2^k + 1 \geq 2^{k+1}, \end{aligned}$$

2026 Таким образом учётная стоимость операции splay для вершины  $x$  будет не пре-  
2027 восходить

$$2028 \quad 3(\Phi_0(r) - \Phi_0(x)) + 1 \leq 3 \left( \log \sum_{i=1}^n \rho(i) - \log \rho(x) \right) + 1.$$

2029 Отсюда суммарное количество операций на  $m$  запросов не будет превосходить  
2030  $3 \sum_{i=1}^n q_i \cdot (\Phi_0(r) - \Phi(i)) + m$ . Положим теперь  $\rho(i) = q_i$  для всех  $i \in [n]$ . При таком  
2031 выберем весов получаем финальную оценку

$$2032 \quad O \left( m + \sum_{i=1}^n q_i \cdot \log \frac{m}{q_i} \right) = O \left( m + m \cdot H \left( \frac{q_1}{m}, \dots, \frac{q_n}{m} \right) \right).$$

2033  $\square$

2034 **Следствие 6.5.1** (Статическая оптимальность splay-деревя). Суммарное количе-  
2035 ство операций на выполнение  $m$  запросов поиска в splay-дереве не больше, чем в лю-  
2036 бом статическом двоичном дереве поиска.

2037 *Доказательство.* Заметим, что оптимальное количество операций достигается  
2038 на соответствующем дереве кода Хаффмена с частотами  $\frac{q_1}{m}, \frac{q_2}{m}, \dots, \frac{q_n}{m}$ .  $\square$

2039 **6.6. Декартово дерево**2040 **6.6.1. Определение и теорема существования**

2041 **Определение 6.6.1.** *Декартово дерево (cartesian tree)* — это бинарное дерево, в  
 2042 каждой вершине которого хранится пара значений (ключ, приоритет), и которое  
 2043 является одновременно деревом поиска по ключу и кучей по приоритету. Впер-  
 2044 вые рассмотрены Вулемином [5] в 1980.

2045 Совершенно не очевидно, что декартово дерево существует для любого набора  
 2046 пар ключей и приоритетов. Поэтому мы начнём со следующей теоремы.

2047 **Теорема 6.6.1.** *Для любого набора входных пар  $(k_1, p_1), (k_2, p_2), \dots, (k_n, p_n)$  суще-*  
 2048 *ствует декартово дерево.*

2049 *Доказательство.* Мы предложим алгоритм построения декартова дерева. Пред-  
 2050 положим, что пары уже отсортированы по ключу, т.е.  $k_i < k_{i+1}, \forall i \in [n - 1]$  (в про-  
 2051 тивном случае отсортируем их). Теперь найдём пару с максимальным приорите-  
 2052 том, пусть это пара  $(k_i, p_i)$ . Назначим вершину  $(k_i, p_i)$  корнем дерева и рекурсивно  
 2053 вызовем алгоритм на  $(k_1, p_1), \dots, (k_{i-1}, p_{i-1})$  и  $(k_{i+1}, p_{i+1}), \dots, (k_n, p_n)$ , получая та-  
 2054 ким образом левое и правое поддерево соответственно. Заметим теперь, что в  
 2055 корне выполняются одновременно свойство дерева поиска и свойство кучи. Сле-  
 довательно, в результате мы получим декартово дерево (см. 6.6.1).  $\square$

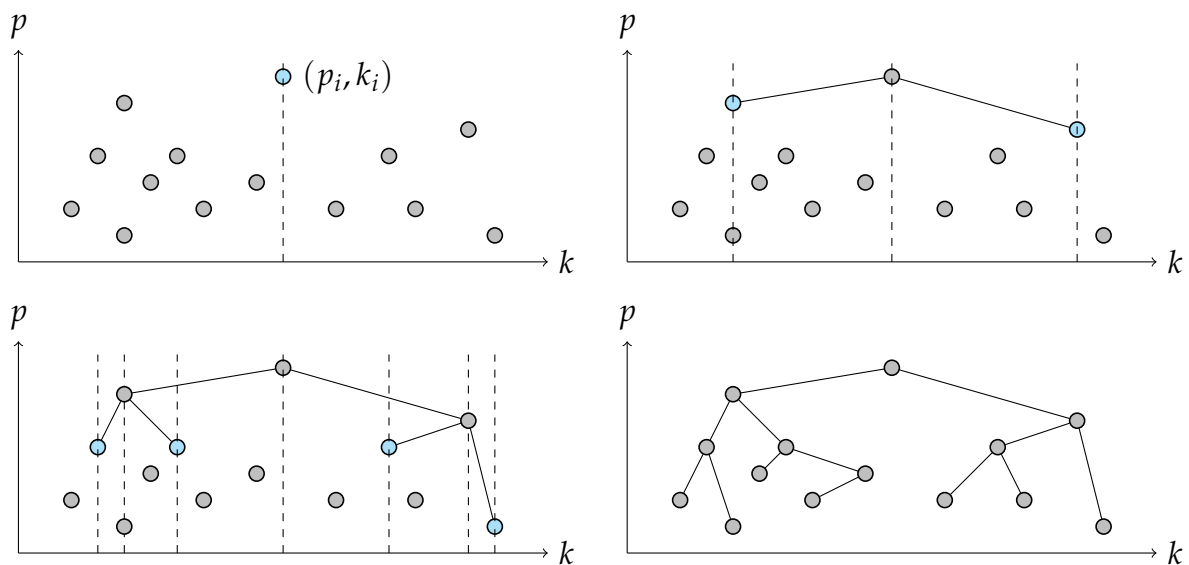


Рис. 6.6.1. Геометрическая интерпретация построения декартова дерева.

2056

2057 Давайте теперь оценим время полученного алгоритма: мы  $n$  раз ищем мак-  
 2058 симум, т.е. получается  $O(n^2)$ . Легко заметить, что этот алгоритм действительно  
 2059 будет работать за квадратичное время, если, например, приоритеты тоже отсор-  
 2060 тированы по возрастанию. Кроме того, этот алгоритм построения позволяет до-  
 2061 казать следующее утверждение.

2062 **Утверждение 6.6.1.** *Если все ключи и приоритеты различны, то декартово дерево*  
 2063 *определено однозначно.*

2064 Действительно, в этом случае на каждой итерации корнем будет выбираться  
 2065 вершина с наибольшим приоритетом, и только эта вершина может быть корнем,  
 2066 иначе нарушится свойство кучи. Аналогично, уникальность ключей гарантиру-  
 2067 ет, что после выбора корня каждая вершина однозначно будет отнесена либо к  
 2068 левой, либо к правой части.

### 2069 6.6.2. Эффективное построение

2070 Нам уже известно, что в общем случае построить дерево поиска быстрее чем  
 2071 за  $\Theta(n \log n)$  не получится. Кроме того, мы умеем строить AVL-дерево и splay-  
 2072 дерево за  $O(n \log n)$ . Для декартова же дерева у нас пока есть только квадратич-  
 2073 ный в худшем случае алгоритм. В этом разделе мы разработаем алгоритм, кото-  
 2074 рый будет работать за линейное время при условии, что входные данные уже от-  
 2075 сортированы по ключу (в противном случае нам придётся предварительно отсор-  
 2076 тировать входные данные, и алгоритм построения будет работать за  $O(n \log n)$ ).

2077 Будем добавлять вершины в дерево в порядке возрастания ключа. Каждая сле-  
 2078 дующая вершина имеет ключ не меньше предыдущих, а следовательно, её мож-  
 2079 но будет добавить самой последней вершиной в самой правой ветке. Однако при  
 2080 этом мы не должны нарушить свойство кучи. Пусть  $v_1, v_2, \dots, v_t$  — это вершины  
 2081 правой ветки в порядке от корня вниз (т.е.  $v_1$  — это корень). Предположим, что мы  
 2082 добавляем вершину  $v = (k_i, p_i)$ . Имеет место один из трёх случаев (см. рис. 6.6.2).  
 2083 Если  $p_i$  меньше приоритета  $v_t$ , то добавим  $v$  правым потомком вершины  $v_t$ . Если,  
 2084 напротив,  $p_i$  больше приоритета корня, то сделаем вершину  $v$  новым корнем, а  
 2085  $v_1$  — её левым потомком. Остался случай, когда  $p_i$  меньше приоритета корня, но  
 2086 больше приоритета  $v_t$ . Найдём пару таких вершин  $v_j$  и  $v_{j+1}$ , что  $p_i$  меньше прио-  
 2087 ритета  $v_j$ , но больше  $v_{j+1}$ . Тогда  $v$  становится правым потомком  $v_j$ , а  $v_{j+1}$  — левым  
 потомком  $v$ .

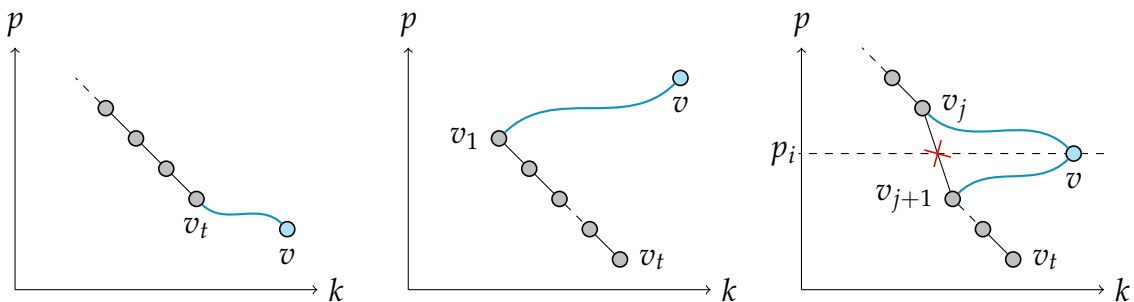


Рис. 6.6.2. Три случая добавления следующей вершины.

2088 Оценим время работы этого алгоритма. Если при добавлении новой вершины  
 2089 мы будем каждый раз перебирать все вершины  $v_1, \dots, v_t$ , то в худшем случае алго-  
 2090 ритм получится квадратичным (действительно, квадратичное время достигается,  
 2091 например, если на вход алгоритму подать следующий набор пар  $(1, n), (2, n -$   
 2092  $1), \dots, (n, 1)$ ). Кажется, что мы ничего не выиграли по сравнению с первым рекур-  
 2093 сивным алгоритмом построения. Однако, если изменить порядок и перебирать  
 2094 вершины  $v_1, \dots, v_t$  с конца (т.е. дополнительно хранить указатель на последнюю  
 2095 вершину в правой ветви и перебирать вершины от последней к корню), то алго-  
 2096 ритм будет иметь линейное(!) время работы.

2098 Давайте разберёмся, почему так происходит. Когда мы добавляем новую вер-  
 2099 шину  $v$  между вершинами  $v_j$  и  $v_{j+1}$  (третий случай на рис. 6.6.2), то вершины

2100  $v_{j+1}, \dots, v_t$  попадают в левое поддереву  $v$  и перестают принадлежать правой ветке. Поэтому вершины  $v_{j+1}, \dots, v_t$  никогда больше не будут встречаться при поиске места для новой вершины. Другими словами, если мы ищем подходящую пару вершин в правой ветке с конца, то каждое просмотренное на этой итерации ребро  $(v_k, v_{k+1})$  больше никогда не будет просмотрено — оно либо будет удалено при добавлении новой вершины, либо перестанет принадлежать правой ветке. При этом на каждой итерации мы добавляем не более одного ребра в правую ветвь. Значит в сумме при поиске подходящей пары вершин мы рассмотрим не более  $n$  рёбер, а значит сложность полученного алгоритма не будет превышать  $O(n)$ .

2109 *Замечание 6.6.1.* В предыдущем параграфе на самом деле описана амортизационная оценка для операции добавления новой вершины с ключом больше, чем все ключи в дереве. Подобные рассуждения будут ещё неоднократно встречаться в курсе. Неформально можно сформулировать следующий принцип: „если на каждой итерации некоторого алгоритма мы увеличиваем некоторую величину  $t$  не более чем на единицу, и совершаем  $k \leq t$  операций, попутно уменьшая  $t$  на  $k$ , то суммарное количество операций не превосходит количества итераций“. Этот неформальный принцип можно доказать «методом ростовщика»: на каждой итерации мы откладываем 1 рубль, а каждая операция тратит 1 рубль, причём нельзя потратить больше, чем к этому моменту отложено. Так как каждая операция в результате оплачена, то количество операций не превосходит количество отложенных рублей, т.е. количества итераций.

2121 *Упражнение 6.6.1.* Докажите этот принцип при помощи метода потенциалов.

### 2122 6.6.3. Операции

2123 Осталось научиться добавлять и удалять вершины. Нам потребуются две вспомогательные операции — `cartesian-split` и `cartesian-merge`.

```
2125 # Разделяет дерево на два по ключу x
2126 def cartesian-split(x, v = root):
2127     if v == None:
2128         return (None, None)
2129
2130     if key(v) < x:
2131         (T1, T2) = cartesian-split(x, v.right)
2132         v.right = T1
2133         return (v, T2)
2134     else:
2135         (T1, T2) = cartesian-split(x, v.left)
2136         v.left = T2
2137         return (T1, v)
```

```
2138 # Слияние двух деревьев
2139 # ключи T1 меньше ключей T2
2140 def cartesian-merge(T1, T2):
2141     if T1 == None:
2142         return T2
2143     if T2 == None:
2144         return T1
```

2145

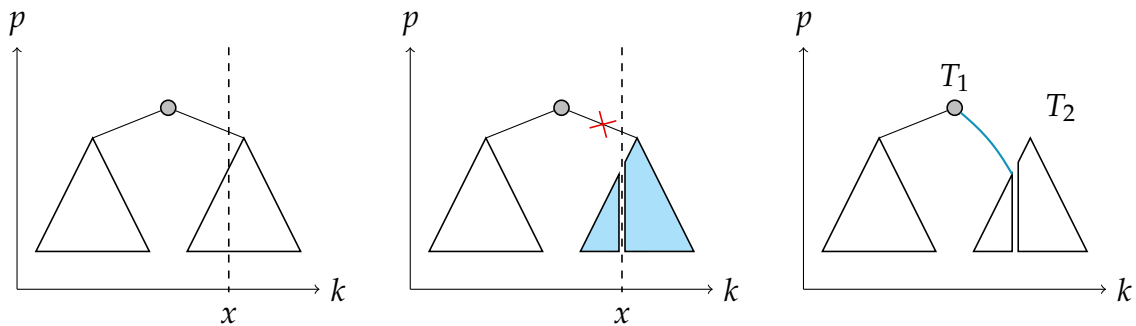


Рис. 6.6.3. Работа операции cartesian-split(k).

```

2146 # корень правого дерева имеет больший приоритет
2147 if pri(T1) < pri(T2):
2148     T = cartesian-merge(T1, T2.left)
2149     T2.left = T
2150     return T2
2151 else
2152     T = cartesian-merge(T1.right, T2)
2153     T1.right = T
2154     return T1

```

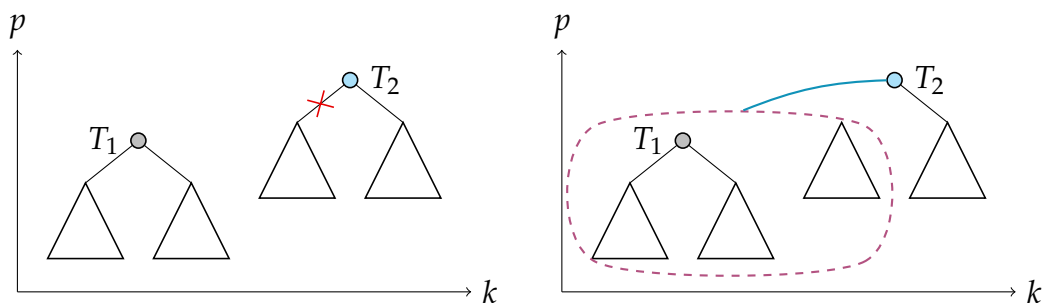


Рис. 6.6.4. Работа операции cartesian-merge(T1, T2).

