

# Глава 1

## Использование Стека Вызовов

### 1.1. Проверка скобочных выражений

При решении задач с использованием стека, или при помощи автоматов со стековой памятью, часто можно обойтись без организации стека как самостоятельной структуры данных. В качестве примера рассмотрим две реализации программы проверки правильности скобочного выражения (с несколькими типами скобок).

Для начала рассмотрим *синтаксис* скобочных выражений. Правильное скобочное выражение  $S$  подчиняется следующим правилам:

$$\begin{array}{l} S \leftarrow \emptyset \\ | \quad ' ( ' S ' ) ' S \\ | \quad ' \{ ' S ' \} ' S \\ | \quad ' [ ' S ' ] ' S \end{array}$$

Этот набор правил называется *грамматикой* и его следует читать так:

Правильное скобочное выражение  $S$  является либо пустой строкой, либо правильным скобочным выражением, заключенным в парные скобки (различных типов), за которым следует правильное скобочное выражение.

$S$  в данном наборе называется нетерминалом, а все варианты скобок – алфавитом или терминалами – символами, которые не допускают дальнейшего определения.

Обычно, для проверки алгоритму отдается строка, и на выходе алгоритм выдает бинарный ответ — ДА или НЕТ (является ли строка правильным скобочным выражением или нет). Классическим решением данной задачи является организация автомата со стековой памятью, который последовательно считывает символы строки и сохраняет в стековой памяти путь из открытых скобок. Это решение следует из того, что указанный набор правил можно прочесть так:

Правильное скобочное выражение — это такая последовательность символов, в которой каждая открывающая скобка должна быть закрыта либо сразу, либо после правильного скобочного выражения, а каждая закрывающая скобка должна следовать либо сразу после парной открывающей, либо после правильного скобочного выражения.

### 1.2. Использование внешнего стека

Рассмотрим классическое решение. Сначала опишем стековые процедуры `push` и `pop`:

- 1  $\langle \text{Реализация стека 1} \rangle \equiv$   
     $\langle \text{Глобальные переменные стека 2a} \rangle$   
     $\langle \text{Реализация push 2b} \rangle$   
     $\langle \text{Реализация pop 2c} \rangle$

Для простоты положим стек как массив символов ограниченной, но достаточной длины, подходящей под лимиты задачи (если они есть, конечно):

2a  $\langle$ Глобальные переменные стека 2a $\rangle \equiv$  (1)

```
#define MAX_STACK 1024
char stack[MAX_STACK];
int stack_ptr = 0;
```

Defines:

MAX\_STACK, used in chunk 2b.

stack, used in chunk 2.

stack\_ptr, used in chunk 2.

Функция **push** кладет значение на вершину стека и увеличивает указатель **stack\_ptr** на единицу.

На всякий случай, проверяем, что мы не переполнили стек:

2b  $\langle$ Реализация push 2b $\rangle \equiv$  (1)

```
int push(char x) {
    if ( stack_ptr >= MAX_STACK)
        return -1;
    stack[stack_ptr++] = x;
    return 0;
}
```

Defines:

push, used in chunk 3b.

Uses MAX\_STACK 2a, stack 2a, and stack\_ptr 2a.

Функция **pop** проверяет, что стек не пуст, возвращает значение с вершины стека (через указатель) и уменьшает **stack\_ptr**.

2c  $\langle$ Реализация pop 2c $\rangle \equiv$  (1)

```
int pop(char* x) {
    if ( stack_ptr <= 0 )
        return -1;
    *x = stack[--stack_ptr];
    return 0;
}
```

Defines:

pop, used in chunk 3.

Uses stack 2a and stack\_ptr 2a.

Теперь, допустим, у нас уже есть строка *s* для проверки. Просто перемещаемся последовательно по строке и, в зависимости от читаемого символа, если он является открывающей скобкой – кладем (парную ей) на стек, если закрывающей — сверяем с вершиной. При любом несовпадении, либо по досрочном окончании строки или стека — делаем заключение о неправильности скобочного выражения.

3a  $\langle \text{Функция проверки выражения со стеком 3a} \rangle \equiv$

```
int check(char* s) {
    int err = 0;
    char x;
    while(*s) {
        switch (*s) {
             $\langle \text{Открывающие скобки - кладем на стек 3b} \rangle$ 
             $\langle \text{Закрывающие скобки сверяем с вершиной 3c} \rangle$ 
            default:
                break;
        }
        if ( err ) return -1;
        s++;
    }
     $\langle \text{Проверяем, что стек пуст 3d} \rangle$ 
    return 0;
}
```

Defines:

check, never used.

В случае, когда читаемый символ *\*s* является открывающей скобкой, кладем на стек (с проверкой выхода за пределы) закрывающую пару.

3b  $\langle \text{Открывающие скобки - кладем на стек 3b} \rangle \equiv$  (3a)

```
case '(':
    err = push(')'); break;
case '{':
    err = push('}'); break;
case '[':
    err = push(']'); break;
```

Uses push 2b.

Если читаемый символ — закрывающая скобка, сверяем с вершиной (если она доступна). Выходим сразу, если что-то не так.

3c  $\langle \text{Закрывающие скобки сверяем с вершиной 3c} \rangle \equiv$  (3a)

```
case ')': case '}': case ']':
    if ( (err = pop(&x)) || x != *s )
        return -1;
    break;
```

Uses pop 2c.

Все остальные символы мы просто игнорируем, перемещаясь по строке дальше. Как только мы дошли до конца строки:

3d  $\langle \text{Проверяем, что стек пуст 3d} \rangle \equiv$  (3a)

```
if ( !pop(&x) ) return -1;
```

Uses pop 2c.

Данный код, а точнее код функций `push` и `pop` будет естественным образом усложняться, если размер стека заранее неизвестен и ограничений тоже никаких нет.

### 1.3. Использование стека вызовов

Посмотрим, как можно обойтись без использования внешней структуры и методов `push` и `pop`. Оказывается, древовидная структура рекурсивных вызовов функций языка Си соответствует древовидному синтаксису наших выражений, а переменные, хранящиеся в стеке вызовов соответствуют стеку открытых скобок. Если в строке встречается открывающая скобка — запоминаем парную ей в локальной переменной нашей функции. Если закрывающая — возвращаемся, а вызывающая сторона проверит совпадение.

4a  $\langle \text{Рекурсивная проверка 4a} \rangle \equiv$  (5a)

```

char check_rec(char** s) {
    char wait_for;
    char c = *(*s)++;
    switch(c) {
         $\langle \text{Открывающие скобки - запоминаем 4b} \rangle$ 
         $\langle \text{Закрывающие скобки - возвращаемся 4c} \rangle$ 
        default:
            return check_rec(s);
    }
     $\langle \text{Проверяем совпадение 4d} \rangle$ 
}

```

Defines:

`check_rec`, used in chunks 4d and 5a.

В этой реализации мы используем указатель на указатель `char** s` для передачи между рекурсивными вызовами “потока” символов. Каждый вызов может модифицировать указатель на символ потока. Чаще всего нам потребуется читать следующий символ при рекурсивном вызове, поэтому мы сразу перемещаем указатель после чтения в переменную `c`.

4b  $\langle \text{Открывающие скобки - запоминаем 4b} \rangle \equiv$  (4a)

```

case '{': wait_for = '}'; break;
case '[': wait_for = ']'; break;
case '(': wait_for = ')'; break;

```

Закрытием считаем не только скобки, но и конец строки (символ `'\0'`). Поскольку функция проверки должна возвращать 0 в случае правильности выражения и прочие значения — в противном случае, этот вариант вполне нам подойдет.

4c  $\langle \text{Закрывающие скобки - возвращаемся 4c} \rangle \equiv$  (4a)

```

case '}': case ']': case ')': case '\0': return c;

```

После успешной проверки — двигаемся дальше. Общим результатом проверки будет являться результат проверки выражения, начинающегося со следующего символа (поток `s` уже подготовлен). Если `check_rec` вернул не то, что мы ожидаем — возвращаемся с ненулевым результатом, который будет интерпретирован, как ошибка.

4d  $\langle \text{Проверяем совпадение 4d} \rangle \equiv$  (4a)

```

return check_rec(s) != wait_for ? -1 : check_rec(s);

```

Uses `check_rec` 4a 5b.

Как видим, второе решение более изящно и ближе к оригинальному описанию проверяемого синтаксиса.

Сведем все это воедино, с простой оберткой для тестирования:

5a  $\langle skobkicheck.c\ 5a \rangle \equiv$

```
#include <stdio.h>
⟨Рекурсивная проверка 4a⟩
int main() {
    char buf[8192];
    char* s = fgets(buf, sizeof(buf), stdin);
    if ( !check_rec(&s) )
        printf("OK\n");
    else
        printf("NOT OK\n");
    return 0;
}
```

Defines:

`main`, never used.

Uses `check_rec` 4a 5b.

Возможно, более наглядным (но не самым лучшим) и приближенным к математическому определению будет вариант, когда мы не будем выносить проверку совпадения за оператор `switch`, а будем проверять и двигаться дальше непосредственно в `case-e`:

5b  $\langle \text{Еще одна реализация check\_rec 5b} \rangle \equiv$

```
char check_rec(char** s) {
    char c = *(*s)++;
    switch (c) {
        case '{': return check_rec(s) != '}' ? -1 : check_rec(s);
        case '[': return check_rec(s) != ']' ? -1 : check_rec(s);
        case '(': return check_rec(s) != ')' ? -1 : check_rec(s);
        case '}': case ']': case ')': case '\0': return c;
        default: return check_rec(s);
    }
}
```

Defines:

`check_rec`, used in chunks 4d and 5a.

*Красота!*

## Задание

1. Дополните решение возможность анализировать выражение из четырех типов скобочных пар: {}, (), <>, [].
2. Дополните грамматику и анализатор так, чтобы в строке могли встечаться строки, заключенные в двойные кавычки, между которыми текст проверяться не должен. Пример правильного выражения: (")(")(")<>"("). Подсказка: набор правил теперь должен включать еще вспомогательный нетерминал  $T$ , обозначающий строку без кавычек внутри:

$$\begin{array}{lcl}
 T & \leftarrow & \emptyset \\
 & | & A T \\
 S & \leftarrow & \emptyset \\
 & | & ' ' T ' ' S \\
 & | & ' ( S ' ) ' S \\
 & | & ' \{ S \} ' S \\
 & | & ' [ S ] ' S
 \end{array}$$

где  $A$  – произвольный символ, не равный ". Для дополнительного символа нетерминала  $T$  нужно будет сделать отдельную анализирующую функцию.

## Глава 2

# Потоковая обработка данных

Достаточно часто встречается задача обработки данных, которые поступают на вход вычислительной системы последовательно. Причём возможности повторного их получения нет, либо она существенно ограничена. В этом случае говорят, что имеют дело с *потоком* данных. Примерами потока могут служить символы, введенные с клавиатуры, сигнал, принимаемый из эфира или микрофона, равно как и с любых других датчиков, сетевые пакеты и т.п. Обработка потока происходит т.н. *однопроходными алгоритмами*, поскольку они проходят всю последовательность входных данных ровно один раз и второй проход (или произвольный доступ) к данным им не требуется.

Однопроходный алгоритм также может быть применен, конечно, и к файлам с произвольным доступом, однако "однопроходность" подчеркивает относительную простоту модели вычислений.

Результатом работы однопроходного алгоритма может быть также *поток* данных и/или выдача результата вычислений по окончании входного потока.

Обычно выделяют три класса обработчиков: *фильтры* (filter), *отображения* (map) и *редукцию* (reduce), хотя все они являются вариациями *свёртки* (fold). Все потоковые вычисления можно свести к композиции таких обработчиков.

Фильтры принимают на вход поток и возвращают новый поток, но содержащий только элементы удовлетворяющие некоторому условию. Например, фильтр всех четных чисел:

$$f(\{1, 2, 3, 4, 5, 6, \dots\}) = \{2, 4, 6, \dots\}$$

Отображения возвращают новый поток, который содержит результат применения некоторой функции к каждому элементу входного потока. Например, возведение в квадрат:

$$m(\{1, 2, 3, 4, 5, 6, \dots\}) = \{1, 4, 9, 16, 25, 36, \dots\}$$

Редукция строит *агрегат* из элементов потока, если задана операция "сложения"  $a \oplus b$  на множестве элементов потока, обладающая свойством дистрибутивности. например, суммирует числа:

$$r(\{1, 2, 3, 4, 5, 6\}) = 31$$

Можно, конечно, заставить редукцию выдавать поток (она готова для каждого входного элемента предоставить результат), но обычно она "ждет" окончания входного потока и выдает агрегат.

Обобщенная свертка похожа на редукцию, но она может накапливать результат (строить агрегат) типа отличного от типа элементов потока. Вот такую, наиболее общую, форму алгоритмов мы и рассмотрим.

### 2.1. Однопроходные алгоритмы без памяти

Потоком  $S$  назовем однородную последовательность данных  $S = \{a_0, a_1, \dots, a_i, \dots\}$ , принадлежащих некоторому типу  $a_i \in A$ .

Подразумевается, что длина потока заранее неизвестна, хотя конечна.

Считается, что для потока заданы две функции, имеющие семантику *головы* и *хвоста* —  $h$  и  $t$ . Взятие элемента из потока  $h : S \rightarrow A$ , которая возвращает элемент  $A$ , если поток не пуст, и функция перехода к следующему состоянию потока  $t : S \rightarrow S$ . Для полноты описания лучше иметь эти

функции разделенными, хотя, часто оказывается, что  $h$  и  $t$  совмещены ( $t$  вызывается автоматически, а состояние потока  $S$  скрыто, как, например, при использовании `stdin` в `getchar` или `scanf`).

Поскольку длина потока неизвестна, и он может закончиться в любой момент, результат нужно готовить сразу, с момента поступления первого элемента. Материал для построения результата называется состоянием  $Q$ , которое может меняться после почтения каждого элемента. Сам результат  $R$ , который вообще-то может быть менее богат, чем состояние, должен получаться из него посредством функции  $g : Q \rightarrow R$ .

Рассмотрим типичный обработчик потока. Чаще всего он имеет форму  $p : S \times Q \rightarrow R$ . На входе мы имеем поток  $S$  (для которого заданы функции  $h$  и  $t$ ), и, важно заметить, какое-то состояние  $Q$ . Само тело обработчика выражается в виде функции  $f : Q \times A \rightarrow Q$ , которая, собственно, и занимается модификацией состояния в зависимости от обрабатываемого элемента.

$$p(s, q) = \begin{cases} s = \emptyset : & g(q) \\ \text{иначе} : & p(t(s), f(q, h(s))) \end{cases} \quad (2.1)$$

Т.е. в случае, если поток пуст ( $h$  не может больше взять элемент из потока), вычисляем результат исходя из накопленного состояния  $q$ . Иначе, если элемент в потоке существует, то производим вычисление и рекурсивно продолжаем обработку потока. Запуск обработчика подразумевает, помимо потока, также и начальное состояние. Назовем его  $q_0$ . В каких-то случаях оно может оказаться тривиальным, а где-то его придется еще специальным образом "готовить".