2155 На основе этих двух операций легко реализовать операции добавления и удале-  
 2156 ния вершин.

```

2157 # Добавление вершины
2158 def cartesian-insert((k,p), T):
2159     v = node(k, p)
2160     (T1, T2) = cartesian-split(k, T)
2161     T1 = cartesian-merge(T1, v)
2162     return cartesian-merge(T1, T2)
2163
2164 # Удаление всех вершин с ключом k
2165 def cartesian-remove(k, T):
2166     (T1, T2) = cartesian-split(k, T)
2167     (T2l, T2r) = cartesian-split(k+1, T2)
2168     return cartesian-merge(T1, T2r)

```

2169 Как и операции для других деревьев, все рассмотренные операции с декарто-  
 2170 вым деревом работают за время  $O(h)$ , где  $h$  — высота дерева. Однако, как мы уже

2171 знаем, у декартова дерева с  $n$  вершинами высота может достигать  $n$ . Кроме того,  
 2172 нет никакой возможности балансировать декартово дерево, т.к. легко построить  
 2173 пример входных данных, на которых декартово дерево определено однозначно  
 2174 и имеет линейную глубину. Что же делать?

#### 2175 6.6.4. Дуча

2176 **Определение 6.6.2.** Дуча (*treap*, *курево*, *дерамиды*) — это дерево поиска постро-  
 2177 енное на декартовом дереве, в котором приоритеты выбираются независимо и  
 2178 случайно.

2179 Идея использовать случайные значения для приоритетов в декартовом дереве  
 2180 появилась только в 1996 году и принадлежит Сиделю и Арагону [6].

2181 **Теорема 6.6.2.** Математическое ожидание высоты дучи ограничено  $O(\log n)$ .

2182 *Доказательство.* Мы здесь воспользуемся теоремой о математическом ожида-  
 2183 нии глубины рекурсии для быстрой сортировки. Какая связь у декартова дерева  
 2184 с быстрой сортировкой? Давайте ещё раз посмотрим на квадратичный алгоритм  
 2185 построения декартова дерева, при помощи которого мы доказали теорему суще-  
 2186 ствования. Предположим, что на вход нам дали только ключи, а приоритеты мы  
 2187 выбираем независимо и случайно. С какой вероятностью каждая вершина будет  
 2188 выбрана корнем? Каждая вершина имеет одинаковую вероятность стать корнем,  
 2189 и это свойство сохраняется для рекурсивных вызовов: на каждом подотрезке, со-  
 2190 ответствующем рекурсивному вызову алгоритма, все вершины имеют одинако-  
 2191 вую вероятность стать корнем соответствующего поддерева. Этот процесс очень  
 2192 похож на то, что происходит в быстрой сортировке: каждая вершина имеет оди-  
 2193 наковую вероятность стать опорным элементом при первом вызове, и это свой-  
 2194 ство сохраняется для рекурсивных вызовов. Эти процессы не просто похожи, а  
 2195 верно и более сильное свойство: для любого фиксированного набора ключей су-  
 2196 ществует взаимнооднозначное соответствие между всеми возможными дучами  
 2197 и деревьями рекурсии быстрой сортировки на этом наборе. Более того, это соот-  
 2198 ветствие сохраняет вероятности (т.е. некоторая конкретная дуча будет иметь ту  
 2199 же вероятность, что и соответствующее ей дерево рекурсии быстрой сортировки).  
 2200 Таким образом оценка на математическое ожидание глубины рекурсии быстрой  
 2201 сортировки влечёт за собой такую же оценку на математическое ожидание высо-  
 2202 ты декартова дерева.  $\square$

2203 **Следствие 6.6.1.** Операции *find*, *cartesian-insert* и *cartesian-remove* для дучи  
 2204 работают за  $O(\log n)$  в среднем.

## 2205 6.7. В-деревья

2206 Другой подход к реализации деревьев поиска — это так называемые *В-деревья*,  
 2207 у которых вершины могут иметь произвольное количество потомков. В-деревья  
 2208 часто используются для поиска по данным, хранящихся во внешней памяти (на-  
 2209 пример, в файловых системах или базах данных): за счёт увеличения арности  
 2210 внутренних вершин, происходит уменьшение глубины дерева, что позволяет умень-  
 2211 шить количество запросов на чтение ко внешней памяти (стоимость которых обыч-  
 2212 но значительно больше, чем стоимость операций с внутренней памятью). Хоро-  
 2213 шо известный всем пример В-дерева — это система каталогов файловой системы:

2214 В каждом каталоге может быть произвольное число файлов (листьев) и подкаталогов (внутренних вершин). В различных реализациях В-деревьев, элементы мо-

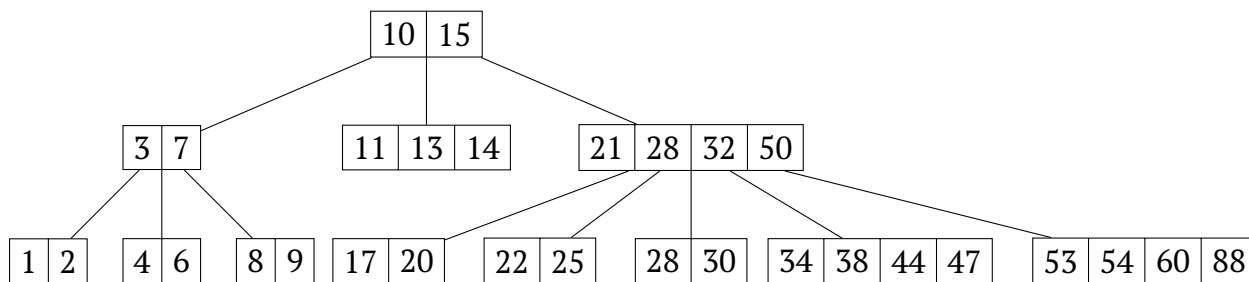


Рис. 6.7.1. В-дерево.

2215

2216 гут храниться как во всех вершинах, так и только в листьях (тогда в вершинах хра-  
 2217 нятся ссылки на максимальные или минимальные элементы в соответствующих  
 2218 поддеревьях). «Компромиссным» вариантом между двоичными деревьями и В-  
 2219 деревьями является *2-3 дерево*, реализующее идею сбалансированного В-дерева,  
 2220 у которого все внутренние вершины имеют два или три потомка.



## 2221 Глава 7

# 2222 Хеширование

2223 В этом разделе мы снова поговорим про реализации абстрактных типов дан-  
2224 ных: множество и словарь (стр. 51). Ранее мы предполагали, что на ключах за-  
2225 дан линейный порядок, и этого было достаточно, чтобы выполнять сортировку и  
2226 применять двоичный поиск или реализовывать дерево поиска. Для разговора о  
2227 хешировании нам потребуется другое предположение — мы будем предполагать,  
2228 что существует функция, сопоставляющая каждому ключу некоторое неотрица-  
2229 тельное целое число.

## 2230 7.1. Прямая адресация

2231 Начнём с ситуации, когда ключи сами по себе являются целыми числами из  
2232 отрезка  $[0, m - 1]$ . Если число  $m$  не очень большое, то мы можем хранить элемен-  
2233 ты в массиве длины  $m$ , или, другими словами, применить *таблицу с прямой ад-*  
2234 *ресацией (direct-address table)*. Для реализации множества в каждой ячейке будем  
2235 хранить булево значение, которое определяет, присутствует элемент в множе-  
2236 стве или нет. Для реализации словаря в каждой ячейке будем хранить ссылку на  
2237 соответствующее значение или пустую ссылку, обозначающую отсутствие ключ-  
2238 ча. При такой реализации используемая память будет пропорциональна  $m$ , а все  
2239 операции поиска, добавления и удаления будут работать за  $O(1)$ . Подобную кон-  
2240 струкцию мы уже встречали в алгоритме сортировки подсчётом — там в каждой  
2241 ячейке хранился счётчик, обозначающий сколько раз элемент встречается в мас-  
2242 сиве. Узким местом данного подхода является избыточное использование памя-  
2243 ти: например, при хранении множества 32-битных ключей, нам в любом случае  
2244 потребуются зарезервировать  $m = 2^{32}$  ячеек памяти вне зависимости от размера  
2245 множества, а для хранения 64-битных ключей данный метод и вовсе неприме-  
2246 ним.

## 2247 7.2. Хеш-таблица

2248 Основной проблемой прямой адресации было то, что общее количество ячеек  
2249  $m$  в таблице никак не зависело от количества элементов, которые в этой таблице  
2250 хранятся. Для эффективного использования памяти мы хотели бы, чтобы размер  
2251 таблицы был не более  $O(n)$ , где  $n$  — количество элементов, которые хранятся в  
2252 таблице. Но как этого достичь, если диапазон ключей очень большой или ключи  
2253 и вовсе не являются целыми числами?

### 2254 7.2.1. Хеш-функция

2255 Мы будем предполагать, что на множестве ключей  $K$  задана некоторая хеш-  
2256 функция (hash function)

$$2257 h : K \rightarrow \{0, \dots, m - 1\}.$$

2258 Значение  $h(k)$  будем называть хеш-кодом (hash code) ключа  $k$ . Более того, мы бу-  
2259 дем предполагать, что мы умеем определять хеш-функции для разных значений  
2260  $m$  и таким образом влиять на количество пустых ячеек в таблице. В дальнейшем  
2261 при анализе сложности мы потребуем, чтобы хеш-функции обладали некоторы-  
2262 ми дополнительными свойствами, которые позволят получить хорошие оценки  
2263 в среднем. Пока же мы будем считать, что „хорошая“ хеш-функция должна хоро-  
2264 шо „перемешивать“ результаты, т.е., например, значение хеш-функции от двух  
2265 близких по значению ключей могут быть произвольно далеки друг от друга.

### 2266 7.2.2. Методы разрешения коллизий

2267 Наличие хеш-функции  $h$  позволяет нам определить хеш-таблицу для ключей  
2268 из множества  $K$ . Заведём массив  $H$  размера  $m$ . В этом массиве каждому ключу  
2269  $k \in K$  будет соответствовать ячейка с номером  $h(k)$ . По принципу Дирихле, ес-  
2270 ли размер множества ключей больше, чем  $m$ , то обязательно найдутся два ключа,  
2271 которым будет соответствовать одна и та же ячейка. Каким образом хранить два  
2272 элемента в одной ячейке? Такая ситуация, при которой  $h(k_1) = h(k_2)$  при  $k_1 \neq k_2$ ,  
2273 называется *коллизией*. Коллизии неизбежно случаются, и поэтому хеш-таблицы  
2274 обязательно имеют *метод разрешения коллизий*. Существуют множество подхо-  
2275 дов к этой проблеме, мы рассмотрим только два наиболее популярных из них.

2276 **Метод цепочек.** Метод цепочек (chaining) предлагает в каждой ячейке таблицы  
2277  $H$  хранить список и записывать в него все ключи, которые попали в эту ячейку.  
При такой реализации метода разрешения коллизий время поиска в ячейке

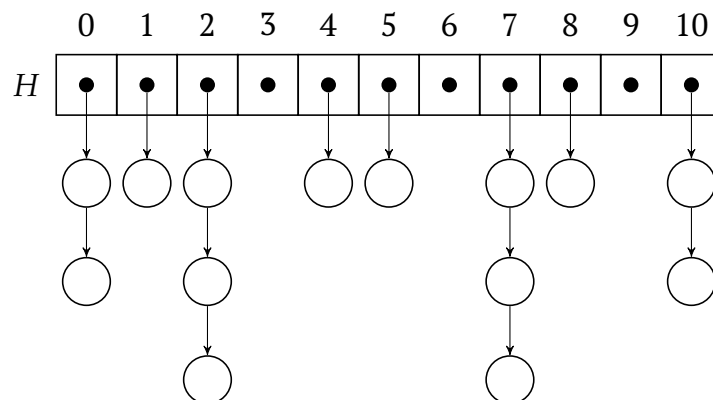


Рис. 7.2.1. Хеш-таблица с разрешения коллизий при помощи метода цепочек,  $m = 11$ .

2278 составляется  $O(l)$ , где  $l$  — длина цепочки. Операция добавления элемента так  
2279 же занимает время  $O(l)$ , если нам нужно проверить, нет ли уже этого ключа в  
2280 списке, и  $O(1)$ , если в хеш-таблице разрешаются повторяющиеся ключи. Удале-  
2281 ние занимает  $O(l)$ . Становится ясно, что в дальнейшем нас будут интересовать  
2282 хеш-функции, которые имеют небольшое количество коллизий, и тем самым в  
2283

2284 хеш-таблице будут не очень длинные цепочки. Как мы увидим дальше, этот ме-  
 2285 тод асимптотически оптимален при правильном выборе хеш-функции. Однако у  
 2286 него есть небольшой недостаток — для хранения списка требуется дополнитель-  
 2287 ное место для хранения ссылок.

2288 **Открытая адресация.** Для полноты картины рассмотрим общую схему мето-  
 2289 дов, альтернативных методу цепочек. При *открытой адресации* в каждой ячейке  
 2290 может храниться только один ключ, поэтому требуется  $n \leq m$ , но каждому ключу  
 2291 соответствует не одна ячейка, а последовательность ячеек. Для этого выбирается  
 2292 хеш-функция от двух аргументов:

$$2293 \quad h : K \times [0, m - 1] \rightarrow [0, m - 1],$$

2294 где второй аргумент задаёт номер в последовательности. Таким образом ключу  $k$   
 2295 соответствуют ячейки с номерами

$$2296 \quad h(k, 0), h(k, 1), \dots, h(k, m - 1).$$

2297 Разумно потребовать, чтобы при фиксированном  $k$  эта последовательность была  
 2298 бы некоторой перестановкой чисел  $\{0, \dots, m - 1\}$ . Для того, чтобы найти ключ  $k$  в  
 2299 хеш-таблице с открытой адресацией, нужно последовательно проверять ячейки  
 2300  $h(k, 0), h(k, 1), \dots$  до тех пор, пока не встретится либо ключ  $k$ , либо пустая ячейка.  
 2301 Добавление ключа реализуется так же: ячейки  $h(k, 0), h(k, 1), \dots$  последовательно  
 2302 проверяются, пока не встретится пустая ячейка. При удалении элемента возни-  
 2303 кают трудности, т.к. удаление может нарушить цепочку, соответствующую неко-  
 2304 торому ключу. Это можно обойти, храня в каждой ячейке дополнительный флаг,  
 2305 обозначающий, что значение из ячейки было удалено — таким образом можно  
 2306 различать пустые ячейки, и ячейки, значение из которых было удалено. Однако  
 2307 при такой реализации длина цепочек не уменьшается даже при удалении эле-  
 2308 ментов, что отрицательно сказывается на производительности.

2309 Такую хеш-функцию с двумя аргументами можно сконструировать из „обыч-  
 2310 ной“ хеш-функции. Самая простая реализация, задающая *линейный поиск* в таб-  
 2311 лице, строит  $h$  по хеш-функции  $h_1$  следующим образом:

$$2312 \quad h(k, i) = (h_1(k) + i) \bmod m.$$

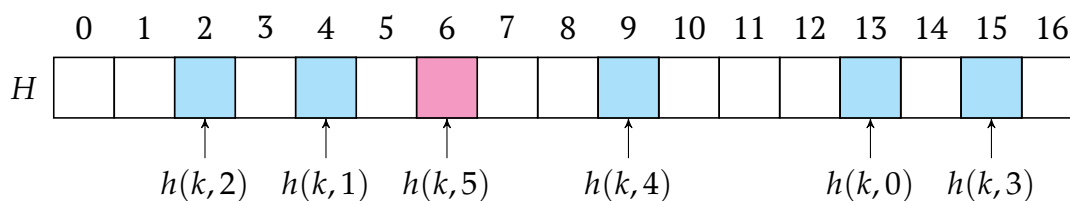
2313 Однако этот способ обладает серьёзным недостатком — цепочки для близких яче-  
 2314 ек таблицы могут „склеиваться“. В более сложном методе *двойного хеширования*  
 2315 функция  $h$  строится на основе двух „обычных“ хеш-функций  $h_1, h_2$  следующим  
 2316 образом:

$$2317 \quad h(k, i) = (h_1(k) + i \cdot h_2(k)) \bmod m.$$

2318 Для этого метода цепочки различных ключей будут „разбросаны“ по таблице с  
 2319 разным шагом, и поэтому не будут склеиваться. Нужно только как-то обеспе-  
 2320 чить, чтобы для фиксированного ключа элементы последовательности не повто-  
 2321 рялись. Для этого удобно выбрать  $m$  простым, и тогда нам подойдёт любая хеш-  
 2322 функция  $h_2(k) : K \rightarrow [1, m - 1]$ .

### 2323 7.2.3. Примеры хеш-функций

2324 В общем случае для того, чтобы выбрать хорошую хеш-функцию, нам нужно  
 2325 что-то знать про распределение ключей. В этом разделе мы рассмотрим некото-  
 2326 рые конструкции, которые могут давать неплохие результаты.

Рис. 7.2.2. Просмотр ячеек при открытой адресации,  $m = 17$ .

- 2327 • **Равномерно распределённые вещественные числа.** Будем считать, что  
 2328 вещественные числа лежат в полуинтервале  $[0, 1)$ . Если это не так, то вос-  
 2329 пользуемся масштабированием. Тогда в качестве хеш-функции можно ис-  
 2330 пользоваться

$$2331 \quad h(k) = \lfloor mk \rfloor.$$

2332 Такая хеш-функция просто разбивает  $[0, 1)$  на  $m$  полуинтервалов одинако-  
 2333 вой длины и сопоставляет каждому числу номер интервала, в который оно  
 2334 попадает.

- 2335 • **Целые числа.** Самый простой вид хеш-функция для целых чисел даёт метод  
 2336 деления с остатком:

$$2337 \quad h(k) = k \bmod m.$$

2338 Более сложный *мультипликативный метод* обладает значительно лучшими  
 2339 перемешивающими свойствами:

$$2340 \quad h(k) = \lfloor m \cdot \{ck\} \rfloor,$$

2341 где  $\{\cdot\}$  обозначает дробную часть числа, а  $c$  — некоторая вещественная кон-  
 2342 станта. Для лучшего результата константу  $c$  нужно выбирать иррациональ-  
 2343 ной, например,  $c = \sqrt{2}$ , тогда дробная часть  $ck$  не равна нулю при любом  $k$   
 2344 (в памяти компьютера можно задать только некоторое рациональное при-  
 2345 ближение иррационального числа, но и этого будет достаточно). Этот метод  
 2346 можно обобщить и на пару значений, выбрав две различные константы  $c_1$  и  
 2347  $c_2$  и определив

$$2348 \quad h(k_1, k_2) = \lfloor m \cdot \{c_1 k_1 + c_2 k_2\} \rfloor.$$

2349 (константы  $c_1$  и  $c_2$  лучше выбирать линейно независимыми над  $\mathbb{Q}$  иррацио-  
 2350 нальными числами, например,  $\sqrt{2}$  и  $\sqrt{3}$ .)

- 2351 • **Последовательности целых чисел, строки.** Для вычисления хеш-кода по-  
 2352 следовательности целых чисел  $a = (a_0, a_1, \dots, a_{n-1})$  хорошие результаты мо-  
 2353 жет дать *полиномиальный метод*: выберем некоторое взаимнопростое с  $m$   
 2354 число  $x$ , называемое *основанием*, и определим

$$2355 \quad h(a) = (a_0 + a_1 \cdot x + a_2 \cdot x^2 + \dots + a_{n-1} \cdot x^{n-1}) \bmod m.$$

- 2356 • **Кортежи разнотипных значений.** Довольно частая задача, которая воз-  
 2357 никает при написании программ, — это придумать хеш-функцию для неко-  
 2358 торой структуры, поля которой имеют разные типы. В этом случае разумно  
 2359 сначала вычислить хеш-коды для каждого из полей, а потом скомбиниро-  
 2360 вать их при помощи мультипликативного или полиномиального методов.

2361 *Замечание 7.2.1.* Вы могли также встречать или даже применять для каких-то за-  
2362 *дач криптографические хеш-функции*, например, такие как MD5 или SHA-1. Как  
2363 следует из названия, такие хеш-функции предназначены для криптографических  
2364 протоколов и не используются в хеш-таблицах. Во-первых потому, что они име-  
2365 ют очень большое пространство хеш-кодов, значительно больше, чем нужно для  
2366 хеш-таблиц (128 битов для MD5 и 160 битов для SHA-1). А во-вторых, так как на  
2367 криптографические хеш-функции накладываются более сложные требования, то  
2368 они довольно сложно устроены, и как следствие этого, долго вычисляются (зна-  
2369 чительно дольше, чем рассмотренные выше конструкции). С другой стороны, в  
2370 некотором приближении можно считать, что схема хранения файлов в репозитори-  
2371 ии системы контроля версий Git — это хеш-таблица, основанная на криптогра-  
2372 фической хеш-функции SHA-256 (если говорить точнее, то там вместо таблицы  
2373 используется дерево каталогов файловой системы, т.е. B-дерево).

2374 В дальнейшем в этом разделе мы для простоты будем предполагать, что слож-  
2375 ность вычисления хеш-функции равна  $O(1)$ . Это верно для чисел, но неверно, на-  
2376 пример, для строк и последовательностей. Поэтому в случаях, когда вычисление  
2377 хеш-функции требует неконстантного времени, в оценке нужно будет дополни-  
2378 тельно учесть, сколько раз тот или иной алгоритм вычисляет хеш-функцию.

#### 2379 7.2.4. Детали реализации

2380 При реализации хеш-таблицы часто заранее неизвестно, сколько ключей в  
2381 ней будет храниться. Поэтому разумно начинать с некоторой небольшой кон-  
2382 станты, а потом расширять таблицу аналогично тому, как мы делали это для рас-  
2383 ширяющегося массива: если коэффициент заполнения  $\alpha = n/m$  стал больше неко-  
2384 торого порога, то нужно создать новую таблицу размера  $\gamma m$  для некоторой кон-  
2385 станты  $\gamma > 1$ , выбрать новую хеш-функцию  $K \rightarrow \{0, \dots, \gamma m - 1\}$  и добавить в  
2386 неё все ключи из старой. При таком подходе будут верны оценки, полученные  
2387 нами для расширяющегося массива, т.е. на расширение таблицы за всё время ра-  
2388 боты мы дополнительно потратим не более  $O(n)$  копирований и  $O(n)$  вычисле-  
2389 ний хеш-функции. Увеличивать таблицу можно не только при слишком большом  
2390 значении  $\alpha$ , но и если длина самой длинной цепочки стала больше некоторого  
2391 порога.

### 2392 7.3. Гипотеза равномерного хеширования

2393 Для удобства изложения введём обозначение для *коэффициента заполнения*  
2394 хеш-таблицы  $\alpha = \frac{n}{m}$ . Как мы увидели ранее, сложность поиска ключа  $k$  в хеш-  
2395 таблице пропорциональна длине цепочки, соответствующей  $k$ . Поэтому жела-  
2396 тельным свойством хеш-функции является „равномерность“ распределения хеш-  
2397 кодов — тогда в среднем длина цепочки будет равна коэффициенту заполнения,  
2398 и, соответственно, сложность поиска в среднем будет ограничена  $O(1 + \alpha)$ . Од-  
2399 нако на самом деле равномерность распределения хеш-кодов зависит не только  
2400 от функции, но и от распределения ключей. Для любой хеш-функции можно по-  
2401 добрать такое распределение на ключах, для которого распределение хеш-кодов  
2402 будет очень далеко от равномерного: выберем некоторое  $j \in \{0, \dots, m - 1\}$  и рас-  
2403 смотрим равномерное распределение на ключах с хеш-кодом  $j$ . Рассуждая в тер-  
2404 минах злонамеренного противника мы могли бы сказать, что противник всегда  
2405 может заставить хеш-таблицу работать медленно — для этого ему нужно записы-

2406 вать в хеш-таблицу только ключи с хеш-кодом  $j$  для некоторого фиксированного  
2407  $j$ .

2408 На практике же иногда у нас могут быть причины предполагать, что на име-  
2409 ющемся распределении ключей выбранная нами хеш-функция задаёт распреде-  
2410 ление хеш-кодов, близкое к равномерному. Давайте оценим сложность операций  
2411 с хеш-таблицей в этом предположении. Для его формализации мы будем поль-  
2412 зоваться следующей гипотезой, которая избавит нас от необходимости делать  
2413 какие-то предположения про распределение ключей.

2414 *Предположение 7.3.1* (Гипотеза равномерного хеширования). Значение хеш-функ-  
2415 ции  $h : K \rightarrow \{0, \dots, m - 1\}$  от ключа  $k$  является случайной величиной, равномерно  
2416 распределённой на множестве  $\{0, \dots, m - 1\}$ . При этом значения хеш-функции на  
2417 различных ключах  $k_1$  и  $k_2$  независимы.

2418 *Замечание 7.3.1.* На практике значение хеш-функции не является случайной ве-  
2419 личиной, т.к. хеш-функции не используют случайных битов. Тем не менее, мы  
2420 могли бы попробовать запрограммировать хеш-функцию, удовлетворяющую этой  
2421 гипотезе. Такая функция могла бы быть устроена следующим образом: каждый  
2422 раз, когда хеш-функция вычисляется от нового ключа  $k$ , значением  $h(k)$  выби-  
2423 рается случайный элемент множества  $\{0, \dots, m - 1\}$ , и эта информация запоми-  
2424 нается. Если такую функцию реализовать на практике, то возникнет проблема с  
2425 тем, что для эффективного хранения информации о выбранных значениях нам  
2426 потребуется хеш-таблица.

2427 **Теорема 7.3.1.** *В предположении гипотезы равномерного хеширования среднее вре-*  
2428 *мя безуспешного поиска при использовании метода цепочек равно  $\Theta(1 + \alpha)$ .*

2429 *Доказательство.* В процессе безуспешного поиска ключа  $k$  мы должны полно-  
2430 стью просмотреть цепочку в ячейке  $h(k)$ . Для каждого  $i \in \{0, \dots, m - 1\}$  обозна-  
2431 чим длину цепочки в ячейке  $i$  как  $l_i$ . По гипотезе равномерного хеширования зна-  
2432 чение  $h(k)$  с равной вероятностью равно каждому из индексов  $\{0, \dots, m - 1\}$ . Сле-  
2433 довательно, математическое ожидание длины цепочки в ячейке  $h(k)$  равно

$$2434 \quad \mathbb{E}[l_{h(k)}] = \sum_{i=0}^{m-1} \frac{1}{m} \cdot l_i = \frac{n}{m} = \alpha.$$

2435 Отсюда среднее время поиска равно  $\Theta(1 + \alpha)$  (единица возникает, т.к. мы всегда  
2436 сделаем хотя бы одну операцию, даже если цепочка пустая.)  $\square$

2437 В предыдущей теореме нам достаточно было применить гипотезу равномер-  
2438 ного хеширования только для ключа, который мы ищем. В следующей теореме  
2439 нам нужно будет применить эту гипотезу ко всем ключам в таблице.

2440 **Теорема 7.3.2.** *В предположении гипотезы равномерного хеширования среднее вре-*  
2441 *мя успешного поиска при использовании метода цепочек равно  $\Theta(1 + \alpha)$ .*

2442 *Замечание 7.3.2.* В данной теореме идёт речь о поиске ключа, который хранит-  
2443 ся в таблице, поэтому среднее время вычисляется не только по вероятностному  
2444 пространству, соответствующему хеш-функции, но и по множеству ключей, на-  
2445 ходящихся в таблице.