Также существует другой вариант обработчика:

$$\begin{aligned} p(s, q) &= g(p'(s, q)) \\ p'(s, q) &= \begin{cases} s = \emptyset : & q \\ \text{иначе} : & f(p'(t(s), q), h(s)) \end{cases} \end{aligned} \quad (2.2)$$

Здесь нам пришлось ввести вспомогательную функцию  $p'$ , которая, собственно, делает вычисления, а сама  $p$  является *оберткой*.

Дело в том, что находясь в середине потока, можно смотреть на него в двух направлениях. С одной стороны, можно считать себя первым и смотреть на оставшийся "хвост" потока (2.1). Такой обработчик называется *право-рекурсивным* или *хвостовым*. С другой стороны, можно смотреть в начало и считать, что на нас обработка потока должна закончиться, и основная работа уже проделана. В этом случае обработка (2.2) называется *лево-рекурсивной*.

Заметьте, что в обоих случаях последовательность чтения элементов из потока одинакова, и левая рекурсия не означает, что поток читается с конца.

Многие потоки не имеют семантики "хвоста". Из них можно только брать элементы, а хвост сдвигается сам. Т.е. для потока определена только функция  $n : S \rightarrow A$  (next, следующий), а факт окончания потока как-то должен отражаться в возвращаемом значении. Например, можно дополнить множество  $A$  пустым значением  $\emptyset$  и считать его концом потока. В этом случае обработчик выглядит чуть проще:

$$p(s, q) = \begin{cases} n(s) = \emptyset : & g(q) \\ \text{иначе} : & p(s, f(q, n(s))) \end{cases} \quad (2.3)$$

однако, теряется т.н. *ссылочная прозрачность*, когда вызов функции  $n(s)$  можно заменить ее возвращаемым значением, поскольку она меняет *внутреннее состояние* потока.

Алгоритм называется *однопроходным без памяти*, если он имеет хвостовую форму и размер состояния  $\|Q\| = O(1)$  – является константой и не зависит количества прочитанных элементов. Хвостовая форма в данном определении важна, поскольку левая рекурсия подразумевает, что функция  $f$  "ждет" результата рекурсивного вызова  $p'$ , а для этого потребуется стековая память (вспомните функцию переворота потока символов).

Например, простой подсчет суммы потока действительных чисел можно разобрать на следующие "запчасти":

- В качестве функции чтения элемента из потока  $n$  задействуем функцию `scanf`. При работе со стандартным потоком ввода `stdin` сам поток таскать с собой не надо, функция `scanf` и так о нем знает, поэтому его можно из нашей конструкции исключить.



- Состояние  $Q$  – это действительное число, которое, собственно, и является суммой.
- Функция  $f(q, a) = q + a$  – простая сумма двух чисел.
- Функция  $g(q)$  – пусть печатает результат. Будем использовать для этого `printf`.

Давайте оформим это в виде программы на Си.

10a  $\langle \text{обработчик-суммирование } 10a \rangle \equiv$  (10b)

```
int sum(double q) {
    double a;
    if ( scanf("%lf", &a) == 1 )
        return sum(a+q);
    return printf("%lf\n", q);
}
```

Defines:

sum, used in chunks 10 and 11b.

Заметили, где здесь прячутся функции  $n$ ,  $p$ ,  $f$ ,  $g$ ? Поскольку `scanf` возвращает количество прочитанных шаблонов, можно использовать ее для проверки окончания потока (если `scanf` не может прочитать число двойной точности, то вернет 0, а если поток `stdin` закончится (EOF), то -1). Сами прочитанные значения она записывает по переданным адресам. В нашем случае, пришлось выделить переменную `double a` для того, чтобы хранить результат чтения из потока.

Целиком программа выглядит так:

10b  $\langle \text{sum.c } 10b \rangle \equiv$

```
#include <stdio.h>
<обработчик-суммирование 10a>

int main() {
    return sum(0);
}
```

Defines:

main, never used.

Uses sum 10a.

Начальное состояние суммирования, в нашем случае, — это 0.

Попробуем что-нибудь посложнее. Давайте посчитаем минимум, максимум и среднее от введенных чисел. Сначала разберем, что входит в состояние  $Q$ . Очевидно, это композиция из текущих минимума, максимума, а также пары – суммы и количества прочитанных чисел, из которых  $g$  потом посчитает среднее. Однако, если поток совершенно пустой, мы не сможем посчитать ничего. Т.е.  $g$  в случае, если количество прочитанных чисел равно нулю, должна напечатать ошибку. Состояние  $Q$  можно формировать в виде структуры `struct`, а можно передавать его между вызовами в виде нескольких аргументов:

10c  $\langle \text{подсчет статистики } 10c \rangle \equiv$  (11a)

```
int stat(double min, double max, double sum, int n) {
    double a;
    if ( scanf("%lf", &a) == 1 )
        return stat(a < min ? a : min,
                    a > max ? a : max,
                    sum + a,
                    n + 1);
    if ( n == 0 )
        return printf("Error: Empty stream!\n");
    else
        return printf("min=%lf, max=%lf, avg=%lf\n", min, max, sum/n);
}
```

Defines:

stat, used in chunk 11a.

Uses sum 10a.

Остается вопрос, что считать начальным состоянием  $q_0$ ? Если с суммой и количеством все понятно, что в качестве начального минимума можно предложить аналог  $+\text{inf}$ , а в качестве начального максимума —  $-\text{inf}$ . Поскольку любое другое число будет, соответственно, меньше начального минимума и больше начального максимума. К сожалению, таких стандартных констант для типа `double` нет, хотя тип `double` включает в себя эти значения и они могут появляться в процессе вычислений. Но, можно воспользоваться свойством функции стандартной библиотеки `atof`, которая преобразовывает строку в число с плавающей точкой двойной точности, но может также принимать значения `"inf"` и `"-inf"`, возвращая нам соответствующие "бесконечности".

```
11a  <stat.c 11a>≡
      #include<stdio.h>
      // нужно для atof
      #include<stdlib.h>

      <подсчет статистики 10c>

      int main() {
          return stat(atof("inf"), atof("-inf"), 0, 0);
      }
```

Defines:

`main`, never used.

Uses `stat 10c`.

Другим способом избежать проблем с поиском подходящего  $q_0$  для минимума и максимума является является его неявная ивалидация. Т.е. не будем полагаться на  $q_0$  вообще, а начнем считать минимум/максимум с первого элемента потока:

```
11b  <подсчет статистики с невалидным начальным состоянием 11b>≡ (11c)
      int stat2(double min, double max, double sum, int n) {
          double a;
          if ( scanf("%lf", &a) == 1 )
              return stat2( (!n || a < min)? a : min,
                             (!n || a > max)? a : max,
                             sum + a,
                             n + 1);

          if ( n == 0 )
              return printf("Error: Empty stream!\n");
          else
              return printf("min=%lf, max=%lf, avg=%lf\n", min, max, sum/n);
      }
```

Defines:

`stat2`, used in chunk 11c.

Uses `sum 10a`.

Соответственно, вся программа может выглядеть так:

```
11c  <stat2.c 11c>≡
      #include <stdio.h>

      <подсчет статистики с невалидным начальным состоянием 11b>

      int main() {
          return stat2(-100313, 231987, 0, 0);
      }
```

Defines:

`main`, never used.

Uses `stat2 11b`.

Здесь намеренно указаны произвольные начальные значения для минимума и максимума, чтобы показать, что они ни на что не влияют.

## Задание

1. Пусть тип элементов потока – это тройка действительных чисел  $\langle x, y, m \rangle$  обозначающих координату и массу точечного тела. Найти центр масс и общую массу точек.
2. Пусть тип элементов потока – это множество измерений пары значений напряжения и тока через резистор  $\langle U, I \rangle$ . Найдите сопротивление резистора и точность (средне-квадратичное отклонение).

## 2.2. Однопроходные алгоритмы с памятью

Но, бывают задачи, когда  $O(1)$  для  $Q$  явно недостаточно. Рассмотрим, например, подсчет значения медианы множества введенных чисел. Медиана – это такое число из множества, которое делит данное множество на 2 равные (приблизительно) части – числа, которые больше медианы и числа, которые меньше (или равны) медиане. Очевидно, что расчет требует хранения всех чисел (или, хотя бы, их распределения), а на это уже потребуется  $O(n)$  памяти. Если мы на каждом этапе чтения потока будем хранить отсортированный список уже введенных значений, то медианой будет элемент из середины этого списка. Поддерживать отсортированный список достаточно просто – нужно вставлять новый элемент в положенное ему порядком место. Вставка проще всего производится если список односвязный.

Односвязный список – это чуть более сложная структура чем массив, является динамической, но позволяет быструю вставку элемента в начало списка, а также предоставляет возможности ”неограниченного” роста. Т.е. список, в отличие от массива, не занимает памяти, если в нем нет элементов и не требует априорных знаний о возможном размере.

Каждый элемент списка представляется в виде структуры:

12  $\langle \text{определение списка } 12 \rangle \equiv$  (14b)

```
typedef struct _list {
    double e;
    struct _list* next;
} List;
```

Defines:

List, used in chunks 13, 14, 26–29, and 44–48.

которая содержит в себе сам элемент (поле `e`) и указатель на следующий элемент списка (поле `next`). Если элемент является последним, то указатель `next` принимает специальное значение `NULL`, т.е. *нулевой указатель*. Заметьте, что односвязный список является рекурсивной структурой.

Идея поддержания отсортированного списка очень проста. Допустим, у нас есть указатель на голову списка `lst` и нам нужно вставить в список число `e`. Если список пуст (`lst == NULL`), либо когда в голове списка содержится элемент больший или равный `e` (`lst->e >= e`), мы создаем новый элемент `new` и заполняем его поле `new->e = e`. Новый элемент `new` теперь должен стать новой головой, а старый список (если он был) необходимо прицепить ему в качестве хвоста: `new->next = lst`. Т.е. в хвосте списка содержатся элементы со значением  $\geq e$ . Если же в голове списка элемент меньше вставляемого, то нужно подыскать ему место в хвосте списка. Для этого нужно вызвать процедуру вставки еще раз, но уже для хвоста. Поскольку у хвоста может "отрасти" новая голова, то процедура вставки всегда должна возвращать голову нового списка.

13a  $\langle \text{вставка элемента в отсортированный список 13a} \rangle \equiv$  (14b)

```
List* insert_sort(List* lst, double e) {
    if ( lst == NULL || lst->e >= e ) {
        List* new = malloc(sizeof(List));
        new->e = e;
        new->next = lst;
        return new;
    } else {
        lst->next = insert_sort(lst->next, e);
        return lst;
    }
}
```

Defines:

`insert_sort`, used in chunk 14b.

Uses List 12 26a 44a.

Интересной выглядит задача поиска середины списка для вычисления медианы. Можно заняться подсчетом общего числа элементов в списке, потом поделить его пополам и пробежаться от начала списка, отсчитывая половину от общего числа, по достижении – вернуть полученный элемент. Но можно поступить проще, так сказать, однопроходно. Если ввести два указателя и первый, назовем его `slow` каждый раз смещать на один элемент, а второй – `fast` на два элемента списка вперед, то к тому моменту, когда `fast` доберется до конца списка, `slow` как раз окажется посередине:

13b  $\langle \text{поиск середины 13b} \rangle \equiv$  (14b)

```
List* median(List* slow, List* fast) {
    if ( !fast || !fast->next )
        return slow;
    return median(slow->next, fast->next->next);
}
```

Defines:

`median`, used in chunk 14b.

Uses List 12 26a 44a.

Ну и, поскольку, память у нас под список выделяется динамическая, неплохо бы ее в конце всех вычислений очистить:

14a  $\langle \text{очистка списка 14a} \rangle \equiv$  (14b)

```
void list_free(List* lst) {
    if (lst) {
        List* next = lst->next;
        free(lst);
        list_free(next);
    }
}
```

Defines:

`list_free`, used in chunk 14b.

Uses List 12 26a 44a.

Давайте теперь соединим все запчасти в программу:

14b  $\langle \text{stat3.c 14b} \rangle \equiv$

```
#include<stdio.h>
#include<stdlib.h>

 $\langle \text{определение списка 12} \rangle$ 
 $\langle \text{вставка элемента в отсортированный список 13a} \rangle$ 
 $\langle \text{поиск середины 13b} \rangle$ 
 $\langle \text{очистка списка 14a} \rangle$ 

int stat3(List* lst) {
    double a;
    if ( scanf("%ld", &a) == 1 )
        return stat3(insert_sort(lst, a));
    if ( ! lst )
        printf("Error: list is empty!\n");
    else
        printf("median: %lf\n", median(lst, lst)->e)
    list_free(lst) // освобождаем память
    return 0;
}

int main() {
    return stat3(NULL);
}
```

Defines:

`main`, never used.

`stat3`, never used.

Uses insert\_sort 13a, List 12 26a 44a, list\_free 14a, and median 13b.

## Задание

1. Немного оптимизируйте список, сохраняя *распределение* элементов. Т.е. в структуру добавьте поле с количеством – если элемент в потоке встречался  $n$  раз. Понятно, что поиск медианы будет уже не таким красивым, а в состоянии придется "тащить" общее количество прочитанных элементов.
2. Напишите "фильтр" уникальных элементов потока. Программа должна печатать число из потока, если оно встретилось в первый раз. Придется хранить список всех встреченных чисел и каждый раз проверять, не находится ли следующий элемент из потока уже в данном списке. Если число уже есть в списке – читать следующий элемент и ничего не печатать.

## Глава 3

# Численные методы

Немного потренируемся писать алгоритмы численных методов решения двух видов задач, которые часто возникают на практике. Поиск корня нелинейного уравнения и поиск минимума нелинейной функции.

Сразу оговорюсь, что данная глава не претендует на полноту математического описания данных методов и даже не на обзор. Если с поиском корней (нулей) нелинейной функции более-менее все понятно, то поиск минимума – это практически творческая задача.

### 3.1. Поиск корня нелинейного уравнения

Существует несколько распространенных способов решения уравнения вида  $f(x) = 0$ . Все они *итерационные*, т.е. получают значение  $x$  в результате повторяющихся итераций, на каждой итерации  $x_n$  постепенно приближается к искомому значению. Это означает, что за конечное время можно получить только приближенный результат. Поэтому количество итераций определяется необходимой точностью ответа. Методы различаются требованиями к функции  $f$  и *сходимостью* – скоростью, с которой увеличивается точность на каждой итерации.

Но, так или иначе, **все методы** требуют аналитической подготовки к решению. В любом случае, требуется выбора окрестности  $\delta$  и доказательства, что эта окрестность содержит корень, а также, что функция достаточно хорошо себя ведет в этой окрестности, возможно даже потребуются аналитическое вычисление производной функции  $f'(x)$ .

Например, известный метод Ньютона (касательных) обладает самой лучшей сходимостью. Однако, налагает серьезные ограничения на вид функции и ее производные первого и второго порядка (не говоря уже об их существовании) внутри окрестности корня. И в некоторых случаях, например, когда производная в корне обращается в 0 (кратный корень), скорость теряется.

Другие методы хуже в плане сходимости, но уже менее требовательны к виду функции. Самым простым и "робастным" (от robust – надежный) является метод поиска корня делением пополам – *дихотомии*. Для его успешной работы только требуется, чтобы в начальной окрестности поиска – некотором интервале  $[a, b]$  – производная  $f'(x)$  не меняла бы знак, а на концах этого интервала знак функции был бы разным  $f(a) \cdot f(b) < 0$ .

Вот его мы и рассмотрим.

Итак, нам для вычисления корня потребуется тело функции  $f$ , начальный интервал  $[a, b]$  и предельная точность  $\epsilon$ . Максимальное количество итераций можно прикинуть сразу – это  $\log_2|b - a| - \log_2\epsilon$ , поскольку идея метода состоит в том, чтобы сокращать область поиска каждую итерацию в 2 раза, а минимальным размером области является  $\epsilon$ .

Для того чтобы понять, где может быть корень, нам потребуется значение функции  $f$  в трех точках:  $f(a)$ ,  $f(b)$  и  $f(c)$ , где  $c = (a + b)/2$  – центр интервала  $[a, b]$ . Обозначим  $f_a$ ,  $f_b$  и  $f_c$  – соответствующие значения функции Тогда суть метода можно выразить следующими строками:



17a  $\langle \text{суть метода 17a} \rangle \equiv$  (17b)

```

if ( b - a < epsilon ) return c;
if ( fa * fc < 0 )  $\langle \text{ищем корень в левой половине 17c} \rangle$ 
else if ( fc * fb < 0 )  $\langle \text{ищем корень в правой половине 17d} \rangle$ 
else if ( fc == 0 ) return c;
else if ( fa == 0 ) return a;
else if ( fb == 0 ) return b;
else  $\langle \text{корень недостижим 17e} \rangle$ 

```

Вариант конень недостижим возникает, когда знаки  $fa$ ,  $fb$  и  $fc$  совпадают. В этом случае мы перестаем "бороться" за корень и просим у пользователя уточнить область поиска.

Сама функция поиска (алгоритм) будет оформлена в виде рекурсивной функции `sol`. Она должна зависеть от функции  $f$ , которая, в свою очередь, принимает на вход `double` и возвращает `double`: `double f(double)`, начала и конца области поиска `double a`, `b`, и точности `epsilon`.

Рекурсия, собственно, будет происходить в ветках `ищем корень в левой половине` и `ищем корень в правой половине`. Но стоит заметить, что в рекурсивном вызове из трех вычисленных значений  $fa$ ,  $fb$  и  $fc$ , по крайней мере, два потребуются точно. Поэтому имеет смысл их не вычислять заново, а передать в виде параметров в функцию. Т.е. добавим еще аргументы `double fa`, `fb`:

17b  $\langle \text{функция поиска 17b} \rangle \equiv$  (18b)

```

double sol(double f(double), double a, double b, double epsilon,
           double fa, double fb) {
    /* вычисляем только c и fc остальное дано "свыше" */
    double c = (a+b)/2;
    double fc = f(c);

     $\langle \text{суть метода 17a} \rangle$ 
}

```

Defines:

`sol`, used in chunks 17 and 18b.

Теперь распишем ветки с рекурсивным вызовом для поиска в интервалах  $[a, c]$  и  $[c, b]$ . Они выглядят достаточно просто:

17c  $\langle \text{ищем корень в левой половине 17c} \rangle \equiv$  (17a)

```

return sol(f, a, c, epsilon, fa, fc);

```

Uses `sol 17b`.

17d  $\langle \text{ищем корень в правой половине 17d} \rangle \equiv$  (17a)

```

return sol(f, c, b, epsilon, fc, fb);

```

Uses `sol 17b`.

Тут мы не забываем передать уже вычисленные значения  $f$ .

Чтобы указать вызывающей стороне, что корень не найден (точнее, это мы подозреваем, что корня может не быть, хотя ему никто не запрещает находиться в интервале), будем возвращать специальное значение `nan`.

17e  $\langle \text{корень недостижим 17e} \rangle \equiv$  (17a)

```

return atof("nan");

```

Теперь попробуем применить наш алгоритм к какому-нибудь уравнению, которое не имеет алгебраических решений (например, с участием трансцендентных функций):

$$e^x - 5x = 0$$

Оформим его в виде функции:

18a  $\langle$ исследуемая функция 18a $\rangle \equiv$  (18b)

```
double func(double x) {
    return exp(x) - 5*x;
}
```

Defines:

func, used in chunk 18b.

Напишем программу-обертку, которая будет принимать в качестве аргументов диапазон поиска и точность:

18b  $\langle$ sol.c 18b $\rangle \equiv$

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>

 $\langle$ исследуемая функция 18a $\rangle$ 

 $\langle$ функция поиска 17b $\rangle$ 

int main(int argc, char* argv[]) {
    double a = atof(argv[1]);
    double b = atof(argv[2]);
    double epsilon = atof(argv[3]);
    printf("Result = %.20lf\n", sol(func, a, b, epsilon, func(a), func(b)));
}
```

Defines:

main, never used.

Uses func 18a and sol 17b.

Для линковки программы нам потребуется математическая библиотека `libm`, в которой даны определения функций `exp`, `sin`, `cos` и многих других. Они не входят в стандартную библиотеку `libc`, их придется подключать в явном виде. Линковка с математической библиотекой выглядит вот так:

```
gcc -o sol -O3 -lm sol.c
```

Здесь опция `-lm` добавляет тело нужной нам библиотеки. Опция `-O3` относится не к линковке, а к компиляции и указывает компилятору строить оптимизированный код по скорости и по размеру.

Как раньше говорилось, поиск корня всегда требует дополнительной аналитической работы по определению начальной области поиска. Для этого нужно себе хорошо представлять, как ведет себя функция  $f(x)$ . Например, если мы рассматриваем  $f(x) = e^x - 5x$ , то при  $x \rightarrow \pm\infty$  значение  $f(x) \rightarrow +\infty$ . Также можно легко посчитать, что  $f(x)$  достигает минимума в точке  $\log 5$ , а само значение минимума  $-5 - 5 \log 5 < 0$ . Следовательно, функция имеет 2 корня слева и справа от  $x = \log 5 \approx 1.6$ .

Попробуем поискать первый корень с точностью до  $10^{-6}$ :

```
sol 0 1.6 1e-6
```

Получаем вот такой результат:

```
Result = 0.25917091369628908470
```

## Задание

1. Добавьте в функцию поиска `sol` подсчет числа шагов, требующихся для достижения нужной точности. Сверьте с теоретическими оценками.
2. Реализуйте функцию поиска методом Ньютона `solnewton`:  $x_{n+1} = x_n - f(x_n)/f'(x_n)$ . Для этого добавьте реализацию функции  $f'(x)$  и передайте ее в качестве дополнительного параметра в функцию `solnewton`. При этом, в отличие от `sol`, можно избавиться от параметров `b`, `fa` и `fb`, поскольку данному методу достаточно одной начальной точки приближения. Сравните количество шагов, требующихся для достижения нужной точности функциям `sol` и `solnewton` на одинаковых  $f(x)$ .

## 3.2. Поиск минимума

Обобщить методы поиска корня на уравнения с несколькими переменными практически невозможно, поскольку в многомерном случае "корнем", вообще говоря, может являться целая поверхность.

Но вот минимум (экстремум) функции  $f(x)$  вполне можно (численно) искать в случае, если  $x$  – вектор. Если же функция однопараметрическая и при том дифференцируема аналитически, то поиск минимума можно свести к поиску корня  $f'(x) = 0$ .

Одним из простых (но не самых эффективных) способов поиска экстремума является метод *градиентного спуска*. Идея метода проста. Если мы стоим на склоне горы, то чтобы попасть в долину, скорее всего, нужно двигаться в сторону, противоположную направлению *градиента*. Градиентом функции  $\nabla f(\mathbf{x})$  в точке  $\mathbf{x}$  называется вектор, составленный из частных производных функции по компонентам вектора:

$$\nabla f(\mathbf{x}) = \left\{ \frac{\partial}{\partial x_1} f, \frac{\partial}{\partial x_2} f, \dots \right\}$$

Этот вектор лежит в (гипер)плоскости касательной к функции  $f$  в точке  $x$ . И направление этого вектора соответствует направлению наибольшего роста функции в данной точке.

Признаком окончания движения может служить тот факт, что градиент обнуляется. Т.е. в точке, по крайней мере локального минимума  $\nabla f(\mathbf{x}) = 0$ .

Понятно, что как и в случае с корнями, данный метод может привести только к одному (локальному) минимуму, в долине которого находится наше начальное приближение. А поиск абсолютного (глобального) минимума (экстремума) в общем случае — задача вообще нерешаемая. Обычно ищется несколько минимумов с множеством случайно (равномерно, на сколько это возможно) разбросанных начальных приближений и из них выбирается самый минимальный. Но чем больше пространство поиска, тем меньше вероятность, что найдется "глобальный" минимум.

Итак, направление движения приблизительно определено, остается вопрос: на сколько шагать? С одной стороны, шагать нужно быстро, чтобы скорее добраться до минимума. С другой стороны, если шагать далеко, то можно "проскочить" минимум. Поэтому выбирается достаточно маленький коэффициент  $\alpha$ , а координаты следующего приближения вычисляются как

$$\mathbf{x}_{n+1} = \mathbf{x}_n - \alpha \nabla f(\mathbf{x}_n)$$

Значение  $\alpha$  можно подобрать эмпирически в качестве константы, а можно вычислять на каждом шаге. Стоит заметить, что размерность  $[\alpha]$  совпадает с размерностью обратной второй производной  $(\partial^2 f / \partial x^2)^{-1}$ . Поэтому было бы неплохо заручиться оценкой предельных значений *Гессiana* функции  $f$ , которые могут достигаться в точках экстремума. И чем они выше, тем острее пики и впадины и тем меньше должны быть  $\alpha$ , чтобы обеспечить сходимость. Иначе, мы рискуем постоянно проскакивать минимумы.

Методов вычисления значений  $\alpha$  существует множество, однако ни один из них не лишен "произвола" в выборе параметров и так или иначе связан с предположительными знаниями о виде  $f(\mathbf{x})$ .

Для простоты остановимся пока на постоянном коэффициенте спуска  $\alpha = const$

Давайте пока представим итерацию метода градиентного спуска в виде функции `grad`. Для вычисления итерации нам потребуется вектор  $\mathbf{x}$ . Если мы хотим иметь дело с произвольной размерностью, то придется ее передавать, скажем, в параметре `px`. Также нам будут нужны параметры нашего

”произвола”  $\alpha$ ,  $\epsilon$  – предельная точность,  $\delta$  – размер области, на которой будет вычисляться градиент, и, собственно, сама функция  $f$ , только она уже будет принимать на вход вектор (массив `double*` и размерность `int`)

20a  $\langle \text{градиентный спуск 20a} \rangle \equiv$  (23)

```
int grad(double *x, int nx, double alpha, double delta,
        double epsilon, double f(double*, int)){
     $\langle \text{локальные переменные 20b} \rangle$ 
     $\langle \text{вычисляем градиент в точке x 20d} \rangle$ 
     $\langle \text{проверка условий останова 20e} \rangle$ 
     $\langle \text{вычисляем следующее приближение 20f} \rangle$ 
    printf("x[0]=%lf, x[1]=%lf, error=%lf\n", x[0], x[1], f(x, nx));
    return grad(x, nx, alpha, delta, epsilon, f);
}
```

Defines:

`grad`, used in chunk 23.

Перед тем, как вычислить градиент, нам потребуется вычислить значение  $f(x)$ .

20b  $\langle \text{локальные переменные 20b} \rangle \equiv$  (20a) 20c  $\triangleright$

```
double fx = f(x, nx);
```

Сам градиент нам нужно где-то запомнить, для этого аллоцируем массив (динамически), а также будем считать квадрат модуля градиента:

20c  $\langle \text{локальные переменные 20b} \rangle + \equiv$  (20a)  $\triangleleft$  20b

```
double *g = calloc(sizeof(double), nx);
double g_module = 0;
```

20d  $\langle \text{вычисляем градиент в точке x 20d} \rangle \equiv$  (20a)

```
int i;
for (i = 0; i < nx; i++) {
    x[i] += delta;
    g[i] = (f(x, nx) - fx)/delta;
    x[i] -= delta;
    g_module += g[i]*g[i];
}
```

Понятно, что мы вычисляем не градиент, а *приближение градиента* и точность его зависит от выбора параметра  $\delta$ . Мы переиспользуем компоненты вектора  $x$  для выбора шагов в направлении каждой из координат.

20e  $\langle \text{проверка условий останова 20e} \rangle \equiv$  (20a)

```
if (g_module < epsilon*epsilon) {
    free(g);
    return 0;
}
```

Если модуль градиента достигает некоторого предельного значения, то останавливаемся. Можно остановиться и по другим критериям, например, когда  $\alpha \|\nabla f\| < \epsilon$ .

20f  $\langle \text{вычисляем следующее приближение 20f} \rangle \equiv$  (20a)

```
for (i = 0; i < nx; i++) {
    x[i] -= alpha * g[i];
}
free(g);
```

Отлично. Приступаем к самому интересному. Предположим, что мы ставим некоторый эксперимент и хотим подобрать параметры модели, описывающей какое-то явление. В распоряжении у нас есть параметризованная модель  $= F(\vec{\theta}, \mathbf{x})$  которая предсказывает скалярное значение  $y$  от некоторого вектора значений наблюдаемых параметров  $\mathbf{x}$  и (ненаблюдаемых) параметров модели  $\vec{\theta}$ . А также некоторое количество измерений (естественно, с погрешностью):

$$S = \{(\mathbf{x}_1, y_1), (\mathbf{x}_2, y_2), \dots, (\mathbf{x}_N, y_N)\}$$

Задача подбора параметров модели, наилучшим образом описывающей зависимость  $y_i$  от  $\mathbf{x}_i$  обычно сводится к минимизации ошибки  $\varepsilon$ , которая определяется как

$$\varepsilon_{\vec{\theta}} = \frac{1}{N} \sum_i (F(\vec{\theta}, \mathbf{x}_i) - y_i)^2$$

Это средний квадрат отклонения значений, полученных моделью, от реально измеряемых.

Собственно, набор оптимальных параметров  $\vec{\theta} : \varepsilon \rightarrow \min$  мы и будем искать методом градиентного спуска.

Для этого нам потребуются результаты  $N$  измерений наблюдаемых параметров  $x_i$  и величины  $y$ . Их можно хранить в виде нескольких массивов: `double x1[N], double x2[N], ... double y[N]`, а можно сделать и массив структур. Для простоты измерения будут храниться глобально.

Нужна будет функция вычисления ошибки `mse` (от Mean Squared Error – квадрат стандартной ошибки), которая будет принимать на вход массив параметров модели `double* p` их количество `int np`

21a  $\langle \text{функция вычисления ошибки 21a} \rangle \equiv$  (23)