2446 *Доказательство.* Пусть в хеш-таблице хранятся ключи  $k_1, k_2, \dots, k_n$ , и пусть ну-  
2447 мерация ключей соответствует порядку, в котором ключи добавлялись в хеш-  
2448 таблицу. Для каждой пары  $i, j \in [n]$  определим случайную величину  $X_{i,j}$ , равную 1,

2449 если случилась коллизия  $h(k_i) = h(k_j)$ , и равную 0 в противном случае. Заметим,  
 2450 что  $E[X_{i,i}] = 1$ , и по гипотезе равномерного хеширования  $E[X_{i,j}] = \frac{1}{m}$  для  $i \neq j$   
 2451 (здесь мы пользуемся независимостью  $h(k_i)$  и  $h(k_j)$ ).

2452 Предположим теперь, что мы ищем ключ  $k_i$ . Сколько элементов нам нужно  
 2453 просмотреть для поиска  $k_i$  в среднем? Будем считать, что при добавлении эле-  
 2454 мента в ячейку он добавляется в начало списка. Тогда количество элементов, ко-  
 2455 торые мы просмотрим при поиске  $k_i$ , равно  $X_{i,i} + \dots + X_{i,n}$ . Тогда среднее время  
 2456 поиска ключа в таблице (среднее по всем ключам присутствующим в таблице и  
 2457 по распределению случайных величин  $\{X_{i,j}\}$ ) равно

$$\begin{aligned} \frac{1}{n} \sum_{i=1}^n \mathbb{E} \left[ \sum_{j=i}^n X_{i,j} \right] &= \frac{1}{n} \sum_{i=1}^n \left( 1 + \sum_{j=i+1}^n \mathbb{E}[X_{i,j}] \right) \\ &= \frac{1}{n} \sum_{i=1}^n \left( 1 + \sum_{j=i+1}^n \frac{1}{m} \right) \\ &= 1 + \frac{1}{mn} \sum_{i=1}^n (n-i) \\ &= 1 + \frac{1}{mn} \left( n^2 - \frac{n(n+1)}{2} \right) \\ &= \Theta(1 + \alpha). \end{aligned}$$

2459

□

## 2460 7.4. Универсальное хеширование

2461 Гипотеза равномерного хеширования позволяет объяснить, почему при «хо-  
 2462 рошем» (т.е., равномерном) распределении хеш-кодов поиск в хеш-таблице бу-  
 2463 дет работать быстро  $O(1 + \alpha)$ . Однако в общем случае не понятно, будет ли рас-  
 2464 пределение хеш-кодов для конкретной хеш-функции действительно равномер-  
 2465 ным, т.к. не известно распределение на ключах. Эту проблему можно решить, ес-  
 2466 ли рассматривать не одну хеш-функцию, а целое семейство хеш-функций с до-  
 2467 полнительным свойством.

2468 **Определение 7.4.1.** Множество функций  $\mathcal{H}$  из множества ключей  $K$  в  $\{0, \dots, m -$   
 2469  $1\}$  называется *универсальным семейством хеш-функций (УСХФ)*, если для любой  
 2470 пары различных ключей  $k_1$  и  $k_2$

$$2471 \Pr_{h \leftarrow \mathcal{H}} [h(k_1) = h(k_2)] \leq \frac{1}{m}.$$

2472 *Замечание 7.4.1.* Отметим, что аналогичный факт верен в предположении гипо-  
 2473 тезы равномерного хеширования, но вероятность определяется по другому ве-  
 2474 роятностному пространству. Если в случае гипотезы равномерного хеширования  
 2475 значение хеш-функции на разных ключах выбираются случайно и независимо,  
 2476 то в случае универсального семейства хеш-функций случайным является только  
 2477 выбор хеш-функции, и она в свою очередь однозначно определяет хеш-коды для  
 2478 всех ключей.

2479 При реализации хеш-таблицы, основанной на универсальном семействе хеш-  
 2480 функций, перед началом заполнения таблицы мы выбираем случайную функ-  
 2481 цию из универсального семейства и используем её для вычисления хеш-кодов.

2482 За счёт этого обеспечивается хорошая в среднем производительность вне зави-  
 2483 симости от распределения ключей, т.к. теперь усреднение происходит не только  
 2484 по распределению на ключах, но и по выбору хеш-функции из семейства. Дру-  
 2485 гими словами, теперь противнику сложно подобрать распределение, на котором  
 2486 хеш-функция будет иметь много коллизий — он может это сделать для каждой  
 2487 хеш-функции из семейства, но он не знает заранее, какую функцию мы выберем.

2488 **Теорема 7.4.1.** При использовании универсального хеширования среднее время без-  
 2489 успешного поиска ключа равно  $\Theta(1 + \alpha)$ .

2490 *Доказательство.* Пусть в хеш-таблице хранятся ключи  $k_1, k_2, \dots, k_n$ . В процессе  
 2491 безуспешного поиска ключа  $k$  мы должны полностью просмотреть цепочку в ячей-  
 2492 ке  $h(k)$ . Для каждого  $i \in [n]$  определим случайную величину  $X_i$ , равную 1, если  
 2493 случилась коллизия  $h(k_i) = h(k)$ , и равную 0 в противном случае. Из определения  
 2494 универсального семейства следует, что  $\mathbb{E}[X_i] \leq \frac{1}{m}$ . Тогда математическое ожида-  
 2495 ние длины цепочки в ячейке  $h(k)$  равно

$$2496 \quad \mathbb{E} \left[ \sum_{i=1}^n X_i \right] = \sum_{i=1}^n \mathbb{E}[X_i] \leq \frac{n}{m} = \alpha.$$

2497 Поэтому среднее время поиска равно  $\Theta(1 + \alpha)$  (единица возникает, т.к. мы всегда  
 2498 сделаем хотя бы одну операцию.)  $\square$

2499 **Теорема 7.4.2.** При использовании универсального хеширования среднее среднее вре-  
 2500 мя успешного поиска ключа равно  $\Theta(1 + \alpha)$ .

2501 *Доказательство.* Доказательство полностью повторяет доказательство предыду-  
 2502 щей теоремы за исключением того, что  $k$  теперь присутствует в таблице. Поэтому  
 2503 один из членов в получившейся сумме будет равен 1, и мы получим

$$2504 \quad \sum_{i=1}^n \mathbb{E}[X_i] \leq 1 + \frac{n-1}{m} \leq 1 + \alpha.$$

2505  $\square$

2506 Получается, что используя универсальное семейство хеш-функций мы с хоро-  
 2507 шей вероятностью можем обеспечить эффективность хеш-таблицы. Однако пока  
 2508 не ясно, существуют ли такие семейства и как легко их построить. Оказывается,  
 2509 что такие семейства есть и построить их не сложно. Мы покажем, как построить  
 2510 такое семейство для целых чисел.

**Теорема 7.4.3.** Пусть  $K = \{0, \dots, n\}$ . Выберем простое  $p > n$ . Тогда семейство  
 хеш-функций  $\mathcal{H}$ , состоящее из функций

$$h_{a,b}(k) = ((ak + b) \bmod p) \bmod m$$

2511 для всех  $a \in \{1, \dots, p-1\}$  и  $b \in \{0, \dots, p-1\}$ , является универсальным.

2512 *Доказательство.* Рассмотрим пару ключей  $k_1 \neq k_2$  и посмотрим, как на них дей-  
 2513 ствует «внутреннее» линейное преобразование по модулю  $p$ . Положим

$$2514 \quad t_1 = (ak_1 + b) \bmod p, \quad \text{и} \quad t_2 = (ak_2 + b) \bmod p.$$



2515 Т.к. ключи  $k_1$  и  $k_2$  не превосходят  $p$ , то  $t_1 \neq t_2$ . Действительно, если предположить,  
2516 что  $t_1 = t_2$ , то из соотношения

$$2517 \quad ak_1 + b = ak_2 + b \pmod{p}$$

2518 получается, что

$$2519 \quad a(k_1 - k_2) = 0 \pmod{p}.$$

2520 Произведение двух чисел меньших  $p$  (тут важно, что  $p$  простое) делится на  $p$  тогда  
2521 и только тогда, когда одно из них равно нулю. Мы выбираем  $a \neq 0$ , следовательно  
2522  $k_1 - k_2 = 0$ , а такого быть не может, т.к. мы рассматриваем различные  $k_1$  и  $k_2$ .

2523 Теперь покажем, что при фиксированных  $k_1$  и  $k_2$  различные пары  $(a, b)$  приво-  
2524 дят к разным парам  $(t_1, t_2)$ . Действительно, при заданных  $k_1$  и  $k_2$  по  $(a, b)$  можно  
2525 вычислить  $(t_1, t_2)$ , а при заданных  $(t_1, t_2)$  можно однозначно восстановить  $(a, b)$ :  
2526 из соотношений

$$2527 \quad \begin{cases} t_1 = (ak_1 + b) \pmod{p}, \\ t_2 = (ak_2 + b) \pmod{p}, \end{cases}$$

2528 получаем  $a = (t_1 - t_2) \cdot (k_1 - k_2)^{-1} \pmod{p}$  (тут обратное берётся по модулю  $p$ ), а  
2529 затем, уже зная  $a$ , вычисляем  $b = (t_1 - ak_1) \pmod{p}$ . Получается, что между парами  
2530  $(a, b)$  и  $(t_1, t_2)$ ,  $t_1 \neq t_2$  есть взаимнооднозначное соответствие. При всех возмож-  
2531 ных  $p(p - 1)$  парах  $(a, b)$  каждая такая пара  $(t_1, t_2)$  встречается равно один раз.  
2532 Т.е. если выбирать функцию  $h_{a,b}$  случайно и равномерно из  $\mathcal{H}$ , то распределе-  
2533 ние соответствующих пар  $(t_1, t_2)$  тоже будет случайно и равномерно на всех парах  
2534  $t_1, t_2 \in \{0, p - 1\}$ ,  $t_1 \neq t_2$ .

2535 Таким образом вероятность того, что  $h_{a,b}(k_1) = h_{a,b}(k_2)$  равна вероятности то-  
2536 го, что случайные различные  $t_1, t_2 \in \{0, p - 1\}$  окажутся равными по модулю  $m$ .  
2537 Для фиксированного  $t_1$  количество  $t_2 \neq t_1$ , которые дают тот же остаток по моду-  
2538 лю  $m$ , не превосходит

$$2539 \quad \left\lceil \frac{p}{m} \right\rceil - 1 \leq \frac{p + m - 1}{m} - 1 = \frac{p - 1}{m}.$$

2540 Каждое  $t_2$  встречается с вероятностью  $1/p$ , следовательно, вероятность получить  
2541 коллизию равна

$$2542 \quad \frac{1}{p} \cdot \frac{p - 1}{m} \leq \frac{1}{m},$$

2543 что и требовалось доказать. □

2544 Используя универсальное хеширование для целых чисел можно построить уни-  
2545 версальное семейство хеш-функций для целочисленных последовательностей и  
2546 строк, основанное на идее полиномиального хеширования. Пусть все элементы  
2547 последовательности являются целыми числами из  $\{0, \dots, d - 1\}$ . Выберем про-  
2548 стое  $p \geq \max(d, m)$  и построим универсальное семейство хеш-функций вида  $\{0, \dots, p -$   
2549  $1\} \rightarrow \{0, \dots, m - 1\}$ , обозначив его  $\mathcal{H}_p$ . Тогда семейство хеш-функций, состоящее  
2550 из всех функций вида:

$$2551 \quad h_x((a_0, a_1, \dots, a_{n-1})) = h((a_0x_0 + a_1x_1 + \dots + a_{n-1}x^{n-1}) \pmod{p}),$$

2552 где  $x \in \{1, \dots, p - 1\}$  и  $h \in \mathcal{H}_p$ , является универсальным.

## 2553 7.5. Совершенное хеширование

2554 В этом разделе мы покажем, что если в хеш-таблице хранится неизменяемое  
 2555 множество элементов (т.е. мы имеем дело со статической задачей поиска), то  
 2556 можно эффективным образом добиться отсутствия коллизий, и тогда сложность  
 2557 доступа будет  $O(1)$  в худшем. Такая конструкция называется *совершенным хеши-*  
 2558 *рованием*. Эту конструкцию можно применять даже в том случае, если множество  
 2559 ключей в таблице изменяется, но множество всех возможных ключей достаточ-  
 2560 но мало (тогда размер таблицы будет пропорционален мощности множества всех  
 2561 возможных ключей).

2562 Пусть размер множества ключей  $K$ , которые мы хотим хранить в таблице, ра-  
 2563 вен  $n$ . Структура совершенной хеш-таблицы состоит из двух уровней (см. рис. 7.5.1).  
 2564 Первый уровень образован хеш-таблицей  $H$  размера  $n$  с хеш-функцией  $h$ . В каж-  
 2565 дой ячейке этой таблицы хранятся ссылка на хеш-таблицу второго уровня и опи-  
 2566 сание хеш-функции для неё, для ячейки с номером  $i$  будем обозначать их  $T_i$  и  $g_i$   
 2567 соответственно. Таким образом, если  $h(k) = i$ , то ключ  $k$  хранится в хеш-таблице  
 2568  $T_i$  в ячейке с номером  $g_i(k)$ . Осталось определить размер хеш-таблиц второго  
 2569 уровня. Для каждого  $i \in \{0, \dots, n-1\}$  через  $n_i$  обозначим количество ключей,  
 для которых  $h(k) = i$ , тогда размер  $T_i$  равен  $n_i^2$ .

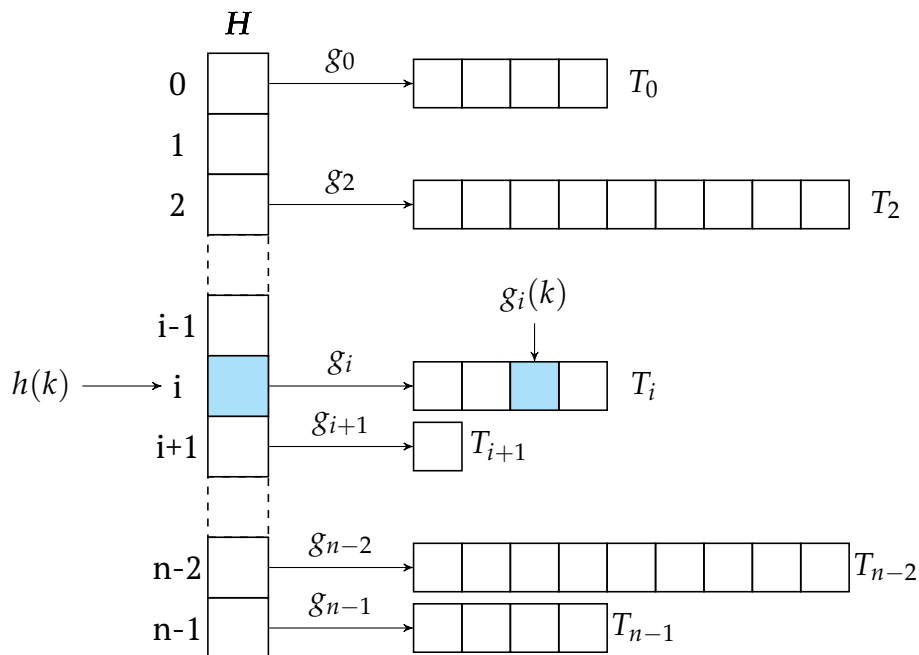


Рис. 7.5.1. Схема совершенного хеширования. Некоторые ячейки хеш-таблицы  $H$  могут оказаться пустыми.