```
double mse(double* p, int np) {
    double e = 0;
    int i;
    for ( i = 0; i < N; i++ ) {
        double d;
         $\langle \text{вычислить отклонение } d \text{ 22b} \rangle$ 
        e += d*d;
    }
    return e/N;
}
```

Defines:

`mse`, used in chunk 23.

Uses N 22a.

Сама модель нам потребуется исключительно здесь и больше нигде, поэтому ее можно не передавать в качестве аргумента и реализовать на месте, прямо внутри `mse`, либо оформить в виде глобальной функции.

Функцию `mse`, в свою очередь, мы будем передавать в качестве аргумента `f` нашей универсальной функции `grad`, которая будет заниматься подбором параметров.

Теперь пришло время определиться с моделью. Допустим, мы измеряем ток  $I$ , протекающий через какой-то элемент цепи, в зависимости от напряжения  $U$ . В результате измерений мы получили множество пар точек  $(I_i, U_i)$ . Занесем их в нашу структуру:

21b  $\langle \text{структура измерения 21b} \rangle \equiv$  (23)

```
typedef struct {
    double I;
    double U;
} Measure;
```

Defines:

`Measure`, used in chunk 22a.

Аллоцируем массив для наших измерений:

22a  $\langle \text{массив измерений 22a} \rangle \equiv$  (23)

```
#define MAX_MEASURES 100
int N;
Measure measures[MAX_MEASURES];
```

Defines:

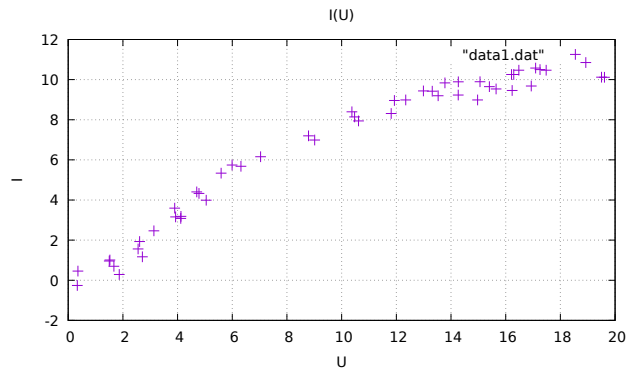
MAX\_MEASURES, used in chunk 23.

measures, used in chunks 22b and 23.

N, used in chunks 21a and 23.

Uses Measure 21b.

Значения измерений на плоскости ложатся как-то так:



На графике видно, что ток растет нелинейно с увеличением  $u$ , и при некоторых значениях  $u$  наступает "насыщение". Попробуем аппроксимировать данную зависимость моделью  $I(u) = I_0(1 - e^{-u/U_0})$  (если она, конечно, имеет физический смысл). Эта модель двухпараметрическая. Здесь параметр  $I_0$  определяет "предельный" ток, а параметр  $U_0$  – начало наступления режима насыщения. Добавим вычисления модели в функцию вычисления ошибки для  $i$ -го измерения:

22b  $\langle \text{вычислить отклонение d 22b} \rangle \equiv$  (21a)

```
d = p[0]*(1 - exp(-measures[i].U/p[1])) - measures[i].I;
```

Uses measures 22a.

Здесь мы используем параметры  $p[0]$  и  $p[1]$  из массива для описания модели.

Напишем программу-обертку для ввода значений измерений и начальных параметров модели, с которых мы начнем спуск, а также параметры градиентного спуска  $\alpha$  и  $\epsilon$ . Параметры модели  $I_0$  и  $U_0$  мы будем передавать в качестве аргументов, а экспериментальные данные будем подавать на `stdin` в виде множества пар действительных значений, разделенных пробелом.

23

*<grad.c 23>*≡

```
#include <stdio.h>
#include <stdlib.h>
#include <math.h>
```

*<структура измерения 21b>*

*<массив измерений 22a>*

*<функция вычисления ошибки 21a>*

*<градиентный спуск 20a>*

```
int main(int argc, char* argv[]) {
    double alpha, epsilon;
    double p[2]; /* параметры модели */
    p[0] = atof(argv[1]);
    p[1] = atof(argv[2]);
    alpha = atof(argv[3]);
    epsilon = atof(argv[4]);
    /* Читаем измерения */
    N = 0;
    while (N < MAX_MEASURES
        && scanf("%lf %lf", &measures[N].U, &measures[N].I) == 2) N++;
    /* Запускаем градиентный спуск */
    grad(p, 2, alpha, epsilon, epsilon, mse);
    /* печатаем результат и финальную ошибку */
    printf("Params: %0.10lf, %0.10lf, error=%lf\n", p[0], p[1], mse(p, 2));

}
```

Defines:

`main`, never used.

Uses `grad` 20a, `MAX_MEASURES` 22a, `measures` 22a, `mse` 21a, and `N` 22a.

Компилировать данную программу нужно также с опцией -O3:

```
gcc -o grad -O3 grad.c -lm
```

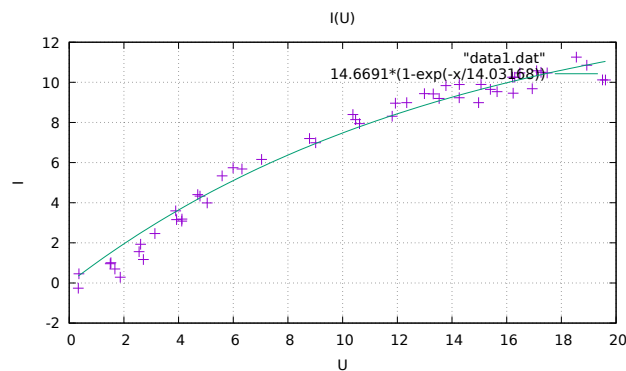
Попробуем запустить ее с приближительными начальными параметрами (в файле `gen.dat` находятся наши измерения).

```
./grad 10 10 0.00001 0.000001 <gen.dat
```

В результате получаем:

Params: 14.6691, 14.03168, error=0.274655

Смотрим на график:



## Задание

- Попробуйте различные данные и различные модели. Оцените количество шагов до получения нужной точности, всегда ли она достигается? Как на сходимость влияют параметры  $\alpha$ ,  $\delta$ ,  $\epsilon$ ?
- Добавьте в алгоритм `grad` "отжиг" – постепенное уменьшение  $\alpha$  с каждой итерацией. Введите параметр отжига  $m$ :  $\alpha_{n+1} = m\alpha_n$  и проверьте сходимость с ним.



## Глава 4

# Сортировка

### 4.1. Сортировка слиянием односвязного списка

Рассмотрим реализацию быстрой сортировки односвязного списка. Она хороша тем, что имеет наименьшее время работы  $O(N \log N)$ , не требует дополнительной памяти и является стабильной (сохраняет порядок равных элементов).

**Определение односвязного списка.** Это структура, которая включает в себя полезную нагрузку (в нашем случае — строка `str`) и указатель на хвост списка `next`.

26a  $\langle linked\ list\ definition\ 26a \rangle \equiv$  (30)

```
typedef struct _List {
    char* str;
    struct _List* next;
} List;
```

Defines:

List, used in chunks 13, 14, 26–29, and 44–48.

Общая идея быстрой сортировки односвязного списка `lst` выглядит вот так:

26b  $\langle sort\ definition\ 26b \rangle \equiv$  (30)

```
void llsort(List** lst) {
    List* a = NULL;
    List* b = NULL;
    if( !*lst || !(*lst)->next )
        return;
     $\langle split\ 26c \rangle$ 
     $\langle sort\ splits\ 26d \rangle$ 
     $\langle merge\ 26e \rangle$ 
}
```

Defines:

llsort, used in chunks 26d and 29d.

Uses List 12 26a 44a.

Для начала мы убеждаемся, что список `lst` имеет смысл сортировать: он не пуст и имеет более одного элемента. В противном случае — выходим, список уже отсортирован.

Затем идет разделение списка на 2 (почти) равные части `a` и `b`:

26c  $\langle split\ 26c \rangle \equiv$  (26b)

```
llsplit(&a, &b, *lst);
```

Uses llsplit 27b.

Затем, эти части сортируются

26d  $\langle sort\ splits\ 26d \rangle \equiv$  (26b)

```
llsort(&a);
llsort(&b);
```

Uses llsort 26b.

И уже отсортированные части склеиваются с сохранением порядка:

26e  $\langle merge\ 26e \rangle \equiv$  (26b)

```
llmerge(lst, a, b);
```

Uses llmerge 27a.

Теперь введем определения функций. Начнем со склейки двух односвязных списков с сохранением порядка `llmerge`. Сначала действуем в предположении, что оба списка не пусты, поэтому мы сравниваем их головные элементы. В результирующий список `lst` записываем наименьшую голову одного из списков и рекурсивно продолжаем склейку между с его хвостом и другим списком, в качестве результата подставляя адрес следующего элемента результирующего списка. Если один из списков пуст, то просто оставляем в качестве результата оставшийся список.

27a  $\langle merge\ definition\ 27a \rangle \equiv$  (30)

```
void llmerge(List** lst, List* a, List* b) {
    if ( a && b ) {
        if ( strcmp(a->str, b->str) <= 0 ) {
            *lst = a;
            llmerge(&a->next, b, a->next);
        } else {
            *lst = b;
            llmerge(&b->next, b->next, a);
        }
    } else {
        *lst = a ? a : b;
    }
}
```

Defines:

`llmerge`, used in chunk 26e.

Uses List 12 26a 44a.

Разберем теперь идею разделения односвязного списка на две равные части. Опять будем решать задачу рекурсивно. Предположим, что у нас уже есть два разделенных списка `a` и `b` — адреса их концов. Если разделяемый список `lst` не пуст, поставляем его голову в конец одного из списков (для определенности, пусть это будет `a`), и далее, рекурсивно вызываем разделение оставшегося хвоста, но при этом меняем `a` и `b` местами. Теперь концом первого списка `a` будет конец второго, а концом второго — конец первого (естественно, с учетом того, что он уже немножко попрос).

27b  $\langle split\ definition\ 27b \rangle \equiv$  (30)

```
void llsplit(List** a, List** b, List* lst) {
    if ( lst ) {
        *a = lst;
        llsplit(b, &(*a)->next, lst->next);
    } else {
        *a = NULL;
        *b = NULL;
    }
}
```

Defines:

`llsplit`, used in chunk 26c.

Uses List 12 26a 44a.

Есть одно замечание, касающееся порядка списков **a** и **b** в рекурсивных вызовах `llsplit` и `llmerge`. Если при разделении мы постоянно меняли **a** и **b** местами, то тоже самое следует сделать и при склейке, чтобы сохранялся порядок равных элементов. Вот, в принципе, и все. Merge-sort на односвязных списках в трех простых функциях.

Не забываем про функции для создания и удаления списков. При добавлении элемента в список потребуется дополнительная память, поэтому нужно обработать ситуации, когда памяти не хватает.

28a `<insert definition 28a>`≡ (30)

```
int llinstert(List** lst, const char* str) {
    List* new;
    if ( new = malloc(sizeof(List)) ) {
        if ( new->str = strdup(str) ) {
            new->next = *lst;
            *lst = new;
            return 0;
        }
        free(new);
        return -1;
    }
    return -2;
}
```

Defines:

`llinsert`, used in chunk 29c.

Uses List 12 26a 44a.

Освобождение памяти, занятой списком происходит в обратном порядке. Мы сохраняем значения поля `next` на случай, если после освобождении головы `lst` оно будет разрушено. Можно, конечно, сначала рекурсивно вызвать `llfree` на хвосте, потом уже освободить голову, но при этом рекурсия перестанет быть хвостовой и может быть неоптимизирована компиляторами.

28b `<free definition 28b>`≡ (30)

```
void llfree(List* lst) {
    if ( lst ) {
        List* next = lst->next;
        free(lst->str);
        free(lst);
        llfree(next);
    }
}
```

Defines:

`llfree`, used in chunk 29f.

Uses List 12 26a 44a.

Заметьте, что `llfree` хоть и удаляет, но не делает полезную модификацию `lst`, поэтому, в отличие от других функций - модификаторов, мы передаем ей только указатель на `List`.

**Тестирование** Для тестирования нам потребуется функция вывода элементов односвязного списка `llprint`. Также выполним ее в виде рекурсии.

28c `<llprint definition 28c>`≡ (30)

```
void llprint(List* lst) {
    if (lst) {
        printf("%s", lst->str);
        llprint(lst->next);
    }
}
```

Defines:

`llprint`, used in chunk 29e.

Uses List 12 26a 44a.

Основная программа для тестирования состоит из четырех основных этапов:

1.  $\langle \text{read text into list } 29c \rangle$  — чтение из стандартного потока ввода,
2.  $\langle \text{sort list } 29d \rangle$  — сортировка,
3.  $\langle \text{output sorted list } 29e \rangle$  — вывод отсортированного списка в стандартный поток вывода
4.  $\langle \text{free list } 29f \rangle$  — освобождение ресурсов, связанных со списком.

29a  $\langle \text{main function } 29a \rangle \equiv$  (30)

```
int main() {
     $\langle \text{variables } 29b \rangle$ 
     $\langle \text{read text into list } 29c \rangle$ 
     $\langle \text{sort list } 29d \rangle$ 
     $\langle \text{output sorted list } 29e \rangle$ 
     $\langle \text{free list } 29f \rangle$ 
    return -1;
}
```

Defines:

`main`, never used.

Для работы нам потребуется указатель на список и буфер для чтения строк. Выбор размера буфера, конечно, зависит от природы строк. Здесь мы предполагаем, что размер строки (в байтах) не превышает 4 Кб. Для произвольных строк с неограниченной длиной такая реализация **не подойдет!**

29b  $\langle \text{variables } 29b \rangle \equiv$  (29a)

```
char buf[4096];
List* lst = NULL;
```

Uses List 12 26a 44a.

Чтение текста из потока по строкам — это функция `fgets`. Она принимает на вход помимо адреса буфера также максимальный размер этого буфера. Если строка превысит этот размер, она обрежется и следующий вызов `fgets` вернет в буфере строку, начинающуюся с места, в котором прервался предыдущий вызов:

29c  $\langle \text{read text into list } 29c \rangle \equiv$  (29a)

```
while( fgets(buf, sizeof(buf), stdin) ) {
    llinsert(&lst, buf);
}
```

Uses `llinsert` 28a.

Сортировка — это просто вызов `llsort`:

29d  $\langle \text{sort list } 29d \rangle \equiv$  (29a)

```
llsort(&lst);
```

Uses `llsort` 26b.

Выводим отсортированный список на экран:

29e  $\langle \text{output sorted list } 29e \rangle \equiv$  (29a)

```
llprint(lst);
```

Uses `llprint` 28c.

И освобождаем его:

29f  $\langle \text{free list } 29f \rangle \equiv$  (29a)

```
llfree(lst);
```

Uses `llfree` 28b.

Целиком программа выглядит вот так:

```
30  <llsort.c 30>≡  
    #include<stdio.h>  
    #include<stdlib.h>  
    #include<string.h>  
    <linked list definition 26a>  
    <merge definition 27a>  
    <split definition 27b>  
    <sort definition 26b>  
    <insert definition 28a>  
    <free definition 28b>  
    <llprint definition 28c>  
    <main function 29a>
```

## 4.2. Задания

1. Проверьте на практике и докажите, что данный алгоритм является стабильным (сохраняет порядок равных элементов). Подсказка: *модифицируйте функцию сравнения так, чтобы она сравнивала первые  $n$  элементов.*
2. Докажите, что время сортировки не зависит от значений элементов и составляет  $O(N \log N)$ .
3. Измените алгоритм разделения `llsplit` так, чтобы он делал только одно изменение (разрезал список пополам в середине). Сохранится ли в этом случае свойство стабильности?

## Глава 5

# Двоичные деревья поиска

Рассмотрим новую структуру данных. *Двоичное Дерево Поиска*. Данная структура позволяет хранить (упорядоченное) множество элементов со временем доступа на запись, поиск, удаление в среднем  $O(\log N)$ .

Каждый узел дерева  $N$  представляет собой либо пустоту  $\emptyset$ , либо тройку  $\langle L \ x \ R \rangle$ , где  $x$  — значение узла,  $L$  и  $R$  — левое и правое поддеревья соответственно.

**Высота дерева.** Каждый узел дерева характеризуется своей высотой  $h$ :

$$\begin{aligned} h(N) &= h(\langle L \ x \ R \rangle) = 1 + \max(h(L), h(R)) \\ h(\emptyset) &= 0 \end{aligned}$$

Узел  $N$  с соответствующей ему высотой  $h$  будем обозначать  $N_h$ .

**Множество узлов.** Множеством узлов дерева  $T$  называется такое множество элементов  $x$ , которое состоит из значения узла  $T$  и значений его поддеревьев.

$$x \in T \Leftrightarrow T = \langle L \ y \ R \rangle \Rightarrow x = y \vee x \in L \vee x \in R$$

### 5.1. AVL - Дерево

**Определение.** AVL-деревом называется двоичное дерево поиска  $T$ , в котором для каждого поддерева  $N = \langle L_{h_l} \ x \ R_{h_r} \rangle$  верно неравенство  $|h_l - h_r| \leq 1$ . AVL-дерево является самобалансирующимся. Т.е. все функции его модификации сохраняют это условие.

**Функция вставки.** Вставку элемента  $x$  в дерево  $T$  будем обозначать функцией  $I(T, x)$ . Обычная вставка, без балансировки может выглядеть следующим образом:

$$\begin{aligned} I(T, x) : \\ I(\emptyset, x) &= \langle \emptyset \ x \ \emptyset \rangle \\ I(\langle L \ y \ R \rangle, x) &= \\ &\quad \begin{aligned} x \leq y &\Rightarrow \langle I(L, x) \ y \ R \rangle \\ x > y &\Rightarrow \langle L \ x \ I(R, x) \rangle \end{aligned} \end{aligned}$$

В самобалансирующемся дереве функция вставки должна в себя включать балансировку. Случай, когда  $T = \emptyset$  является тривиальным, мы его рассматривать не будем. Рассмотрим те случаи, когда условие баланса нарушается и требуется перебалансировка. Вариантов разбалансировки немного. Она может возникать когда высоты двух поддеревьев уже различаются на единицу, например  $T = \langle L_h \ R_{h+1} \rangle$  (для простоты будем опускать значение узла дерева), и добавление нового узла увеличивает дисбаланс:  $T \rightarrow \langle L_h \ R_{h+2} \rangle$ . В этом случае, для восстановления баланса, должен быть выбран новый корень из правого поддерева  $R$ .

Для этого допустим, что правое поддереву состоит из двух своих поддеревьев:  $R_{h+2} = \langle P_m Q_n \rangle$ , где, по условию,  $\max(m, n) = h + 1$ . Тогда, функция балансировки  $\mathcal{B}(T)$  — трансформация дерева  $T$  в новое дерево, в корне которого будет узел  $R$  будет выглядеть так:

$$\begin{aligned} \mathcal{B}(\langle L_h R_{h+2} \rangle) &= \mathcal{B}(\langle L_h \langle P_m Q_n \rangle \rangle) = \\ &= \langle \langle L_h P_m \rangle Q_n \rangle \end{aligned}$$

Если  $m < n$ , тогда  $m = h$ , а  $n = h + 1$ . В этом случае, выбор узла  $R$  в качестве корня нового дерева сохраняет баланс:

$$\langle \langle L_h P_h \rangle_{h+1} Q_{h+1} \rangle_{h+2}$$

В противном случае, дерево останется разбалансированным:

$$\langle \langle L_h P_{h+1} \rangle_{h+2} Q_h \rangle_{h+3}$$

Этот случай указывает нам, что мы “отхватили” от правого поддереву  $R$  слишком большой кусок и отдали его в левое поддерево. Попробуем забрать поменьше. Разделим дополнительно дерево  $P_{h+1} = \langle M_m N_n \rangle$  и передадим в левое поддерево только его левую часть  $M$ :

$$\begin{aligned} \mathcal{B}(T) &= \mathcal{B}(\langle L_h R_{h+2} \rangle) = \\ &= \mathcal{B}(\langle L_h \langle P_{h+1} Q_h \rangle \rangle) = \\ &= \mathcal{B}(\langle L_h \langle \langle M_m N_n \rangle_{h+1} Q_h \rangle \rangle) = \\ &= \langle \langle L_h M_m \rangle \langle N_n Q_h \rangle \rangle \end{aligned}$$

Поскольку  $\max(m, n) = h$ , то и  $\max(h, m) = h$  и  $\max(n, h) = h$ , а следовательно, высоты левого  $\langle L M \rangle$  и правого  $\langle N Q \rangle$  поддеревьев стали равны  $\max(h, m) + 1 = \max(n, h) + 1 = h + 1$  и все дерево оказалось опять сбалансированным.

Если обобщить все варианты разбалансировки слева и справа, то получится 4 варианта реализации функции  $\mathcal{B}(T)$ :

$$\begin{aligned} \mathcal{B}(T) = \\ T = \langle L_h R_{h+2} \rangle \Rightarrow \\ \begin{aligned} R &= \langle P_h Q_{h+1} \rangle \Rightarrow \langle \langle L_h P_h \rangle Q_{h+1} \rangle \\ R &= \langle \langle M_m N_n \rangle Q_h \rangle \Rightarrow \langle \langle L_h M_m \rangle \langle N_n Q_h \rangle \rangle \end{aligned} \\ T = \langle L_{h+2} R_h \rangle \Rightarrow \\ \begin{aligned} L &= \langle P_{h+1} Q_h \rangle \Rightarrow \langle P_{h+1} \langle Q_h R_h \rangle \rangle \\ L &= \langle K_h \langle M_m N_n \rangle \rangle \Rightarrow \langle \langle K_h M_m \rangle \langle N_n R_h \rangle \rangle \end{aligned} \end{aligned}$$

Примечательным является то, что можно осуществлять балансировку, имея только балансы деревьев. Будем считать баланс дерева  $b$  как разницу высот поддеревьев: положительным, если высота правого поддереву больше левого, отрицательным в обратном случае, и равным нулю, если высоты равны. Обозначать баланс дерева будем верхним индексом  $T^b$ , чтобы не путать с высотой.

$$\begin{aligned} \mathcal{B}(T) = \\ T = \langle L R \rangle^2 \Rightarrow \\ \begin{aligned} R &= \langle P Q \rangle^1 \Rightarrow \langle \langle L P \rangle^0 Q \rangle^0 \\ R &= \langle \langle M N \rangle^b Q \rangle^{-1} \Rightarrow \langle \langle L M \rangle^{\min(-b, 0)} \langle N Q \rangle^{\max(-b, 0)} \rangle^0 \end{aligned} \\ T = \langle L R \rangle^{-2} \Rightarrow \\ \begin{aligned} L &= \langle P Q \rangle^{-1} \Rightarrow \langle P \langle Q R \rangle^0 \rangle^0 \\ L &= \langle K \langle M N \rangle^b \rangle^1 \Rightarrow \langle \langle K M \rangle^{\min(-b, 0)} \langle N R \rangle^{\max(-b, 0)} \rangle^0 \end{aligned} \end{aligned}$$

## 5.2. Реализация

Будем делать нашу структуру максимально абстрактной, так, чтобы она могла хранить объекты произвольной природы. Главное, чтобы на множестве этих объектов была определена операция сравнения, а также, если эти объекты сложные, то еще потребуются операции копирования и освобождения памяти (конструктор и деструктор).



Итак, вставка узла. Функция `_tree_insert` принимает на вход Структуру дерева с интерфейсными функциями, указатель на узел и указатель на значение, которое должно быть помещено в это дерево. Данная функция возвращает 0, если высота дерева не изменилась и 1, если высота дерева увеличилась. Это поможет изменить баланс вышестоящему узлу. Результат рекуррентного вызова `_tree_insert` будем сохранять в переменную `dbal` модифицировать при помощи нее баланс получившегося узла. В этой функции мы будем использовать функциональную структуру, которая не будет модифицироваться, а будет замешаться на новую, через конструктор `_Tg`

34a *⟨Вставка нового узла 34a⟩*≡ (41b)

```

int _tree_insert(Tree* t, Node** n, void* data) {
    Node* x = *n;
    int dbal;
    if ( !x ) {
        *n = _Tr(NULL, t->cpy(data), NULL, 0);
        return 1;
    }
    if ( t->cmp(data, x->data) <= 0 )
        dbal = -_tree_insert(t, &x->l, data);
    else
        dbal = _tree_insert(t, &x->r, data);

    *n = _Tr(x->l, x->data, x->r, x->bal + dbal);
    dbal = ABS((*n)->bal) > ABS(x->bal) ? 1 : 0;
    free(x);
    return dbal;
}

```

Defines:

`_tree_insert`, used in chunk 34b.

Uses `_Tr` 34c, `ABS` 41b, `Node` 35b 65c, and `Tree` 35c.

Саму эту функцию использовать напрямую не очень удобно, поэтому мы её обернем интерфейсным методом:

34b *⟨Интерфейс 34b⟩*≡ (41b) 36b▷

```

int tree_insert(Tree* t, void* data) {
    return _tree_insert(t, &t->root, data);
}

```

Defines:

`tree_insert`, used in chunks 41a and 42.

Uses `_tree_insert` 34a and `Tree` 35c.

Конструктор дерева выглядит, с одной стороны, просто, поскольку заполняет поля структуры узла `Node`. Но этот же конструктор также должен следить за балансом. Идея такая, что если требуется балансировка, то конструктор будет рекурсивно вызывать себя для перестроения поддеревьев.

34c *⟨Создание нового узла с балансировкой 34c⟩*≡ (41b)

```

Node* _Tr(Node* l, void* data, Node* r, int bal) {
    Node* res = NULL;
    ⟨Какая-то балансирующая магия 35a⟩

    res = malloc(sizeof(Node));
    *res = (Node){l, r, data, bal};
    return res;
}

```

Defines:

`_Tr`, used in chunks 34a, 35a, and 38–40.

Uses `Node` 35b 65c.

Собственно, сама балансировка, по большей части, повторяет структуру функции *B*:

35a  $\langle \text{Какая-то балансирующая магия 35a} \rangle \equiv$  (34c)

```

if ( bal == 2 ) {
    Node R = *r;
    if ( R.bal < 0 ) {
        Node RL = *R.l;
        res = _Tr(_Tr(l, data, RL.l, MIN(-RL.bal, 0)), RL.data,
                  _Tr(RL.r, R.data, R.r, MAX(-RL.bal, 0)), 0);
        free(R.l);
    } else
        res = _Tr(_Tr(l, data, R.l, 0), R.data, R.r, 0);
    free(r);
} else if ( bal == -2 ) {
    Node L = *l;
    if ( L.bal > 0 ) {
        Node LR = *L.r;
        res = _Tr(_Tr(L.l, L.data, LR.l, MIN(-LR.bal, 0)), LR.data,
                  _Tr(LR.r, data, r, MAX(-LR.bal, 0)), 0);
        free(l->r);
    } else
        res = _Tr(L.l, L.data, _Tr(L.r, data, r, 0), 0);
    free(l);
}
if (res) return res;

```

Uses *\_Tr* 34c, *MAX* 41b, *MIN* 41b 80, and *Node* 35b 65c.

Заметьте, что балансировщик удаляет старые варианты узлов и строит новые, которые используют существующие поддеревья. Дадим теперь определения.

35b  $\langle \text{Определение узла 35b} \rangle \equiv$  (41b)

```

typedef struct _Node {
    struct _Node* l;
    struct _Node* r;
    void* data;
    int bal;
} Node;

```

Defines:

*Node*, used in chunks 34–40, 67, 68a, 70, 71, 73, 77, and 78.

Сами интерфейсные функции по работе со значениями в узлах (копирование *cpy*, сравнение *cmp* и освобождение *free*) ложатся в структуру *Ttree*, заодно в ней будем хранить корень дерева.

35c  $\langle \text{Определение интерфейса дерева 35c} \rangle \equiv$  (41a)

```

typedef struct _Tree {
    void* (*cpy)(void* data);
    int (*cmp)(void* x, void* y);
    void (*free)(void* data);
    struct _Node* root;
} Tree;

```

Defines:

*Tree*, used in chunks 34, 36, 37, and 40–42.

### 5.3. Поиск и разрушение

Поиск осуществляется рекурсивно, результат получается через обратный вызов `cb`. Если поиск нужно прекратить `cb` должен вернуть 1.

36a  $\langle \text{Функция поиска 36a} \rangle \equiv$  (41b)

```
int _tree_search(Tree* t, Node* n, void* key, int (*cb)(void*)) {
    int cmp;
    if ( !n ) return 0;
    cmp = t->cmp(key, n->data);
    if ( cmp == 0 )
        if ( cb(n->data) ) return 1;
    if ( cmp <= 0 )
        return _tree_search(t, n->l, key, cb);
    else
        return _tree_search(t, n->r, key, cb);
}
```

Defines:

`_tree_search`, used in chunk 36b.

Uses `Node` 35b 65c and `Tree` 35c.

36b  $\langle \text{Интерфейс 34b} \rangle + \equiv$  (41b)  $\triangleleft 34b \ 36d \triangleright$

```
int tree_search(Tree* t, void* key, int (*cb)(void*)) {
    return _tree_search(t, t->root, key, cb);
}
```

Defines:

`tree_search`, used in chunk 41a.

Uses `_tree_search` 36a and `Tree` 35c.

Также нам будет полезна функция обхода всего дерева (для отладки и печати). Также делаем ее в паре с интерфейсной. Останавливаемся, если `cb` возвращает ненулевое значение. Дополнительно в `cb` передается уровень узла дерева:

36c  $\langle \text{Функция обхода дерева 36c} \rangle \equiv$  (41b)