2570 Утверждается, что для фиксированного множества ключей при использова-  
 2571 нии универсального хеширования можно выбрать функции  $h, g_0, g_1, \dots, g_{n-1}$  та-  
 2572 ким образом, чтобы не было коллизий (т.е. в каждой ячейке хеш-таблиц второго  
 2573 уровня хранилось не более одного ключа) и при этом суммарный размер таблиц  
 2574 не превышал  $O(n)$ . Начнём с того, что покажем, как выбрать хеш-функции без  
 2575 коллизий хеш-таблиц второго уровня.  
 2576

2577 **Теорема 7.5.1.** При использовании универсального хеширования для хеш-таблицы  
 2578 размера  $m = n^2$  вероятность возникновения коллизий не превышает  $1/2$ .

2579 *Доказательство.* Если всего  $n$  ключей, то различных неупорядоченных пар  $\binom{n}{2}$ .  
 2580 Вероятность коллизии для каждой пары по определению универсального семей-  
 2581 ства хеш-функций не превосходит  $1/m$ . Пусть случайная величина  $X$  равна коли-  
 2582 честву коллизий. Тогда её математическое ожидание

$$2583 \quad \mathbb{E}[X] \leq \binom{n}{2} \cdot \frac{1}{m} = \frac{n(n-1)}{2} \cdot \frac{1}{m} = \frac{n^2 - n}{2} \cdot \frac{1}{n^2} < \frac{1}{2}.$$

2584 Теперь нужно воспользоваться неравенством Маркова (теорема ??):

$$2585 \quad \Pr[X \geq 1] \leq \frac{\mathbb{E}[X]}{1} < \frac{1}{2}.$$

2586

□

2587 Эта теорема показывает, что как минимум половина хеш-функций из универ-  
 2588 сального семейства не будут иметь коллизий для заданного множества ключей,  
 2589 если размер таблицы равен квадрату количества элементов. Поэтому в среднем  
 2590 для каждой таблицы второго уровня нам нужно проверить две случайные хеш-  
 2591 функции из универсального семейства, чтобы найти такую, которая не имеет  
 2592 коллизий. Остаётся показать, что можно выбрать такую  $h$ , что суммарный размер  
 2593 таблиц второго уровня будет ограничен  $O(n)$ . Вспомним, что  $n_i$  — это количество  
 2594 ключей, которые попадают в ячейку  $i$  таблицы  $H$ . Тогда суммарный размер хеш-  
 2595 таблиц второго уровня равен  $\sum_{i=0}^{n-1} n_i^2$ .

2596 **Теорема 7.5.2.** При использовании универсального хеширования для хеш-таблицы  
 2597 размера  $m = n$

$$2598 \quad \Pr \left[ \sum_{i=0}^{n-1} n_i^2 \geq 4n \right] \leq 1/2.$$

2599 *Доказательство.* Сначала покажем, что

$$2600 \quad \mathbb{E} \left[ \sum_{i=0}^{n-1} n_i^2 \right] < 2n.$$

2601 Для этого воспользуемся соотношением  $n_i^2 = n_i + 2\binom{n_i}{2}$ :

$$2602 \quad \mathbb{E} \left[ \sum_{i=0}^{n-1} n_i^2 \right] = \mathbb{E} \left[ \sum_{i=0}^{n-1} n_i + 2 \sum_{i=0}^{n-1} \binom{n_i}{2} \right] = n + 2 \mathbb{E} \left[ \sum_{i=0}^{n-1} \binom{n_i}{2} \right].$$

2603 Заметим, что  $\sum_{i=0}^{n-1} \binom{n_i}{2}$  — это общее число пар ключей, которые дают коллизию.  
 2604 Поэтому математическое ожидание этой величины можно оценить аналогично  
 2605 предыдущей теореме: для каждой пары ключей вероятность коллизии не превос-  
 2606 ходит  $1/m = 1/n$ , следовательно математическое ожидание общего числа колли-  
 2607 зий не превосходит  $\binom{n}{2} \cdot \frac{1}{n}$ . Таким образом получаем, что

$$2608 \quad \mathbb{E} \left[ \sum_{i=0}^{n-1} n_i^2 \right] \leq n + 2 \binom{n}{2} \frac{1}{n} = n + 2 \cdot \frac{n(n-1)}{2} \cdot \frac{1}{n} = n + n - 1 < 2n.$$

2609 Остаётся применить неравенство Маркова (теорема ??):

$$2610 \quad \Pr \left[ \sum_{i=0}^{n-1} n_i^2 \geq 4n \right] \leq \frac{\mathbb{E} \left[ \sum_{i=0}^{n-1} n_i^2 \right]}{4n} < \frac{2n}{4n} = \frac{1}{2}.$$

2611

□

2612       Теперь можно описать (вероятностный) алгоритм построения таблицы для со-  
2613       вершенного хеширования. Пусть нам дано множество  $n$  ключей. Сначала постро-  
2614       им хеш-таблицу  $H$  и выберем хеш-функцию  $h : K \rightarrow \{0, \dots, n - 1\}$  так, чтобы  
2615       сумма квадратов количества элементов в ячейках  $H$  не превосходила  $4n$ . По тео-  
2616       реме 7.5.2 для этого нам в среднем потребуется проверить две случайные хеш-  
2617       функции из универсального семейства. Когда такая  $h$  найдена, то для каждой  
2618       ячейки  $i$  хеш-таблицы  $H$  мы заводим хеш-таблицу второго уровня  $T_i$  размера  $n_i^2$   
2619       и выбираем для неё хеш-функцию  $g_i : K \rightarrow \{0, \dots, n_i^2 - 1\}$  так, чтобы она не име-  
2620       ла коллизий на множестве ключей, для которых  $h(k) = i$ . По теореме 7.5.1 для  
2621       этого нам тоже в среднем потребуется проверить две случайные хеш-функции  
2622       из универсального семейства. В результате, в среднем предложенный алгоритм  
2623       построения работает за  $O(n)$ .

## 2624 Глава 8

# 2625 Числовые алгоритмы

### 2626 8.1. Арифметические операции с $n$ -битовыми числами

2627 В данном разделе мы обсудим алгоритмы для выполнения арифметических  
2628 операций с целыми числами произвольной длины. Числа будут задаваться их  
2629 представлением в двоичной системе счисления, т.е. можно считать, что число  
2630 задаётся битовым массивом его двоичных цифр. Числа в таком представлении  
2631 легко умножать и делить на два, для этого нужно, соответственно, добавить ноль  
2632 в конец массива или удалить последний элемент массива. Ещё проще вычислить  
2633 остаток при делении на два — остаток записан в последнем элементе массива.  
2634 Для сравнения двух таких чисел нужно сначала сравнить их длины (игнорируя  
2635 ведущие нули, если по какой-то причине они присутствуют в массиве), а если  
2636 длины оказались равны, то просто сравнить эти два массива как строки.

2637 Сложение и вычитание  $n$ -битовых целых легко реализовать за  $O(n)$  используя  
2638 школьный метод сложения в столбик. Умножение мы тоже можем реализовать  
2639 в столбик за  $O(n^2)$ , но значительно для битового представления проще реали-  
2640 зовать умножение в виде следующего рекурсивного алгоритма. В качестве идеи  
2641 алгоритма возьмём следующее соотношение:

$$2642 \quad a \cdot b = \begin{cases} 2(a \cdot \lfloor b/2 \rfloor), & b \text{ — чётное,} \\ 2(a \cdot \lfloor b/2 \rfloor) + a, & b \text{ — нечётное.} \end{cases}$$

```
2643 # Умножение двух битовых чисел
2644 def multiply(a, b):
2645     if b == 0:
2646         return 0
2647
2648     # рекурсивно вычисляем 2(a · ⌊b/2⌋)
2649     c = 2 * multiply(a, b // 2)
2650
2651     # обработка нечётного b
2652     if b % 2 == 1:
2653         c = c + a
2654
2655     return c
```

2656 Проанализируем сложность этого алгоритма, если длина  $a$  составляет  $n$  битов,  
2657 длина  $b$  —  $m$  битов, и  $n \geq m$ . В этом коде проверка на равенство нулю, умно-  
2658 жение и деление на два, а также взятие остатка по модулю два, можно реали-  
2659 зовать за  $O(1)$ , пользуясь тем, что числа заданы битовыми массивами. Остаётся

2660 только сложение  $c$  и  $a$ , которое занимает время пропорциональное сумме длин  
 2661  $|c| + |a| \leq 3n$ . Таким образом каждый вызов `multiply` без учёта вложенных ре-  
 2662 курсивных вызовов выполняется за  $O(n)$ . Глубина рекурсии равна количеству бит  
 2663 в записи  $b$ , поэтому в сумме получается  $O(n \cdot m)$ .

2664 Теперь поговорим о делении. Запрограммировать «школьный» алгоритм де-  
 2665 ления в столбик довольно сложно (попробуйте сами!). Значительно проще ре-  
 2666 ализовать следующий рекурсивный для деления битовых чисел с остатком (идея  
 2667 алгоритма аналогична идее алгоритма умножения, рассмотренного выше).

```

2668 # Деление двух битовых чисел с остатком
2669 # Возвращает пару (частное, остаток)
2670 def divide(a, b):
2671     # база рекурсивного алгоритма
2672     if a < b:
2673         return (0, a)
2674
2675     # рекурсивное вычисляем  $\lfloor a/2 \rfloor / b$  и  $\lfloor a/2 \rfloor \bmod 2$ 
2676     (q, r) = divide(a // 2, b)
2677
2678     # умножаем результаты на 2
2679     q = q * 2
2680     r = r * 2
2681
2682     # если a был нечётным, то к остатку ещё нужно добавить 1
2683     if a % 2 == 1:
2684         r = r + 1
2685
2686     # исправляем результат, если остаток стал больше делителя
2687     if r > b:
2688         q = q + 1
2689         r = r - b
2690
2691     return (q, r)

```

2692 Оценим сложность этого алгоритма для  $n$ -битового  $a$  и  $m$ -битового  $b$ , где  $n \geq m$ .  
 2693 Умножение на два, деление на два, а также взятие остатка по модулю два, ре-  
 2694 ализуются за  $O(1)$ . Сравнение  $a$  и  $b$ , сложение и вычитание можно реализовать за  
 2695  $O(n + m) = O(n)$ , т.е. каждый вызов `divide` без учёта рекурсивных вызовов вы-  
 2696 полняется за  $O(n)$ . Общее количество рекурсивных не более в  $n - m$ , в результате  
 2697 получается  $O(n \cdot (n - m + 1))$  (если  $n = m$ , то алгоритм выполнится за  $O(n)$ ).

2698 Для наших целей будет полезно научиться также вычислять *наибольший общий*  
 2699 *делитель (НОД, greatest common divisor, GCD)* двух чисел. Для этой задачи суще-  
 2700 ствует широко известный алгоритм Евклида, который появился задолго по по-  
 2701 явлению вычислительных машин. Идея алгоритма Евклида содержится в следую-  
 2702 щем соотношении:

$$2703 \quad \text{gcd}(a, b) = \text{gcd}(b, a \bmod b).$$

2704 Действительно, пусть  $a = bq + r$ , где  $r = a \bmod b$ . Если какое-то число делит  $b$  и  $r$ ,  
 2705 то оно обязательно делит  $a$ . Если какое-то число делит  $a$  и  $b$ , то оно также делит  
 2706 и  $r = a - bq$ . Таким образом множество общих делителей пары  $(a, b)$  совпадает с  
 2707 множеством общих делителей  $(b, r)$ , а следовательно наибольший делитель у этих

2708 пар тоже общий. Таким образом можно использовать следующий алгоритм для  
2709 вычисления наибольшего общего делителя двух чисел:

```
2710 # алгоритм Евклида
2711 # вычисление наибольшего общего делителя
2712 def gcd(a,b):
2713     if b == 0:
2714         return a
2715
2716     return gcd(b, a % b)
```

2717 Однако для дальнейших целей на будет полезно научиться также находить так  
2718 называемые *коэффициенты Безу*.

2719 **Теорема 8.1.1** (Теорема Безу). Для любых целых  $a$  и  $b$  из того, что  $\text{gcd}(a,b) = d$ ,  
2720 следует, что существуют такие целые  $x$  и  $y$ , что

$$2721 \quad ax + by = d.$$

2722 Такие коэффициенты  $x$  и  $y$  называются *коэффициентами Безу*. Для нахождения  
2723 коэффициентов Безу используется *расширенный алгоритм Евклида*.

2724 *Замечание 8.1.1.* Уравнения в целых числах называются *Диофантовыми уравне-*  
2725 *ниями*. В теореме Безу мы имеем дело с линейным Диофантовым уравнением от  
2726 двух переменных  $ax + by = d$ . Можно показать, что у такого уравнения будет бес-  
2727 конечное количество решений (тут важно, что  $d = \text{gcd}(a,b)$ ). Алгоритм Евклида  
2728 находит в некотором смысле наименьшее решение, а именно такое, что  $|x| \leq b$  и  
2729  $|y| \leq a$ .

2730 Пусть пусть  $a = bq + r$ , где  $r = a \bmod b$ , и  $\text{gcd}(a,b) = d$ . Предположим, что мы  
2731 нашли некоторые  $x'$  и  $y'$  такие, что

$$2732 \quad bx' + ry' = d.$$

2733 Покажем, как вычислить  $x$  и  $y$ . Подставим  $r = a - bq$ .

$$2734 \quad bx' + (a - bq)y' = d.$$

2735 Перегруппируем члены этого выражения:

$$2736 \quad ay' + b(x' - qy') = d.$$

2737 Получается, что  $x = y'$ , а  $y = x' - qy'$ . Что приводит нас к следующему алгоритму.

```
2738 # расширенный алгоритм Евклида
2739 # вычисление наибольшего общего делителя и коэффициентов Безу
2740 # возвращает тройку из НОД и коэффициентов Безу
2741 def egcd(a, b):
2742     if b == 0:
2743         return (a, 1, 0)
2744
2745     # вычисляем q и r
2746     (q, r) = divide(a, b)
2747
2748     # рекурсивно вычисляем d, x' и y'
```

```

2749 (d, x_, y_) = egcd(b, r)
2750
2751 # вычисляем x и y
2752 x = y_
2753 y = x_ - q * y_
2754
2755 return (d, x, y)

```

2756 Для оценки сложности этого алгоритма заметим, если вызвать алгоритм для  $a <$   
 2757  $b$ , то первый же рекурсивный вызов «переставит» аргументы местами, и во всех  
 2758 следующих вызовах  $a \geq b$ . Если  $a \geq b$ , то  $r < a/2$ , т.е. один из аргументов на каж-  
 2759 дой итерации уменьшается как минимум вдвое. Действительно, если  $a \geq 2b$ , то  
 2760  $r < b \leq a/2$ , если же  $a < 2b$ , то  $r = a - b < a/2$ . Поэтому можно оценить время  
 2761 работы этого алгоритма для двух для  $n$ -битовых как  $O(n^3)$ , т.к. глубина рекурсии  
 2762 ограничена  $O(n)$  и на каждой итерации требуется вычислить  $\text{divide}(a, b)$ , кото-  
 2763 рый мы умеем вычислять за  $O(n^2)$ .