```
int _tree_walk(Tree* t, Node* n, int level, int(*cb)(void*, int, int)) {
    if ( !n ) return 0;
    if (_tree_walk(t, n->l, level+1, cb)) return 1;
    if ( cb(n->data, level, n->bal) ) return 1;
    return _tree_walk(t, n->r, level+1, cb);
}
```

Defines:

`_tree_walk`, used in chunk 36d.

Uses `Node` 35b 65c and `Tree` 35c.

36d  $\langle \text{Интерфейс 34b} \rangle + \equiv$  (41b)  $\triangleleft 36b \ 37b \triangleright$

```
int tree_walk(Tree* t, int (*cb)(void*, int, int)) {
    return _tree_walk(t, t->root, 0, cb);
}
```

Defines:

`tree_walk`, used in chunks 41a and 42.

Uses `_tree_walk` 36c and `Tree` 35c.

Разрушение дерева (деструктор) строится по все тому же рекуррентному принципу.

37a  $\langle \text{Деструктор 37a} \rangle \equiv$  (41b)

```
void _tree_free(Tree* t, Node* n) {
    Node* l, *r;
    if ( !n ) return;
    t->free(n->data);
    l = n->l;
    r = n->r;
    free(n);
    _tree_free(t, l);
    _tree_free(t, r);
}
```

Defines:

\_tree\_free, used in chunk 37b.

Uses Node 35b 65c and Tree 35c.

37b  $\langle \text{Интерфейс 34b} \rangle + \equiv$  (41b)  $\triangleleft 36d \ 40b \triangleright$

```
void tree_free(Tree* t) {
    _tree_free(t, t->root);
}
```

Defines:

tree\_free, used in chunks 41a and 42.

Uses \_tree\_free 37a and Tree 35c.

## 5.4. Удаление узла

Удаление узла из дерева, тем более, из самобалансирующегося — это, как и любая другая немонотонность — достаточно сложная операция. Опять попробуем построить теоретическую модель. Начнем строить функцию удаления  $\mathcal{D}(T, x)$  с простых и очевидных случаев:

$$\begin{aligned}\mathcal{D}(\emptyset, x) &= \emptyset, 0 \\ \mathcal{D}(\langle \emptyset x \emptyset \rangle, x) &= \emptyset, 1\end{aligned}\tag{5.1}$$

Забегаая вперед, мы, по аналогии с функцией вставки, будем помимо результирующего дерева возвращать также факт изменения высоты дерева, чтобы вызывающая сторона могла бы поправить баланс. В данном случае, высота дерева не меняется. Теперь посмотрим, что будет в случае, если дерево не пусто:

$$\begin{aligned}\mathcal{D}(\langle L y R \rangle^b, x) &= \\ x < y &\Rightarrow \\ &L', d \leftarrow \mathcal{D}(L, x); \\ &\mathcal{B}(\langle L' y R \rangle^{b+d}), |b+d| < |b| \\ x > y &\Rightarrow \\ &R', d \leftarrow \mathcal{D}(R, x); \\ &\mathcal{B}(\langle L y R' \rangle^{b-d}), |b-d| < |b|\end{aligned}\tag{5.2}$$

Случай удаления, когда  $x = y$ , собственно, представляет основную сложность. Что делать с двумя поддеревьями, когда у них нет общего корня. Оказывается, нужно избрать новый корень — это либо самый левый узел правого поддерева  $R$ , либо самый правый узел левого поддерева, у которого может быть только один потомок. Введем вспомогательные функции для выбора нового корня. Обозначим функции отбора самого левого правого узлов соответственно  $\mathcal{E}_l$  и  $\mathcal{E}_r$ . Эти функции будут возвращать отобранный узел, новое дерево и факт изменения высоты дерева (для корректировки баланса вызываемой стороной).

$$\begin{aligned}\mathcal{E}_l(T) &= \\ T = \langle \emptyset x R \rangle &\Rightarrow x, R, 1 \\ T = \langle L x R \rangle^b &\Rightarrow e, L', d \leftarrow \mathcal{E}_l(L); \\ &e, \mathcal{B}(\langle L' x R \rangle^{b+d}), |b+d| < |b| \\ \mathcal{E}_r(T) &= \\ T = \langle L x \emptyset \rangle &\Rightarrow x, L, 1 \\ T = \langle L x R \rangle^b &\Rightarrow e, R', d \leftarrow \mathcal{E}_r(R); \\ &e, \mathcal{B}(\langle L x R' \rangle^{b-d}), |b-d| < |b|\end{aligned}\tag{5.3}$$

Теперь случай функции удаления, когда удаляется корень, будет выглядеть следующим образом:

$$\begin{aligned}\mathcal{D}(\langle L x R \rangle^b, x) &= \\ b > 0 &\Rightarrow e, R', d \leftarrow \mathcal{E}_l(R); \\ &\mathcal{B}(\langle L e R' \rangle^{b-d}), |b-d| < |b| \\ b \leq 0 &\Rightarrow e, L', d \leftarrow \mathcal{E}_r(L); \\ &\mathcal{B}(\langle L' e R \rangle^{b+d}), |b+d| < |b|\end{aligned}\tag{5.4}$$

## 5.5. Реализация удаления

Определим функции отбора самого левого и правого значений  $\mathcal{E}_l$  и  $\mathcal{E}_r$  дерева, как показано в (5.3):

38  $\langle \text{Функция отбора левого значения 38} \rangle \equiv$  (41b)

```
int elect_l(Node** n, void** e) {
    Node* x = *n;
    int dbal = 1;
    if ( !x->l ) {
        *n = x->r;
        *e = x->data;
    }
```

```

    } else {
        dbal = elect_l(&x->l, e);
        *n = _Tr(x->l, x->data, x->r, x->bal+dbal);
        dbal = ABS((*n)->bal) < ABS(x->bal) ? 1: 0;
    }
    free(x);
    return dbal;
}

```

Defines:

`elect_l`, used in chunk 40a.

Uses `_Tr` 34c, `ABS` 41b, and `Node` 35b 65c.

39

*⟨Функция отбора правого значения 39⟩≡*

(41b)

```

int elect_r(Node** n, void** e) {
    Node* x = *n;
    int dbal = 1;
    if ( !x->r ) {
        *n = x->l;
        *e = x->data;
    } else {
        dbal = elect_r(&x->r, e);
        *n = _Tr(x->l, x->data, x->r, x->bal-dbal);
        dbal = ABS((*n)->bal) < ABS(x->bal) ? 1: 0;
    }
    free(x);
    return dbal;
}

```

Defines:

`elect_r`, used in chunk 40a.

Uses `_Tr` 34c, `ABS` 41b, and `Node` 35b 65c.

Теперь сама функция удаления. Она, по сути, является полным отражением (5.1), (5.2) и (5.4) и использует вспомогательные функции выбора нового корня `elect_l` и `elect_r` в случае, если он удаляется.

40a  $\langle \text{Функция удаления 40a} \rangle \equiv$  (41b)

```
int _tree_delete(Tree* t, Node** n, void* data) {
    Node* x = *n;
    int cmp, dbal;
    if ( !x) return 0;
    cmp = t->cmp(data, x->data);

    printf("Comparing %d and %d => %d\n", (int)data, (int)x->data, cmp);
    if ( !x->l && !x->r && !cmp ) {
        printf("Deleted node found, no children\n");
        t->free(x->data);
        free(x);
        *n = NULL;
        return 1;
    }
    if ( cmp < 0 )
        dbal = _tree_delete(t, &x->l, data);
    else if ( cmp > 0 )
        dbal = - _tree_delete(t, &x->r, data);
    else {
        t->free(x->data);
        dbal = x->bal > 0 ? -elect_l(&x->r, &x->data)
                        : elect_r(&x->l, &x->data);
    }
    *n = _Tr(x->l, x->data, x->r, x->bal+dbal);
    dbal = ABS((*n)->bal) < ABS(x->bal) ? 1 : 0;
    free(x);
    return dbal;
}
```

Defines:

`_tree_delete`, used in chunk 40b.

Uses `_Tr` 34c, `ABS` 41b, `elect_l` 38, `elect_r` 39, `Node` 35b 65c, and `Tree` 35c.

40b  $\langle \text{Интерфейс 34b} \rangle + \equiv$  (41b) <37b

```
int tree_delete(Tree* t, void* data) {
    return _tree_delete(t, &t->root, data);
}
```

Defines:

`tree_delete`, used in chunks 41a and 42.

Uses `_tree_delete` 40a and `Tree` 35c.



Соберем прототипы интерфейсных функций в заголовочный файл:

```
41a <tree.h 41a>≡
    #ifndef _TREE_H_
    #define _TREE_H_
    struct _Node;
    <Определение интерфейса дерева 35с>
    int tree_insert(Tree* t, void* data);
    int tree_delete(Tree* t, void* data);
    int tree_walk(Tree* t, int (*cb)(void*, int, int));
    int tree_search(Tree* t, void* key, int (*cb)(void*));
    void tree_free(Tree* t);
    #endif /* _TREE_H_ */
```

Defines:

\_TREE\_H\_, never used.

Uses Tree 35c, tree\_delete 40b, tree\_free 37b, tree\_insert 34b, tree\_search 36b, and tree\_walk 36d.

А их реализации — в библиотечный файл. Заметьте, что мы не раскрыли определение struct \_Node — это внутренности библиотеки.

```
41b <tree.c 41b>≡
    #include<stdio.h>
    #include<stdlib.h>
    #include "tree.h"
    #define ABS(x) ((x)> 0? (x): -(x))
    #define MIN(x, y) ((x) < (y)? (x) : (y))
    #define MAX(x, y) ((x) > (y)? (x) : (y))
    <Определение узла 35b>
    <Создание нового узла с балансировкой 34с>
    <Вставка нового узла 34а>
    <Функция обхода дерева 36с>
    <Функция поиска 36а>
    <Деструктор 37а>
    <Функция отбора левого значения 38>
    <Функция отбора правого значения 39>
    <Функция удаления 40а>
    <Интерфейс 34б>
```

Defines:

ABS, used in chunks 34a and 38–40.

MAX, used in chunk 35a.

MIN, used in chunks 35a and 72b.

Для тестирования сделаем небольшую программу, которая будет вставлять несколько узлов в дерево и печатать его, затем удалять и печатать еще раз:

```
42  <testtree.c 42>≡
    #include<stdio.h>
    #include<stdlib.h>
    #include "tree.h"
    void* int_copy(void* a) {
        return a;
    }
    void int_free(void* a) {
        return;
    }
    int int_cmp(void* a, void* b) {
        return (a < b)? -1 : (a > b)? 1: 0;
    }

    int tree_print(void* data, int level, int bal) {
        int i;
        for ( i = 0; i < level; i++ )
            printf(" ");
        printf("%d (%d)\n", (int)data, bal);
        return 0;
    }

    int main(int argc, char* argv[]) {
        Tree t;
        int k;
        int m;
        int i;
        t.cpy = int_copy;
        t.free = int_free;
        t.cmp = int_cmp;
        t.root = NULL;
        k = atoi(argv[1]);
        m = atoi(argv[2]);
        for (i = 0; i < k; i++ ) {
            tree_insert(&t, (void*)i);
        }
        tree_walk(&t, tree_print);
        printf("Deleting %d\n", m);
        tree_delete(&t, (void*)m);

        printf("-----\n");
        tree_walk(&t, tree_print);

        tree_free(&t);
        return 0;
    }
```

Defines:

```
int_cmp, never used.
int_copy, never used.
int_free, never used.
main, never used.
tree_print, never used.
```

Uses Tree 35c, tree\_delete 40b, tree\_free 37b, tree\_insert 34b, and tree\_walk 36d.

## Глава 6

# Очередь

### 6.1. Реализация очереди на односвязных списках

**Определение.** Очередью (queue, рiре или FIFO) называется абстрактная структура данных, хранящая упорядоченную последовательность (список) элементов, возможно, различной природы и позволяющая удалять элементы только с начала последовательности и добавлять только в конец (либо наоборот). Очередь нам потребуется для реализации других алгоритмов (обход графа в ширину).

Для начала нам потребуется определение элемента списка произвольной природы. Это структура, которая содержит абстрактный указатель `elt` и указатель на следующий элемент списка `next`.

44a  $\langle \text{определение абстрактного списка 44a} \rangle \equiv$  (47 48a)

```
typedef struct _List {
    void* elt;
    struct _List* next;
} List;
```

Defines:

List, used in chunks 13, 14, 26–29, and 44–48.

Основная проблема реализации очереди связана с тем что добавлять (enqueue) и удалять (dequeue) элементы из нее нужно за константное время  $O(1)$ , т.е. не зависящее от размера очереди. Это возможно сделать, если нам известно два указателя — на первый и последний элементы очереди, назовем их **head** и **last**.

Например, функции добавления enqueue и удаления dequeue могут выглядеть следующим образом:

44b  $\langle \text{enqueue реализация 1 44b} \rangle \equiv$

```
int enqueue(List** head, List** last, void*** elt) {
     $\langle \text{создать новый элемент 44c} \rangle$ 
     $\langle \text{обнулить хвост 44d} \rangle$ 
     $\langle \text{вернуть указатель на полезную нагрузку 44e} \rangle$ 
     $\langle \text{добавить элемент в существующий хвост last 44f} \rangle$ 
     $\langle \text{заместить новым элементом старый last 44g} \rangle$ 
     $\langle \text{случай пустой очереди 44h} \rangle$ 
    return 0;
}
```

Defines:

enqueue, used in chunks 45a, 47, and 48a.

Uses List 12 26a 44a.

Очередь сама заботится об аллокации памяти для своей структуры. Но может случиться, что аллокация невозможна — в этом случае выходим с ошибкой.

44c  $\langle \text{создать новый элемент 44c} \rangle \equiv$  (44b 46b)

```
List* new = malloc(sizeof(List));
if ( ! new ) return -1;
```

Uses List 12 26a 44a.

44d  $\langle \text{обнулить хвост 44d} \rangle \equiv$  (44b)

```
new->next = NULL;
```

44e  $\langle \text{вернуть указатель на полезную нагрузку 44e} \rangle \equiv$  (44b 46b)

```
*elt = &new->elt;
```

Данная строчка может показаться странной. Объяснение дано ниже.

44f  $\langle \text{добавить элемент в существующий хвост last 44f} \rangle \equiv$  (44b)

```
if (*last)
    (*last)->next = new;
```

44g  $\langle \text{заместить новым элементом старый last 44g} \rangle \equiv$  (44b)

```
*last = new;
```

44h  $\langle \text{случай пустой очереди 44h} \rangle \equiv$  (44b)

```
if ( ! *head ) *head = new;
```

Эта реализация требует пояснений. Помимо двух двойных указателей в аргументах функции используется тройной(!) указатель на `void`. На самом деле это проявление “абстрактности” данной структуры. Функция `enqueue` не занимается аллокацией памяти для элемента очереди а *возвращает указатель* на ту область памяти, куда вызывающая сторона сможет записать *указатель* на реальные данные.

Чтобы не ошибиться в количестве символов ``*'` можно прибегнуть к мнемоническому правилу: сколько раз в описании аргумента встречается слов *‘возвращает’*, *адрес* или *‘указатель’* столько ``*'` и надо ставить после типа (перед именем). Т.е. в нашем случае — это ``void***'`. Можно это также прочесть так: нам требуется, чтобы функция *вернула* <sup>1</sup> *указатель* <sup>2</sup> на ту область памяти, куда мы смогли бы записать значение `void*`. В результате, опять поучаем тип `void***`. Он же и подсказывает нам использование данной функции:

45a *⟨возможное использование enqueue 45a⟩*≡

```
List* head = NULL;
List* tail = NULL;
void** pdata;
...
if(!enqueue(&head, &tail, &pdata)) {
    *pdata = malloc(...);
}
...
```

Uses `enqueue` 44b 46b and `List` 12 26a 44a.

Таким образом, данная функция *возвращает* нам данные типа `void**`, которые мы будем использовать как *адрес*, по которому мы положим реальные данные.

Это же правило применимо и первым двум аргументам `head` и `last`. Функция `enqueue` *возвращает* нам *указатели* на `List`.

45b *⟨dequeue реализация 1 45b⟩*≡

```
int dequeue(List** head, List** last, void** elt) {
    List* next, *old = *head;
    ⟨обработать случай пустого списка 45c⟩
    ⟨вернуть полезную нагрузку 45d⟩
    ⟨сохранить следующий элемент очереди 45e⟩
    ⟨обработать случай последнего элемента 45f⟩
    ⟨освободить память 45g⟩
    ⟨сместить начало очереди на следующий элемент 46a⟩
    return 0;
}
```

Defines:

`dequeue`, used in chunks 47 and 48a.

Uses `List` 12 26a 44a.

45c *⟨обработать случай пустого списка 45c⟩*≡

```
if ( ! old ) return -1;
```

(45b)

45d *⟨вернуть полезную нагрузку 45d⟩*≡

```
*elt = old->elt;
```

(45b 46e)

45e *⟨сохранить следующий элемент очереди 45e⟩*≡

```
next = old->next;
```

(45b)

45f *⟨обработать случай последнего элемента 45f⟩*≡

```
if ( old == *last )
    *last = NULL;
```

(45b)

45g *⟨освободить память 45g⟩*≡

```
free(old);
```

(45b 46e)

46a  $\langle \text{сместить начало очереди на следующий элемент 46a} \rangle \equiv$  (45b)  
`*head = next;`

Хранить два указателя не совсем удобно, да и не нужно. Можно оставить себе один указатель, “замкнув” список так, чтобы последний элемент указывал не на NULL, а на голову списка, при этом достаточно хранить только `last`, что и будет нашей очередью. Соответственно, наши функции немного видоизменяются:

46b  $\langle \text{enqueue простая реализация 46b} \rangle \equiv$  (47 48b)  
`int enqueue(List** queue, void*** elt) {`  
 $\langle \text{создать новый элемент 44c} \rangle$   
 $\langle \text{вернуть указатель на полезную нагрузку 44e} \rangle$   
 $\langle \text{вставить элемент между последним и первым 46c} \rangle$   
 $\langle \text{сместить указатель очереди на новый элемент 46d} \rangle$   
`return 0;`  
`}`

Defines:

`enqueue`, used in chunks 45a, 47, and 48a.

Uses List 12 26a 44a.

Изменяется только следующие части:

46c  $\langle \text{вставить элемент между последним и первым 46c} \rangle \equiv$  (46b)  
`if ( ! *queue ) {`  
`new->next = new;`  
`} else {`  
`new->next = (*queue)->next;`  
`(*queue)->next = new;`  
`}`

Если очередь состоит из одного элемента — замкнуть его на себя.

46d  $\langle \text{сместить указатель очереди на новый элемент 46d} \rangle \equiv$  (46b)  
`*queue = new;`

Соответствующим образом изменяется функция `dequeue`:

46e  $\langle \text{dequeue простая реализация 46e} \rangle \equiv$  (47 48b)  
`int dequeue(List** queue, void** elt) {`  
`List* old;`  
 $\langle \text{обработать случай пустой очереди 46f} \rangle$   
 $\langle \text{изъять голову из списка 46g} \rangle$   
 $\langle \text{вернуть полезную нагрузку 45d} \rangle$   
 $\langle \text{случай последнего элемента 46h} \rangle$   
 $\langle \text{освободить память 45g} \rangle$   
`return 0;`  
`}`

Defines:

`dequeue`, used in chunks 47 and 48a.

Uses List 12 26a 44a.

46f  $\langle \text{обработать случай пустой очереди 46f} \rangle \equiv$  (46e)  
`if (! *queue ) return -1;`

46g  $\langle \text{изъять голову из списка 46g} \rangle \equiv$  (46e)  
`old = (*queue)->next;`  
`(*queue)->next = old->next;`

46h  $\langle \text{случай последнего элемента 46h} \rangle \equiv$  (46e)  
`if ( old == *queue) *queue = NULL;`

Вот и вся реализация очереди.

В качестве примера использования очереди приведем алгоритм синтеза звука Карплуса-Стронга.

```
47 <karplus.c 47>≡
#include<stdio.h>
#include<stdlib.h>
#include<string.h>
<определение абстрактного списка 44a>
<enqueue простая реализация 46b>
<dequeue простая реализация 46e>
int main() {
    List *queue = NULL;
    void** pdata;
    int v, f = 0, n = 20;
    while(!feof(stdin) && fscanf(stdin, "%d ", &v)) {
        if ( enqueue(&queue, &pdata) )
            break;
        *pdata = (void*)v;
    }

    while(!dequeue(&queue, (void**)&v) && n) {
        printf("%d\n", v);
        if ( v == f ) n--;
        else n = 20;

        f = f/2 + v/2;
        if ( enqueue(&queue, &pdata) )
            break;
        *pdata = (void*)f;
    }

    while (!dequeue(&queue, (void**)&v) ) {
        printf("%d\n", v);
    }
}
```

Defines:

main, never used.

Uses dequeue 45b 46e, enqueue 44b 46b, and List 12 26a 44a.

## 6.2. Задание

Использовать рекурсию.

1. Упростите функцию `enqueue` в предположении, что полезная нагрузка `elt` уже аллоцирована.
2. Разберите реализацию алгоритма синтеза звука удара по струне Карплуса-Стронга. Модифицируйте его для 1) чисел с плавающей точкой, 2) для строк текста так, чтобы буквы из двух строк перемежались между собой (случайным образом) в пропорции 1:1, 3\*) для строк текста с перемежающимися словами.

## 6.3. Реализации очереди в виде библиотеки

Поскольку, в будущем, нам потребуется реализация очереди (для обхода графа в глубину), оформим ее в виде отдельного модуля, который состоит из заголовочного файла `llqueue.h` и файла `llqueue.c`

48a

```
<llqueue.h 48a>≡
#define _LLQUEUE_H_
#define _LLQUEUE_H_
<определение абстрактного списка 44a>
int enqueue(List** q, void*** data);
int dequeue(List** q, void** data);
#endif
```

Defines:

`_LLQUEUE_H_`, never used.

Uses `dequeue` 45b 46e, `enqueue` 44b 46b, and `List` 12 26a 44a.

48b

```
<llqueue.c 48b>≡
#include <stdio.h>
#include <stdlib.h>
#include "llqueue.h"

<enqueue простая реализация 46b>
<dequeue простая реализация 46e>
```



## Глава 7

# Поиск в ширину

Рассмотрим группу задач, которые сводятся к задаче поиска *в ширину* (breadth-first search). Эта группа задач так или иначе связана с нахождением кратчайшего или оптимального пути, минимального числа ходов в какой-либо игре или распространение фронта в неоднородных средах. В любом случае, пространство поиска представляется в виде графа, а решение задачи – это цепочка ребер, отвечающая какому-то заданному условию минимума. Сама природа значений, связанных с ребрами или с узлами не так важна. Важно, чтобы эти значения были неизменны в процессе решения.

Пусть задан граф  $G = (V, E)$ , состоящий из вершин  $V$  и ребер  $E$ . Предположим, что заданы некоторые значения на ребрах  $f(E) \rightarrow X$  из некоторого множества, на котором определен порядок элементов, а также агрегирующий оператор  $X \oplus X \rightarrow X$ , обладающий свойствами дистрибутивности и коммутативности и монотонности на множестве  $X: \forall a, b \in X a \oplus b \geq a$ .

И пусть заданы начальная и конечная вершины  $v_0, v_t \in V$ . Требуется найти такую неповторяющуюся последовательность вершин, последовательно соединенных ребрами:

$$P_{0 \rightarrow t} = \{v_{p_i}, i = 0 \dots k\}, v_{p_0} = v_0, v_{p_k} = v_t : \forall i < k e_{p_i} = (v_{p_i}, v_{p_{i+1}}) \in E,$$

которая дает минимальное значение  $P_{\min} = \operatorname{argmin} Q(P)$ , где  $Q$  вычисляется по формуле:

$$Q(P) = f(v_{p_0}, v_{p_1}) \oplus f(v_{p_1}, v_{p_2}) \oplus \dots \oplus f(v_{p_{k-1}}, v_t), v_{p_i} \in P$$

**Идея.** Основная идея поиска пути, при котором достигается минимум какого-то агрегата (и при этом чтобы вершины пути не повторялись) заключается в сохранении промежуточного значения этого агрегата в вершинах графа, а сами вершины обходить способом, похожим на принцип распространения фронта волны.

Допустим, у нас есть некоторый *фронт*, состоящий из вершин графа  $v_j$ , для которых вычислено значения агрегата  $Q(v_j)$ . Для каждой вершины фронта  $v_j$  выбираются смежные вершины  $v_i$  такие, что на них еще значение агрегата не задано, ( $Q(v_i) = +\infty$ ), либо  $Q(v_i) > Q(v_j) + f(v_j, v_i)$ . Каждая такая вершина добавляется к вершинам фронта, а  $v_j$  – удаляется. Причем первой удаляется вершина, добавленная первой. В результате получается новый фронт, который, по сути, является очередью вершин. Изначально фронт состоит из одной вершины  $v_0$ , и процесс обхода графа и вычисления агрегата продолжается пока фронт не пуст. Результатом процесса будет либо найденное минимальное значение  $Q(v_t)$ , либо отсутствие пути  $v_0, \dots, v_t$ , в этом случае,  $Q(v_t)$  не будет вычислено ( $Q(v_t) = +\infty$ ). Этот алгоритм называется алгоритмом Дейкстры.

### 7.1. Лабиринт

Одной из наиболее распространенных (учебных) задач, которые позволяют продемонстрировать подход поиска в ширину – это задача поиска кратчайшего пути в плоском лабиринте между входом и выходом. Лабиринт – это поле, состоящее из клеток, каждая из которых может быть либо пустой, либо стеной. Путь может состоять только из пустых клеток. Вход и выход, разумеется, также являются пустыми клетками. Размер лабиринта конечен, он окружен (невидимыми) стенами, функция “стоимости” перехода между клетками – единица, а агрегирующий оператор – простое суммирование, так что путь в лабиринте характеризуется просто числом клеток в нем.

Будем читать структуру лабиринта из текстового файла, в котором будет  $m$  строк по  $n$  символов. Символ '.' (точка) будет означать пустую клетку, символ '#' (решетка) — стену, а символы 'i' и 'o' — соответственно, вход и выход. Результатом решения задачи поиска кратчайшего пути будет все то же представление лабиринта в виде матрицы символов, но с обозначенным кратчайшим путем (одним из) при помощи символов 'p' и обозначением количества клеток в нем.

Вот пример задания файла с лабиринтом

```
#####
#...#.....#.....#.....#.....o#
#...#.....#.....#.....#.....#
#.#.#####.....####.###.....#...##...#.....#
#.#.###...#.....#.....###.....#...##...#.....#
#i#.#.....#.....#.....##.....#...##...#.....##
#..##.....###.....#.....####.##.....##.###.....#
#.#.....#.....#.....#.....##.....##
#####
```

Общая структура программы будет выглядеть так:

51a  $\langle \text{labyrinth.c 51a} \rangle \equiv$   
 $\langle \text{Необходимые библиотеки 53b} \rangle$   
 $\langle \text{Определения 53a} \rangle$   
 $\langle \text{Глобальные переменные 51b} \rangle$   
 $\langle \text{Вспомогательные функции 52c} \rangle$   

```
int main(int argc, char* argv[]) {
    char* filename;
    int err = 0;
    if ( argc < 2 ) {
        fprintf(stderr, "usage: labyrinth <file.lab>\n");
        exit(1);
    }
    filename = argv[1];
     $\langle \text{Прочитать файл лабиринта 52a} \rangle$ 
     $\langle \text{Найти кратчайший путь 54g} \rangle$ 
    if(  $\langle \text{Путь найден 55a} \rangle$  ) {
         $\langle \text{Напечатать лабиринт с кратчайшим путем 55c} \rangle$ 
    } else {
         $\langle \text{Решения не существует 56b} \rangle$ 
    }
     $\langle \text{Деинициализация 56c} \rangle$ 
    return 0;
}
```

Defines:

main, never used.

Теперь потихонечку начнем.

Сначала определим формат представления лабиринта. Пусть это будет массив целых чисел Lab, который мы будем использовать одновременно как информацию о геометрии лабиринта (стены и пустые клетки), так и о длине пути от входа.

51b  $\langle \text{Глобальные переменные 51b} \rangle \equiv$  (51a) 51c >  

```
unsigned int* Lab;
```

Defines:

Lab, used in chunks 52c, 53d, and 56c.

Также нам потребуется размер лабиринта, координаты входа и выхода. Пусть это будут целочисленные ширина W и высота H. Вход и выход будут храниться в парах IX, IY и OX, OY соответственно.

51c  $\langle \text{Глобальные переменные 51b} \rangle + \equiv$  (51a) < 51b 52b >  

```
unsigned int W, H;
unsigned int IX, IY;
unsigned int OX, OY;
```

Defines:

H, used in chunks 52c, 53d, and 56a.

IX, used in chunks 52c, 54f, and 56a.

IY, used in chunks 52c, 54f, and 56a.

OX, used in chunks 52c, 55, and 56a.

OY, used in chunks 52c, 55, and 56a.

W, used in chunks 52c, 53d, and 56a.

Прочитаем лабиринт из файла. Если не удастся это сделать – напечатаем ошибку и выйдем. Будем читать файл построчно, длина первой строки будет определять ширину лабиринта, а количество строк – высоту. Каждый раз, когда строка будет зачитана, будем заполнять массив `Lab` в соответствии со следующими правилами: Всякий символ, не равный '.', 'i' или 'o' будет восприниматься как клетка со стеной и иметь значение 0, в противном случае — это пустая клетка, которая будет иметь значение -1 (аналог  $+\infty$  для `unsigned int`). Соответственно, встретившиеся символы 'i' и 'o' дадут нам координаты IX, IY и OX, OY.

52a *⟨Прочитать файл лабиринта 52a⟩*≡ (51a)

```
if ( err = lab_read(filename) ) {
    fprintf(stderr, "Error while reading labyrinth: %d\n", err);
    exit(1);
}
```

Uses `lab_read` 52c.

Сама реализация функции `lab_read` содержит в себе чтение буфера `BUF`, его анализ и заполнение массива `Lab`.

52b *⟨Глобальные переменные 51b⟩*+≡ (51a) <51c 53c>

```
char BUF[4096];
```

Defines:

BUF, used in chunks 52c and 56a.

52c *⟨Вспомогательные функции 52c⟩*≡ (51a) 54e>

```
int lab_read(char* filename) {
    int i, k;
    FILE *f = fopen(filename, "r");
    if ( !f ) return -1;
    W = 0;
    H = 0;
    Lab = NULL;
    while(!feof(f)) {
        if ( !fgets(BUF, sizeof(BUF), f)) break;
        if ( !H ) W = strlen(BUF);
        Lab = realloc(Lab, sizeof(*Lab)*W*(H+1));
        for ( i = 0; i < W && BUF[i]; i++ ) {
            switch(BUF[i]) {
                case 'i': case 'o': case '.':
                    *(Lab + W*H + i) = -1; break;
                default:
                    *(Lab + W*H + i) = 0;
            }
            if (BUF[i] == 'i') { IX = i; IY = H; }
            else if ( BUF[i] == 'o') { OX = i; OY = H; }
        }
        while ( i < W ) *(Lab + W*H + i++) = 0;
        H++;
    }
    fclose(f);
    return 0;
}
```

Defines:

`lab_read`, used in chunk 52a.

Uses BUF 52b, H 51c, IX 51c, IY 51c, Lab 51b, OX 51c, OY 51c, and W 51c.

В результате у нас получится массив `Lab` размера  $W * H$  с элементами 0 и -1.

Теперь перейдем к самому главному — нахождению кратчайшего пути. Для этого нам потребуется очередь. Воспользуемся реализацией из предыдущего раздела. В очереди будем хранить пару координат клетки  $x$  и  $y$ , для которых мы создадим структуру

53a  $\langle \text{Определения 53a} \rangle \equiv$  (51a) 53d

```
typedef struct _cell {
    int x, y;
} Cell;
```

Defines:

`Cell`, used in chunks 53 and 54.

Подключим необходимые заголовочные файлы

53b  $\langle \text{Необходимые библиотеки 53b} \rangle \equiv$  (51a)

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include "llqueue.h"
```

Создадим глобальную переменную с очередью.