2764 Более аккуратный анализ позволяет получить оценку  $O(n^2)$ . Для этого сначала  
 2765 докажем, следующую лемму.

2766 **Лемма 8.1.1.** Для любых целых  $a$  и  $b$ ,  $a \geq b > 0$ , функция  $\text{egcd}(a, b)$  вернёт тройку  
 2767  $(d, x, y)$ , где  $|x| \leq b$  и  $|y| \leq a$ .

2768 *Доказательство.* Будем доказывать по индукции. В качестве базы рассмотрим  
 2769 случай, когда  $a$  делится на  $b$ . Тогда  $r$  будет равен нулю, а следовательно, рекур-  
 2770 сивный вызов  $\text{egcd}(b, r)$  вернёт  $(b, 1, 0)$ . Получаем, что  $|x| = 0 \leq b$  и  $|y| = 1 \leq a$ ,  
 2771 т.е. для этого базового случая лемма верна.

2772 Теперь рассмотрим случай, когда  $a > b > 0$ , и  $a$  не делится на  $b$ . Предпо-  
 2773 ложим, что для всех  $a' < a$  и  $b' < b$  лемма верна. Положим, как в алгоритме,  
 2774  $q = \lfloor a/b \rfloor$  и  $r = a \bmod b$ . Тогда вложенный в  $\text{egcd}(a, b)$  вызов  $\text{egcd}(b, r)$  вернёт  
 2775 тройку  $(d, x', y')$ , для которой по предположению индукции выполняется  $|x'| \leq r$   
 2776 и  $|y'| \leq b$ . Тогда

$$\begin{aligned}
 |x| &= |y'| \leq r = a \bmod b < b, \\
 |y| &= |x' - qy'| \leq |x'| + |qy'| \leq r + qb = a. \quad \square
 \end{aligned}$$

2778 Теперь используя эту лемму покажем, что  $n$ -битового  $a$  и  $m$ -битового  $b$ , где  
 2779  $n \geq m$ , алгоритм выполняется за  $O(n^2)$  операций. Будем тоже доказывать это  
 2780 по индукции. База индукции: для однобитных чисел достаточно константного  
 2781 количества операций. Предположим, что для любых пар чисел, длины которых не  
 2782 превосходят  $m$ , соответственно, алгоритм выполняется за  $O(m^2)$ . Тогда сложность  
 2783 вычисления  $\text{egcd}(a, b)$  складывается из сложности вызова  $\text{divide}$ , рекурсивного  
 2784 вызова  $\text{egcd}$  и сложности вычисления  $y$ , т.е.  $O(n(n - m + 1))$ ,  $O(m^2)$  и  $O(|x'| + |q| \cdot$   
 2785  $|y'|)$ . Последнее слагаемое по лемме 8.1.1 можно ограничить  $O(n + (n - m) \cdot m)$ .  
 2786 Все три слагаемые не превосходят  $O(n^2)$ , так мы и получаем желаемую оценку.

2787 *Замечание 8.1.2.* В данном алгоритме нам потребуется работа с отрицательными  
 2788 числами. Для представления отрицательных чисел потребуется ещё один допол-  
 2789 нительный бит, отвечающий за знак числа, т.е. для хранения  $n$ -битового числа  
 2790 нужен  $n + 1$  бит. Кроме того, в реализацию арифметических операций нужно бу-  
 2791 дет добавить код с обработкой знака.

2792 Ещё один полезный алгоритм — это алгоритм для возведения в степень. На-  
 2793 ивный алгоритм, который для вычисления числа  $a^b$  в цикле  $b$  умножает на  $a$  (и,



2794 следовательно, требует  $b$  умножений), нельзя назвать эффективным. Дело в том,  
 2795 что значение  $n$ -битового числа может достигать  $2^n$ , а значит время предложен-  
 2796 ного алгоритма экспоненциально от длины числа  $b$ . Более эффективный алгоритм  
 2797 можно построить, если действовать аналогично алгоритму умножения, рассмот-  
 2798 ренному выше.

```

2799 # Быстрое возведение в степень
2800 def power(a, b):
2801     if b == 0:
2802         return 1
2803
2804     # рекурсивно вычисляем  $a^{\lfloor b/2 \rfloor}$ 
2805     res = power(a, b // 2)
2806
2807     # возводим результат в квадрат
2808     res = res * res
2809
2810     # обработка нечётного  $b$ 
2811     if b % 2 == 1:
2812         res = res * a
2813
2814     return res
  
```

2815 Давайте проанализируем время работы этого алгоритма для пары  $n$ -битовых чи-  
 2816 сел  $a$  и  $b$ . Глубина рекурсии не превосходит  $n$ , а на каждом уровне рекурсии вы-  
 2817 числяется три умножения. В результате получаем оценку  $O(n^3)$ , что значительно  
 2818 лучше экспоненциальной оценки наивного алгоритма.

2819 К сожалению, в предыдущем рассуждении мы совершили очень грубую ошиб-  
 2820 ку. Мы предположили, что все умножения будут стоить не более  $O(n^2)$ , но эта  
 2821 оценка верна только для  $n$ -битовых чисел. Давайте проанализируем, какой дли-  
 2822 ны будут числа, которые получаются в данном алгоритме. Для простоты будем  
 2823 считать, что  $a = b = 2^n$ . Сколько бит в числе  $a^b$ ? Количество битов в двоичной за-  
 2824 писи числа с точностью до константы равно двоичному логарифму этого числа. В  
 2825 нашем случае  $\log_2 a^b = \log_2 2^{n2^n} = n2^n$ . Таким образом на верхнем уровне рекур-  
 2826 сии будут умножаться два  $n2^{n-1}$ -битных числа. Соответственно, сложность этого  
 2827 алгоритма уж точно не меньше  $\Omega(n^2 2^{2n})$ , т.е. алгоритм имеет экспоненциальную  
 2828 сложность. Эта ошибка показывает, насколько внимательным нужно быть при  
 2829 оценке сложности алгоритмов.

## 2830 8.2. Модульная арифметика

2831 *Модульная арифметика* — это работа с целыми числами в кольце остатков по  
 2832 некоторому модулю  $m$ , т.е. все числа и все результаты операций берутся по моду-  
 2833 лю  $m$ . Реализовать сложение двух чисел по модулю очень просто — нужно вычис-  
 2834 лить результат как с обычными целыми числами, а от результата взять остаток  
 2835 по модулю  $m$ .

```

2836 # Сложение двух чисел по модулю
2837 def plus_mod(a, b, m):
2838     return (a + b) % m
  
```

```

2839
2840 # Умножение двух чисел по модулю
2841 def multiply_mod(a, b, m):
2842     return (a * b) % m

```

2843 Пусть  $m$  —  $n$ -битное число. Тогда оба алгоритма имеют сложность  $O(n^2)$ .

2844 Осталось разобраться с делением по модулю. Деление по модулю — это де-  
 2845 ление в кольце остатков, оно работает не так, как деление обычных целых чисел.  
 2846 Для того, чтобы поделить  $a$  на  $b$  по модулю  $m$ , нужно найти такое число  $q \in [m - 1]$ ,  
 2847 что

$$2848 \quad bq = 1 \pmod{m}.$$

2849 Такое число  $q$  называется *обратным* к  $b$  и обозначается  $b^{-1}$ . Теперь можно опре-  
 2850 делить деление на  $b$  как умножение на  $b^{-1}$ :

$$2851 \quad a/b = ab^{-1} \pmod{m}.$$

2852 Нетрудно проверить, что обратные могут быть только у чисел взаимно простых  
 2853 с  $m$ . Пусть  $\gcd(b, m) = d > 1$ . Тогда  $bq$  делится на  $d$  для любого  $q$ , а следовательно  
 2854 остаток от  $bq$  при делении на  $m$  тоже делится на  $d$ , т.е.  $(bq \pmod{m}) \neq 1$ .

2855 Пусть  $b$  и  $m$  взаимно просты. Как найти обратное к  $b$ ? Тут нам поможет расши-  
 2856 ренный алгоритм Евклида. Давайте найдём коэффициенты Безу для  $b$  и  $m$ . Это  
 2857 будет пара  $x$  и  $y$  удовлетворяющие следующему соотношению  $bx + my = 1$ . Рас-  
 2858 смотрим это соотношение по модулю  $m$ . Тогда в левой части второе слагаемое  
 2859 обнуляется, и мы получаем

$$2860 \quad bx = 1 \pmod{m},$$

2861 т.е.  $x$  — обратное к  $b$ . Осталось заметить, что расширенный алгоритм Евклида  
 2862 гарантирует, что  $|x| \leq m$ . Это даёт следующий алгоритм деления.

```

2863 # Деление на a на b по модулю m
2864 # предполагается, что gcd(b, m) = 1
2865 def divide_mod(a, b, m):
2866     # вычисляем коэффициенты Безу
2867     (d, x, y) = egcd(b, m)
2868
2869     # обработка отрицательного x
2870     if x < 0:
2871         x = x + m
2872
2873     return (a * x) % m

```

2874 Сложность полученного алгоритма  $O(n^2)$ .

2875 Давайте также проанализируем следующий код, вычисляющий  $a^b \pmod{m}$ .

```

2876 # Быстрое возведение в степень по модулю
2877 def power_mod(a, b, m):
2878     if b == 0:
2879         return 1
2880
2881     # рекурсивно вычисляем a^{b/2} mod m
2882     res = power_mod(a, b // 2, m)
2883

```

```
2884     # возводим результат в квадрат
2885     res = res * res
2886
2887     # обработка нечётного b
2888     if b % 2 == 1:
2889         res = res * a
2890
2891     return res % m
```

2892 Поскольку все числа в этом алгоритме не превосходят  $m^3$ , то алгоритм имеет слож-  
2893 ность  $O(n^3)$  (теперь без подвоха).



## 2894 Глава 9

# 2895 Альтернативные модели вычисления

2896 До этого момента мы всегда работали с моделью памяти с произвольным до-  
2897 ступом к памяти. Такая модель позволяет формализовать алгоритмы, которые  
2898 выполняются на одном процессоре, а все данные хранятся в оперативной па-  
2899 мяти. Однако, современные вычислительные системы могут быть значительно  
2900 сложнее. Например, может оказаться, что данных так много, что не влезают в  
2901 оперативную память одного компьютера. Такие данные приходится хранить на  
2902 диске или даже на нескольких дисках, доступ к которым значительно более мед-  
2903 ленный, и это нужно учитывать при разработке алгоритмов. Кроме того, вычис-  
2904 ления могут производиться не на одном процессоре, а сразу на нескольких маши-  
2905 нах, которые взаимодействуют между собой. В этой главе мы рассмотрим несколь-  
2906 ко моделей, которые позволяют описать алгоритмы для работы с большими дан-  
2907 ными и распределённые вычисления. В заключение, поговорим о других моде-  
2908 лях, которые не вошли в этот обзор.

## 2909 9.1. Работа с большими данными

2910 В этом разделе мы поговорим про модель внешней памяти, которая позволя-  
2911 ет формализовать взаимодействие в внешними запоминающими устройствами,  
2912 и подходит для описания алгоритмов над данными, которые не влезают в опера-  
2913 тивную память. Кроме этого мы обсудим важное уточнение этой модели — мо-  
2914 дель *cache-oblivious*.

### 2915 9.1.1. Модель внешней памяти

2916 *Модель внешней памяти (external memory model)* описывает вычислительное устрой-  
2917 ство, состоящее из процессора с внутренней памятью размера  $M$  (соответствует  
2918 оперативной памяти), также называемой *кешем*, и подключенной к нему неогра-  
2919 ниченной *внешней памяти* (соответствует внешнему запоминающему устройству  
2920 или взаимодействию по сети). Внутренняя и внешняя память разбиты на блоки  
2921 размера  $B$ . Все операции с внешней памятью оперируют блоками: за одну опера-  
2922 цию можно прочитать или записать один блок размера  $B$ . Сложность алгоритма в  
2923 этой модели определяется, как количество таких операций со внешней памятью.

2924 Может показаться странным, что сложность алгоритма никак не учитывает  
2925 количество операций процессора. Дело в том, что в данной модели предпола-  
2926 гается, что операции с данными во внутренней памяти на несколько порядков  
2927 быстрее, чем чтение и запись данных из внешней памяти. Поэтому вклад про-  
2928 цессорного времени в сложность алгоритма незначителен по сравнению вкладом  
2929 операций ввода-вывода.

2930 Какой сложности алгоритмов мы хотели бы добиться в этой модели? В идеале,  
2931 мы хотели бы «перенести» сложность алгоритмов из модели произвольной памя-  
2932 ти, заменив оценках  $n$  на  $n/B$ . В этом случае мы будем говорить, что полученный  
2933 алгоритм является эффективным. Например, если в RAM-модели некоторая за-  
2934 дача решается за  $O(n)$ , то в модели внешней памяти мы хотим построить алго-  
2935 ритм, работающий за  $O(n/B)$ . Далее, на примере нескольких задач мы увидим,  
2936 когда это удаётся сделать, а когда — нет.

### 2937 Поиск элемента в массиве

2938 В RAM-модели задача поиска элемента в массиве решается за  $O(n)$ . Напри-  
2939 мер, можно последовательно перебрать все элементы массива и сравнить их с  
2940 искомым. В модели внешней памяти эта задача решается за  $O(n/B)$  совершен-  
2941 но аналогичным образом: массив последовательно считывается во внутреннюю  
2942 память по блокам, и все элементы текущего блока сравниваются с искомым эле-  
2943 ментом. В худшем случае такой алгоритм считает не более  $\lceil n/B \rceil + 1 = O(n/B)$   
2944 блоков.

### 2945 Двоичный поиск в упорядоченном массиве

2946 В RAM-модели задача поиска элемента в упорядоченном массиве решается  
2947 за  $O(\log n)$ . Попробуем реализовать аналогичный алгоритм в модели внешней  
2948 памяти. Несложно заметить, что такой алгоритм всё так же будет работать за  
2949  $\Theta(\log n)$ . Это происходит из-за того, что при двоичном поиске нам нужно прове-  
2950 рять элементы массива, которые далеко находятся друг от друга, следовательно  
2951 на каждый элемент потребуется отдельная операция чтения. Так как чтение про-  
2952 исходит по блокам, то на каждый интересующий нас элемент массива мы неволь-  
2953 но дополнительно считываем  $B - 1$  элемент, которые никак не используются. Это  
2954 не эффективно, но, к сожалению, при двоичном поиске это неизбежно.

2955 Тем не менее, есть способ организовать данные во внешней памяти так, чтоб  
2956 поиск можно было реализовать эффективно. Для этого нужно хранить данные в  
2957 сбалансированном  $B$ -дереве поиска, размер вершины которого, совпадает с раз-  
2958 мером блока. В таком  $B$ -дереве можно искать за  $O(\log_B n)$ .

### 2959 Сортировка слиянием

2960 Сортировка слиянием является хорошим кандидатом на сортировку, которая  
2961 будет хорошо себя вести в модели со внешней памятью. Действительно, в этом  
2962 алгоритме данные в массивах читаются последовательно, что позволяет легко  
2963 модифицировать его для чтения по блокам. Для того, чтобы не задумываться о  
2964 том, достаточно ли внутренней памяти для хранения дерева рекурсии, мы бу-  
2965 дем использовать нерекурсивную версию этого алгоритма. Вся работа с очередь  
2966 массивов, которая нужна для нерекурсивной реализации этого алгоритма, эф-  
2967 фективно реализуется во внешней памяти, т.к. по своей сути является работой с  
2968 последовательными данными.

2969 Осталось разобраться с реализацией слияния. Для слияния двух массивов бу-  
2970 дем поддерживать во внутренней памяти три блока: текущий блок с данными  
2971 первого массива, текущий блок с данными второго массива и блок с результа-  
2972 тами слияния. Когда элементы какого-то из первых двух блоков полностью обрабо-  
2973 таны, этот блок заменяется следующим блоком соответствующего массива. Когда

2974 блок с результатами заполняется, то его содержимое записывается во внешнюю  
2975 память. В результате мы получаем алгоритм, который имеет сложность  $O(\frac{n}{B} \log \frac{n}{B})$ .

2976 Полученный алгоритм можно сделать ещё более эффективным. Недостаток  
2977 предложенной реализации заключается в том, что внутренняя память использу-  
2978 ется неэффективно — по сути используются только три блока. Это можно испра-  
2979 вить, если реализовать сортировку слиянием, которая разбивает массив не на две  
2980 части, а на  $k$  частей, и, соответственно, сливает массивы сразу  $k$  упорядоченных  
2981 массивов ( $k$ -ways mergesort). В RAM-машине такой алгоритм будет иметь слож-  
2982 ность  $\Theta(nk \log_k n)$  — множитель  $k$  в оценке возникает из-за того, что при слиянии  
2983  $k$  массивов на каждой итерации требуется выбрать минимум из  $k$  элементов. Од-  
2984 нако в модели внешней памяти операции процессора не учитываются. Поэтому  
2985 реализация такого алгоритма в этой модели будет иметь сложность  $O(\frac{n}{B} \log_k \frac{n}{B})$   
2986 при условии, что  $k + 1$  блоков влезет во внутреннюю память (дополнительный  
2987 блок нужен для хранения результата слияния). Поэтому возьмём  $k = M/B - 1$ .  
2988 Итоговая сложность в модели внешней памяти  $O(\frac{n}{B} \log_{M/B} \frac{n}{B})$ .

### 2989 9.1.2. Модель cache-oblivious

2990 Модель *cache-oblivious* (можно примерно перевести как модель для алгорит-  
2991 мов не обращающие внимание на кеш) — это модификация модели внешней па-  
2992 мяти, в которой алгоритму неизвестны значения  $M$  и  $B$  (и у алгоритма нет воз-  
2993 можности их узнать). Другими словами, в этой модели такое же вычислитель-  
2994 ное устройство, как и в модели внешней памяти, но всё кеширование и чтение  
2995 по блокам реализовано на более низком уровне, к которому у алгоритма нет до-  
2996 ступа. Т.е. алгоритму известно про кеширование, но не известно какие именно  
2997 параметры кеширования используются. Такое допущение сильно всё усложняет.  
2998 Например, в таких условиях у нас уже не получится применить алгоритм сорти-  
2999 ровки, рассмотренный выше — в нём мы опирались на знание  $M$  и  $B$ . С другой  
3000 стороны, обычный алгоритм линейного поиска элемента в массиве будет рабо-  
3001 тать вне зависимости от значений  $M$  и  $B$ : при последовательном чтении элемен-  
3002 тов из внешней памяти они будут загружаться в память блоками (неизвестно-  
3003 го алгоритму размера) и таким образом каждый элемент будет прочитан толь-  
3004 ко один раз. Следовательно, такой простой алгоритм будет иметь оптимальную  
3005 сложность  $O(n/B)$ . Про такие алгоритмы принято говорить, что они являются  
3006 *cache-friendly*.

3007 Удивительным образом, даже в отсутствии информации о значении  $M$  и  $B$ ,  
3008 удаётся построить оптимальные алгоритмы для множества задач. В том числе  
3009 можно построить алгоритм сортировки за время  $O(\frac{n}{B} \cdot \log_{M/B} \frac{n}{B})$ . Однако алгорит-  
3010 мы получаются значительно сложнее, чем сортировка в модели внешней памя-  
3011 ти (например, там могут потребоваться специальные довольно специфические  
3012 структуры данных, которые позволяют эффективно сливать упорядоченные мас-  
3013 сивы, не зная про параметры кеширования). Изучение этих алгоритмов выходит  
3014 за пределы этого курса.

3015 Нельзя не упомянуть, почему изучение этой модели так важно. Дело в том,  
3016 что имеет место следующий факт (тут он изложен неформально, но его можно  
3017 сформулировать строго и доказать).

3018 **Утверждение 9.1.1.** *Если алгоритм эффективен в модели cache-oblivious, то он*  
3019 *эффективно использует кеширование в системах с иерархическим кешированием.*

3020 Это особенно интересно, т.к. современные вычислительные устройства широ-  
3021 ко используют иерархическое кеширование: у процессора есть несколько уров-

ней кеширования для работы в оперативной памяти, далее идут кеши внешних запоминающих устройств, которые так же могут быть иерархическими, взаимодействие по сети тоже кешируется, и т.д. Получается, что эффективные алгоритмы для этой модели очень полезны на практике. К сожалению, зачастую сложность таких алгоритмов нивелирует все их преимущества.

## 9.2. Параллельные и распределённые вычисления

В этом разделе мы рассмотрим две модели для параллельных и распределённых вычислений.

### 9.2.1. Модель PRAM

Параллельная модель RAM (*parallel RAM, PRAM*) описывает вычислительное устройство с  $P$  процессорами и единой общей разделяемой между всеми процессорами памятью. При этом считается, что для работы с памятью уже реализованы необходимые механизмы синхронизации. В зависимости от типа механизмов синхронизации выделяют три разновидности PRAM:

- с эксклюзивным чтением и записью (EREW, exclusive read — exclusive write),
- с одновременным чтением и эксклюзивной записью (CREW, concurrent read — exclusive write),
- с одновременным чтением и записью (CRCW, concurrent read — concurrent write).

Сложность вычисления определяется максимальное среди всех процессоров количество совершённых операций.

Данная модель позволяет описывать алгоритмы для параллельных вычислений, но на практике возможность применения таких алгоритмов ограничена. Дело в том, что количество процессоров в пределах одного компьютера ограничено в т.ч. по инженерным соображениям (на практике  $P$  обычно исчисляется десятками или сотнями, в редких случаях несколькими тысячами). При необходимости увеличения числа процессоров на практике используются кластеры, состоящие из множества отдельных компьютеров, соединённых по сети, каждый из которых обладает своей собственной памятью. Таким образом, эта модель имеет практическое применение только для  $P = O(1)$ , что мало интересно с точки зрения сложности алгоритмов, т.к. такое допущение не позволяет ускорить вычисления более чем в константу раз.

Если же разрешить количеству процессоров  $P$  быть неконстантным, то в этой модели возникают алгоритмы с неожиданными оценками сложности, которые бесполезны для применения на практике. Для иллюстрации давайте посмотрим на алгоритм поиска максимума для модели CRCW PRAM, который находит максимум в массиве длины массиве из  $n$  за  $O(1)$  для  $P = n^2$ .

### Поиск максимума в модели CRCW PRAM

В CRCW PRAM нескольким процессорам разрешено одновременно обращаться к одной и той же ячейке памяти, как на чтение, так и на запись. При этом гарантируется, что при одновременном чтении все прочитанные данные будут



3063 корректны (соответствовать значениям, которые были записаны в ячейку), а при  
3064 одновременной записи разными процессорами в одну ячейку новое значение  
3065 ячейки будет соответствовать одному из записываемых значений. Алгоритм бу-  
3066 дет состоять из трёх фаз.

3067 1. На первой фазе  $n$  процессоров одновременно заполняют массив `maximum`  
3068 размера  $n$  значениями `True`. Таким образом, изначально каждый элемент  
3069 массива является кандидатом на то, чтобы быть максимумом. При этом каж-  
3070 дый процессор выполняет  $O(1)$  операций.

```
3071 # цикл выполняется параллельно на n процессорах  
3072 for i in range(n):  
3073     maximum[i] = True
```

3074 2. На второй фазе нам потребуются все  $n^2$  процессоров. Каждый процессор  
3075 сравнивает два элемента массива  $i$ , если оказывается, что первый элемент  
3076 меньше второго, то отмечает, что первый элемент не может быть миниму-  
3077 мом. При этом каждый процессор выполняет  $O(1)$  операций.

```
3078 # цикл выполняется параллельно на n^2 процессорах  
3079 for i in range(n):  
3080     for j in range(n):  
3081         if a[i] < a[j]:  
3082             maximum[i] = False
```

3083 3. К третьей фазе в массиве `maximum` значения `True` могут остаться только в  
3084 тех ячейках, которые соответствуют элементам массива, которые не мень-  
3085 ше всех оставшихся. Другими словами, `True` записано в только ячейках, ко-  
3086 торые соответствуют максимумам. Осталось найти такую ячейку и записать  
3087 её номер в переменную `result`. На этой фазе нам потребуется  $n$  процессо-  
3088 ров, каждый из которых опять выполнит только  $O(1)$  операций.

```
3089 # n процессоров вычисляют ответ  
3090 for i in range(n):  
3091     if maximum[i]:  
3092         result = i
```

3093 Полученный алгоритм действительно имеет сложность  $O(1)$ , но совершенно бес-  
3094 полезен на практике, ведь для того, чтобы найти максимум среди тысячи эле-  
3095 ментов, требуется компьютер с миллионом процессоров.

3096 *Упражнение 9.2.1.* Укажите, где в этом алгоритме используются свойства CRCW  
3097 PRAM? Будет ли алгоритм работать в других вариантах PRAM?

### 3098 9.2.2. Модель BSP

3099 Модель *частично-синхронного параллелизма (Bulk-Synchronous Parallel, BSP)* со-  
3100 стоит из

- 3101 •  $P$  процессоров, каждый со своей собственной локальной памятью (1 едини-  
3102 ца времени на операцию),

- 3103 • *коммуникационной среды*: сети для обмена сообщениями и общей внешней  
3104 памяти ( $G$  единиц времени на пересылку одного значения)<sup>1</sup>,
- 3105 • *барьерный механизм синхронизации* ( $L$  единиц времени на синхронизацию).

3106 Вычисления в такой модели устроены следующим образом. Вход алгоритма за-  
3107 писывается во внешнюю память. Процесс вычисления разбит на фазы — *супер-*  
3108 *шаги*. Внутри супершага процессоры взаимодействуют с внешней памятью, про-  
3109 изводят вычисления, используя данные во внутренней памяти, и посылают сооб-  
3110 щения другим процессорам. Все эти действия происходят асинхронно, т.е. один  
3111 процессор ничего не знает про состояние другого процессора. Внутри суперша-  
3112 га процессоры не могут как-то взаимодействовать между собой (один процес-  
3113 сор может послать сообщение другому процессору, но доставка может случить-  
3114 ся только в самом конце супершага). Супершаг заканчивается, когда происходит  
3115 барьерная синхронизация. Гарантируется, что после завершения синхронизации  
3116 все процессоры закончат выполнение своих текущих заданий, а все посланные  
3117 сообщения будут доставлены. После завершения синхронизации начинается сле-  
3118 дующий супершаг.

3119 Такая общая модель позволяет описать почти любую распределённую вычис-  
3120 лительную систему. Тут стоит отметить, что внешняя память и механизм син-  
3121 хронизации могут быть реализованы на основе сети для обмена сообщениями. В  
3122 то же время, сама сеть может быть реализована на основе доступа к общей внеш-  
3123 ней памяти. Таким образом, некоторые составляющие части модели BSP, кото-  
3124 рые могут отсутствовать в какой-то конкретной вычислительной системе, могут  
3125 быть реализованы за счёт других составляющих.

3126 Для каждой конкретной распределённой вычислительной системы можно по-  
3127 добрать значения  $P$ ,  $G$  и  $L$ , которые её описывают. Параметр  $G$  соответствует *от-*  
3128 *ставанию сети (обратной пропускной способности)*, т.е. показывает во сколько раз  
3129 время пересылки одной ячейки памяти превышает время доступа к ячейке памя-  
3130 ти во внутренней памяти процессора, которое принято за единицу. Параметр  $L$   
3131 соответствует *задержке* сети, т.е. ограничивает время, которое в худшем случае  
3132 нужно на доставку сообщения. Например, для суперкомпьютера Cray T3E значе-  
3133 ния параметров будут такими:  $P = 64$ ,  $G \approx 78$ ,  $L \approx 1825$ .

3134 Сложность вычисления естественным образом соответствует количеству по-  
3135 траченных единиц времени (по глобальному времени). Введём следующие обо-  
3136 значения. Для конкретного процессора обозначим

- 3137 •  $\text{comp}(s, p)$  — число операций совершенных процессором  $p$  на супершаге  $s$ ,
- 3138 •  $\text{comm}(s, p)$  — число ячеек переданных и полученных процессором  $p$  на су-  
3139 першаге  $s$ .

3140 Теперь для всего супершага  $s$  можно обозначить

- 3141 •  $\text{comp}(s) = \max_{p \in [P]} \text{comp}(s, p)$  — общее число операций на супершаге  $s$ ,
- 3142 •  $\text{comm}(s) = \max_{p \in [P]} \text{comm}(s, p)$  — общая коммуникация на супершаге  $s$ .

3143 Тогда сложность шага  $s$  можно определить так:

$$3144 \text{cost}(s) = \text{comp}(s) + G \cdot \text{comm}(s) + L.$$

<sup>1</sup>В рамках модели BSP предполагается, что посылка сообщений и взаимодействие с внешней памятью имеют общую природу и стоят одинаково.

3145 Если вычисление содержит sync синхронизаций (и, следовательно, sync супершагов), то общая сложность вычисления определяется так:

$$3147 \quad \text{cost} = \sum_{s \in [\text{sync}]} \text{cost}(s) = \text{comp} + G \cdot \text{comm} + L \cdot \text{sync},$$

3148 где comp и comm — это общее количество вычислений и общее количество коммуникаций, соответственно. Так как мы не знаем конкретных значений  $P$ ,  $G$  и  $L$ ,  
3149 то мы будем оценивать составляющие comm, comp и sync по-отдельности.

3151 При построении алгоритмов в этой модели мы будем предполагать, что размер входных данных всегда значительно превышает количество процессоров, т.е.  $n \gg P$ . Какие алгоритмы в этой модели мы могли бы назвать эффективными? В идеале мы ожидаем, что ускорение от использования распределённых вычислений будет пропорционально количеству процессоров. Поэтому для задач, которые решаются на RAM машине за  $O(f(n))$  с использованием  $O(g(n))$  операций ввода-вывода, мы будем стремиться получить алгоритм, для которого

$$3158 \quad \text{comp} = O(f(n)/P), \quad \text{comm} = O(g(n)/P).$$

3159 Более того, мы хотели бы минимизировать количество синхронизаций. Мы будем стремиться разрабатывать алгоритмы, в которых количество синхронизаций не зависит от размера входа, т.е. sync может быть функцией  $P$ , но не  $n$ , а ещё лучше — константой. Давайте рассмотрим пару примеров таких алгоритмов.

### 3163 Поиск элемента в массиве

3164 Для определённости будем искать первое вхождение элемента в массив в предположении  $n = \Omega(P^2)$ . Алгоритм будет состоять из двух супершагов. На первом супершаге каждый процессор считывает из внешней памяти и проверит свой подмассив размера  $n/P$ . После этого каждый процессор, который нашёл искомым элемент в своём подмассиве, отправит индекс первого вхождения некоторому выделенному процессору, например, процессору с номером 1. Таким образом, для первого супершага мы получаем

$$3171 \quad \text{comp}(1) = O(n/P), \quad \text{comm}(1) = O(n/P).$$

3172 На втором супершаге этот выделенный процессор выберет из присланных индексов наименьший и запишет его в переменную для ответа. Соответственно,

$$3174 \quad \text{comp}(2) = O(P), \quad \text{comm}(2) = O(1).$$

3175 Таким образом, при  $n = \Omega(P^2)$  данный алгоритм имеет следующие оценки:

$$3176 \quad \text{comp} = O(n/P), \quad \text{comm} = O(n/P), \quad \text{sync} = O(1).$$

### 3177 Сортировка методом регулярного сэмплирования

3178 Будем предполагать, что  $n = \Omega(P^3)$ , и что все элементы массива различны.  
3179 Алгоритм будет состоять из четырёх супершагов.

3180 1. Каждый процессор считывает из внешней памяти и сортирует свой подмассив размера  $n/P$ . Из полученного упорядоченного массива выбирается набор из  $P$  порядковых статистик с шагом  $n/P^2$  начиная с индекса 0, будем

3183 называть их *пробами*. Таким образом массив разбивается на  $P$  блоков оди-  
 3184 накового размера, причём каждый блок начинается с элемента-пробы. Все  
 3185 процессоры пересылают свои наборы проб некоторому выделенному про-  
 3186 цессору.

$$3187 \quad \text{comp}(1) = O((n \log n)/P), \quad \text{comm}(1) = O(n/P).$$

3188 2. Выделенный процессор получает  $P^2$  проб, сортирует их и выбирает из них  
 3189  $P$  порядковых статистик с шагом  $P$  начиная с индекса 0, будем называть их  
 3190 *разделителями* и обозначать  $s_1, \dots, s_P$ . Полученный набор разделителей пе-  
 3191 ресылается всем процессорам.

$$3192 \quad \text{comp}(2) = O(P^2 \cdot \log P), \quad \text{comm}(2) = O(P^2).$$

3193 3. Разделители задают разбиение элементов исходного массива на  $P$  (не обяза-  
 3194 тельно равных) частей. Теперь каждому процессору с номером  $i$  будут соот-  
 3195 ветствовать полуинтервал элементов от  $[s_i, s_{i+1})$  (последнему процессору со-  
 3196 ответствует полуинтервал  $[s_P, \infty)$ ). На третьем супершаге каждый процессор  
 3197 просматривает элементы своего подмассива (полученного на первом супер-  
 3198 шаге и уже к этому моменту упорядоченного) и отправляет каждый элемент  
 3199 тому процессору, которой отвечает за соответствующий полуинтервал. Для  
 3200 того, чтобы оценить число сообщений, которое в конце супершага получит  
 3201 каждый процессор, нам придётся доказать следующее утверждение.

3202 **Утверждение 9.2.1.** *Каждый полуинтервал содержит не более  $2n/P$  элемен-*  
 3203 *тов исходного массива.*

3204 *Доказательство.* Для каждого полуинтервала  $T$  мы будем называть блок  $B$   
 3205 *левым*, если его левая граница не больше левой границы  $T$ . Все остальные  
 3206 блоки мы будем называть *правыми*. Верны следующие наблюдения.

- 3207 • У каждого процессора есть не более одного левого блока, который бы  
 3208 пересекался с  $T$ . Действительно, если какой-то левый блок некоторого  
 3209 процессора пересекает  $T$ , то все блоки до него не пересекают  $T$ , а все  
 3210 блоки после него — правые.
- 3211 • Среди всех блоков всех процессоров есть не более  $P$  правых блоков, ко-  
 3212 торые пересекаются с  $T$ . Каждый правый блок пересекающийся с  $T$  на-  
 3213 чинается с некоторой пробы, причём эта проба обязательно принад-  
 3214 лежит  $T$ . Так как каждый полуинтервал содержит не более  $P$  проб, то  
 3215 таких блоков не больше  $P$ .

3216 Следовательно, каждый интервал содержит элементы из не более, чем  $2P$   
 3217 блоков, т.е. в сумме не более  $2P \cdot n/P^2 = 2n/P$  элементов.  $\square$

3218 Это позволяет оценить сложность третьего супершага:

$$3219 \quad \text{comp}(3) = O(n/P), \quad \text{comm}(3) = O(n/P).$$

3220 Оценка на коммуникацию  $O(n/P)$  учитывает и  $n/P$  элементов, которые по-  
 3221 слал каждый процессор, и не более  $2n/P$  элементов, которые каждый про-  
 3222 цессор получил.

3223 4. Каждый процессор сортирует элементы своего интервала и записывает их  
3224 во внешнюю память.

$$3225 \text{comp}(4) = O((n \log n)/P), \quad \text{comm}(4) = O(n/P).$$

3226 В сумме при условии  $n = \Omega(p^3)$  получаем желаемые оценки:

$$3227 \text{comp} = O((n \log n)/P), \quad \text{comm} = O(n/P), \quad \text{sync} = O(1).$$

3228 Можно показать, что для данной эти оценки являются оптимальными для сорти-  
3229 ровок, основанных на сравнениях.

### 3230 9.3. Другие модели

3231 Есть множество других моделей, которые оказались за рамками этого экскур-  
3232 са. Например, есть гибридная модель, объединяющая модель внешней памяти  
3233 и PRAM — *параллельная модель внешней памяти (parallel external memory)*. Мож-  
3234 но описать модель вычисления для формализации алгоритмов на GPU — модель  
3235 описывающая принцип «*одиночный поток команд, множественный поток данных*»,  
3236 *ОКМД (SIMD, single instruction multiple data)*. Есть и другие модели, непохожие на  
3237 то, что мы уже видели. К ним, например, относится модель вычислений для *по-*  
3238 *токовых алгоритмов*, в которой вычисляющее устройство с ограниченной памя-  
3239 тью может читать входные данные исключительно последовательно и только один  
3240 раз. Представьте себе, например, сетевой маршрутизатор, обрабатывающий се-  
3241 тевой трафик в реальном времени, пытающийся выделить наиболее популярные  
3242 IP адреса. Память маршрутизатора несравнимо меньше, чем объём трафика, ко-  
3243 торый через него проходит. Однако, некоторые характеристики можно вычис-  
3244 лить и в такой ограниченной модели, чему посвящено большое число статей и  
3245 даже книг.

3246 Кроме того, стоит упомянуть ещё модели, которые позволяют формализовать  
3247 понятие алгоритма для квантовых компьютеров. Это отдельная интересная тема,  
3248 которая на данный момент всё ещё является исключительно теоретической (су-  
3249 ществующие квантовые компьютеры имеют очень ограниченные возможности).



## 3250 Приложение А

## 3251 Математические факты

### 3252 А.1. Суммирование последовательностей

3253 **Лемма А.1.1.** Сумма первых  $n$  членов арифметической прогрессии

$$3254 \quad a_1 = a, \quad a_i = a_{i-1} + d = a + (i-1)d,$$

3255 вычисляется по формуле

$$3256 \quad S_n = \frac{a_1 + a_n}{2} \cdot n = \frac{2a + (n-1)d}{2} \cdot n.$$

*Доказательство.*

$$3257 \quad 2S_n = \sum_{i=1}^n a_i + \sum_{i=1}^n a_i = \sum_{i=1}^n a_i + \sum_{i=1}^n a_{n-i+1} = \sum_{i=1}^n (a_i + a_{n-i+1}) = \sum_{i=1}^n (2a + (n-1)d).$$

3258 Отсюда получаем, что

$$3259 \quad S_n = \frac{1}{2} \sum_{i=1}^n (2a + (n-1)d) = \frac{n}{2} \cdot (2a + (n-1)d).$$

3260

□

3261 **Лемма А.1.2.** Сумма первых  $n$  членов геометрической прогрессии

$$3262 \quad a_1 = a, \quad a_i = a_{i-1} \cdot q = a \cdot q^{i-1},$$

3263 вычисляется по формуле

$$3264 \quad S_n = \frac{a_{n+1} - a_1}{q - 1} = a \cdot \frac{q^n - 1}{q - 1}.$$

*Доказательство.*

$$S_n = \sum_{i=1}^n a_i = \sum_{i=1}^n a \cdot q^{i-1}.$$

$$3265 \quad qS_n = q \sum_{i=1}^n a \cdot q^{i-1} = \sum_{i=2}^{n+1} a \cdot q^{i-1}.$$

$$qS_n - S_n = \sum_{i=2}^{n+1} a \cdot q^{i-1} - \sum_{i=1}^n a \cdot q^{i-1} = a \cdot q^n - a.$$

3266 Отсюда получаем, что

$$3267 \quad S_n = \frac{a \cdot q^n - a}{q - 1} = a \cdot \frac{q^n - 1}{q - 1}.$$

3268

□

3269 **Лемма А.1.3.** Сумма бесконечной убывающей геометрической прогрессии с  $q \in (0, 1)$   
 3270 вычисляется по формуле  $S = \frac{1}{1-q}$ .

*Доказательство.*

$$3271 \quad S = \lim_{n \rightarrow \infty} a \cdot \frac{q^n - 1}{q - 1} = a \cdot \frac{-1}{q - 1} = \frac{a}{1 - q}.$$

3272

□

3273 **Лемма А.1.4.** Сумма первых  $n$  членов бесконечной убывающей последовательности  
 3274  $a_i = i \cdot q^i$  при  $q \in (0, 1)$  вычисляется по формуле  $S_n = \frac{q}{(1-q)^2}$

3275 *Доказательство.* Рассмотрим функцию  $f(x) = \sum_{i=1}^{\infty} x^{i-1}$  и вычислим её вычис-  
 3276 лим производную.

$$3277 \quad f'(x) = \sum_{i=1}^{\infty} (i-1)x^{i-2} = \sum_{i=1}^{\infty} i \cdot x^{i-1}.$$

3278 С другой стороны по лемме А.1.3  $f(x) = \frac{1}{1-x}$  при  $x \in (0, 1)$ .

$$3279 \quad f'(x) = \left( \frac{1}{1-x} \right)' = \frac{1}{(1-x)^2}.$$

3280 Осталось заметить, что

$$3281 \quad S_n = \sum_{i=1}^{\infty} i \cdot q^i = q \cdot f'(q),$$

3282 а следовательно  $S_n = \frac{q}{(1-q)^2}$ .

□



## Литература

- 3284 [1] Бабенко М.А., Левин М.В. Введение в теорию алгоритмов и структур дан-  
3285 ных. — Москва : МЦНМО, 2016. — ISBN: [9785443923963](#).
- 3286 [2] Алгоритмы. Построение и анализ [пер. с англ.] / Кормен Т., Лейзерсон Ч., Ри-  
3287 вест Р., and Штайн К. — Вильямс, 2009. — ISBN: [9785845908575](#). — Access mode:  
3288 <https://books.google.ru/books?id=UVg1LPacgRcC>.
- 3289 [3] Шень А. Программирование: теоремы и задачи. — Москва : МЦНМО,  
3290 2017. — ISBN: [978-5-4439-0685-0](#). — Access mode: [https://hal.archives-ouvertes.](https://hal.archives-ouvertes.fr/hal-01480636/document)  
3291 [fr/hal-01480636/document](https://hal-01480636/document).
- 3292 [4] Sleator Daniel Dominic, Tarjan Robert Endre. Self-adjusting Binary Search Trees //  
3293 *J. ACM*. — 1985. — July. — Vol. 32, no. 3. — P. 652–686. — Access mode: [http://doi.](http://doi.acm.org/10.1145/3828.3835)  
3294 [acm.org/10.1145/3828.3835](http://doi.acm.org/10.1145/3828.3835).
- 3295 [5] Vuillemin Jean. A Unifying Look at Data Structures // *Commun. ACM*. — 1980. —  
3296 Apr. — Vol. 23, no. 4. — P. 229–239. — Access mode: [http://doi.acm.org/10.1145/](http://doi.acm.org/10.1145/358841.358852)  
3297 [358841.358852](http://doi.acm.org/10.1145/358841.358852).
- 3298 [6] Seidel Raimund, Aragon Cecilia R. Randomized search trees // *Algorithmica*. —  
3299 1996. — Vol. 16. — P. 464–497.