53c  $\langle \text{Глобальные переменные 51b} \rangle + \equiv$  (51a) <52b

```
List* Queue;
```

Defines:

`Queue`, used in chunks 53 and 54.

Теперь определим рекурсивную функцию обхода лабиринта с фронтом в очереди `Queue`. Если очередь пуста – прекращаем работу и проверяем, дошли ли мы до выхода. Немного упростим себе задачу, добавив макросы для проверки правильности шага из текущей клетки `IS_VALID`, нахождения внутри границ лабиринта `IN_BOUNDS`, а также доступ к значению в клетке `CELL`:

53d  $\langle \text{Определения 53a} \rangle + \equiv$  (51a) <53a

```
#define IS_VALID(x, y) ((x)*(x) + (y)*(y) > 0 && (x)*(y) == 0)
#define IN_BOUNDS(X, Y) ((X)>= 0 && (X) < W && (Y)>=0 && (Y) < H)
#define CELL(X, Y) *(Lab + W*(Y) + (X))
```

Defines:

`CELL`, used in chunks 54–56.

`IN_BOUNDS`, used in chunks 54c and 55b.

`IS_VALID`, used in chunks 54c and 55b.

Uses `H` 51c, `Lab` 51b, and `W` 51c.

53e  $\langle \text{Функция обхода 53e} \rangle \equiv$  (54e)

```
int walk() {
    Cell* c;
    Cell** front;
    int dx, dy;
    unsigned int p;
    if ( ! Queue ) {
        printf("Front is empty!\n");
        return 0;
    }
     $\langle \text{Читаем текущую клетку из фронта 54a} \rangle$ 
     $\langle \text{Запоминаем число шагов до текущей клетки 54b} \rangle$ 
     $\langle \text{Перебираем возможные варианты путей и формируем новый фронт 54c} \rangle$ 
    return walk();
}
```

Defines:

`walk`, used in chunk 54g.

Uses `Cell` 53a and `Queue` 53c.

54a *⟨Читаем текущую клетку из фронта 54a⟩≡* (53e)  
`dequeue(&Queue, (void**)&c);`

Uses Queue 53c.

54b *⟨Запоминаем число шагов до текущей клетки 54b⟩≡* (53e)  
`p = CELL(c->x, c->y);`

Uses CELL 53d.

При проверке возможного шага воспользуемся тем, что -1 для `unsigned int` — это максимальное число, а стена — это 0. Число шагов `p` у нас строго положительно, поэтому перезаписывать значение в клетке будем только если оно больше чем `p+1`. Т.е. это точно не стена, а клетка, либо еще не пройденная, либо с большим числом шагов (хотя, в нашем случае такая ситуация невозможна):

54c *⟨Перебираем возможные варианты путей и формируем новый фронт 54c⟩≡* (53e)  
`for( dx = -1; dx <= 1; dx++ ) {  
 for ( dy = -1; dy <= 1; dy++ ) {  
 if ( IS_VALID(dx, dy) && IN_BOUNDS(c->x + dx, c->y + dy) ) {  
 if ( CELL(c->x + dx, c->y + dy) > p + 1 ) {  
 CELL(c->x + dx, c->y + dy) = p + 1;  
⟨Добавить новую клетку во фронт 54d⟩  
 }  
 }  
 }  
}`  
`free(c);`

Uses CELL 53d, IN\_BOUNDS 53d, and IS\_VALID 53d.

Далее следует немного Си-шной магии. Это получение адреса области памяти в которую нужно будет записать указатель на нашу структуру `Cell` и заполнение ее новыми значениями. Заметьте, что заполнить структуру, доступную по указателю можно не прибегая к серии присваиваний отдельных ее полей, а сразу приравнять известным значениям. Следим за количеством “звезд”:

54d *⟨Добавить новую клетку во фронт 54d⟩≡* (54c)  
`enqueue(&Queue, (void***)&front);  
*front = malloc(sizeof(Cell));  
**front = (Cell){c->x + dx, c->y + dy};`

Uses Cell 53a and Queue 53c.

Добавляем нашу функцию в список вспомогательных функций.

54e *⟨Вспомогательные функции 52c⟩+≡* (51a) <52c 54f>  
*⟨Функция обхода 53e⟩*

А само нахождение кратчайшего пути выразим в функции инициализации с последующим вызовом функции обхода `walk`:

54f *⟨Вспомогательные функции 52c⟩+≡* (51a) <54e 55b>  
`void init_walk() {  
 Cell** init;  
 CELL(IX, IY) = 1;  
 enqueue(&Queue, (void***)&init);  
*init = malloc(sizeof(Cell));  
**init = (Cell){IX, IY};  
}`

Defines:

`init_walk`, used in chunk 54g.

Uses CELL 53d, Cell 53a, IX 51c, IY 51c, and Queue 53c.

54g *⟨Найти кратчайший путь 54g⟩≡* (51a)  
`init_walk();  
walk();`

Uses `init_walk` 54f and `walk` 53e.

Как понять, что выход достигнут? Значение в этой клетке должно быть “конечным”.

55a  $\langle \text{Путь найден 55a} \rangle \equiv$  (51a)

```
CELL(OX, OY) < (unsigned int)(-1)
```

Uses CELL 53d, OX 51c, and OY 51c.

Теперь не менее сложная часть — это рисование найденного пути. Для его маркировки будем использовать еще одно “близкое к бесконечности” значение `(unsigned int)(-2)`, которое мы врядли достигнем. Строить найденный путь, очевидно, лучше всего с конца:

55b  $\langle \text{Вспомогательные функции 52c} \rangle + \equiv$  (51a)  $\langle 54f \ 56a \rangle$

```
int walk_back(int x, int y, int p) {
    int dx, dy;
    if (p == 1) return 0;
    for (dx = -1; dx <= 1; dx++ ) {
        for (dy = -1; dy <=1; dy++ ) {
            if ( IS_VALID(dx, dy) && IN_BOUNDS(x+dx, y+dy) ) {
                if ( CELL(x+dx, y+dy) == p - 1 ) {
                    CELL(x+dx, y+dy) = (unsigned int)(-2);
                    return walk_back(x+dx, y+dy, p - 1);
                }
            }
        }
    }
    return -1;
}
```

Defines:

`walk_back`, used in chunk 55c.

Uses CELL 53d, IN\_BOUNDS 53d, and IS\_VALID 53d.

55c  $\langle \text{Напечатать лабиринт с кратчайшим путем 55c} \rangle \equiv$  (51a)

```
walk_back(OX, OY, CELL(OX, OY));
print_lab();
```

Uses CELL 53d, OX 51c, OY 51c, print\_lab 56a, and walk\_back 55b.

Функция печати будет выглядеть обратной функции чтения. Читаем построчно массив Lab и заполняем буфер BUF символами, соответствующими либо пустой клетке, либо стене, либо пути, либо точкам входа и выхода.

56a  $\langle \text{Вспомогательные функции 52c} \rangle + \equiv$  (51a)  $\triangleleft 55b$

```
void print_lab() {
    int x, y;
    for ( y = 0; y < H; y++ ) {
        for ( x = 0; x < W; x++ ) {
            switch(CELL(x, y)) {
                case 0:
                    BUF[x] = '#'; break;
                case (unsigned int)(-2):
                    BUF[x] = '.'; break;
                default:
                    if ( x == IX && y == IY )
                        BUF[x] = 'i';
                    if ( x == OX && y == OY )
                        BUF[x] = 'o';
                    else
                        BUF[x] = ' ';
            }
        }
        BUF[W] = '\0';
        puts(BUF);
    }
}
```

Defines:

print\_lab, used in chunk 55c.

Uses BUF 52b, CELL 53d, H 51c, IX 51c, IY 51c, OX 51c, OY 51c, and W 51c.

В случае отсутствия пути, напечатаем сообщение:

56b  $\langle \text{Решения не существует 56b} \rangle + \equiv$  (51a)

```
fprintf(stderr, "Path not found!\n");
```

Не забудем подчистить за собой массив.

56c  $\langle \text{Деинициализация 56c} \rangle + \equiv$  (51a)

```
free(Lab);
```

Uses Lab 51b.

Ну вот, опять все. Кстати, результат решения задачи будет выглядеть как-то так:

## 7.2. Задания

1. Разрешите шагать в ближайшие клетки по диагонали.
2. Разрешите шагать в ближайшие клетки по диагонали, только измените стоимость шага так, чтобы хождения по осям прибавляло к пути 10, а хождение по диагонали – 14.



## Глава 8

# Хеш-таблицы

В задачах хранения данных с быстрым ( $< O(N)$ ) доступом, помимо деревьев поиска можно использовать хеш-таблицы. Они хороши тем, что в некотором диапазоне объемов хранимой информации способны обеспечить доступ к элементам за время близкое к  $O(1)$ .

**Идея.** Известно, доступ к элементам массива  $a_i$  в памяти по индексу (т.е. фактически, по адресу) происходит очень быстро. Но не все множества элементов можно пронумеровать сквозным индексом. Зато можно сопоставить каждому объекту  $x$  некоторое число  $h$ , которое затем использовать в виде индекса массива  $a_{h(x)}$ . Функция сопоставления  $x \rightarrow h$  должна быть такой, чтобы (хоть сколько-нибудь) различающимся объектам ставить в соответствие разные числа. Если эти числа находятся в разумных пределах ( $h \in 0 \dots N$ ) и хорошо распределены (равномерно, без “выбросов”) то мы практически получим аналог массива с объектами в качестве индексов. Такие функции называют хеш-функциями. А способ хранения объектов в массиве, индексируемом хеш-функциями называют *хешированием*.

Естественно, различных объектов может быть потенциально больше чем  $N$ , и есть вероятность (по теореме о голубях и ящиках), что два разных объекта, которые мы собираемся добавить в наш массив, могут иметь равные значения хеш-функции. Как при этом поступить? Значение хеш-функции не идентифицирует объект. В этом случае, говорят, что произошла *коллизия*, т.е. столкновение.

**Разрешение коллизий.** Способов, как поступить с коллизиями несколько. Первый – это допустить хранение нескольких объектов в одной ячейке, соответствующей одному значению хеш-функции, поместив их в односвязный список. Это т.н. метод цепочек.

Второй метод, не требующий динамической аллокации – это метод т.н. открытой адресации, когда второму объекту, претендующему на уже занятое место в массиве, подыскивается другое место. Вариантов поиска местечка тоже несколько. Но все равно, этот метод имеет один недостаток – сложность реализации его слишком высока, а удаление объекта оказывается невозможным.

**Расширение хеш-таблицы.** Понятно, что оба метода работают только когда количество элементов сопоставимо с размером массива. В противном случае, время поиска (или вставки) приближается к линейному. В случае, если используется метод цепочек, естественно, количество элементов может быть больше, но не сильно. Во втором случае, когда массив пытается максимально заселиться объектами, время поиска начинает существенно расти гораздо раньше (при коэффициенте заполнения порядка 0.7). А если хеш-функция подобрана плохо (пожертвовали качеством ради скорости и простоты реализации), то имеет смысл просто хранить объекты в односвязном списке – он будет работать быстрее только из-за того, что в нем меньше накладных расходов.

В обоих вариантах, если лимитами задачи не ограничено количество элементов, подлежащих хранению, необходимо предусмотреть возможность расширения таблицы. Делать это можно экспоненциально, т.е. сразу расти в 2 раза. При этом есть несколько способов избавиться от старой таблицы. Либо сразу “перехешировать” элементы из старой таблицы в новую, либо перехешировать их пачками при каждом удобном случае, вставлять новые объекты только в новую, а искать в обеих таблицах. Либо не перехешировать, а оставить две таблицы для поиска, а когда заполнится вторая – добавить

3-ю и т.д. В этом случае количество таблиц будет пропорционально логарифму числа элементов с хорошим коэффициентом при асимптотике.

## 8.1. Когда надо и когда не надо использовать хеш-таблицы

**Надо,** когда число объектов велико, но ограничено, а сами объекты небольшие. Чем больше объект, тем дольше вычисление хеш-функции. Когда объекты не надо удалять. Когда существует хорошая (равномерно распределяющая) и быстрая хеш-функция. Когда значение хеш-функции можно закешировать или предвычислить. Когда таблица гарантировано помещается в ОЗУ.

**Не надо,** когда таблица небольшая – линейный поиск/вставка в односвязный список будет иметь меньше накладных расходов, Когда неизвестны ограничения по количеству/размеру объектов – лучше использовать бинарные деревья поиска. Когда таблица точно не поместится в ОЗУ – накладных расходов по работе с диском будет существенно больше, так что разница с В-деревом нивелируется. Когда объекты нужно будет часто удалять. Когда будет требоваться обход всех объектов. Когда будет важно сохранение порядка вставки. Когда потребуется хранить множество объектов с одинаковыми ключами.

## 8.2. Реализация

Рассмотрим самый простой случай реализации хеш-таблицы со списками. Для начала, нужно понять, объекты какой природы мы будем хранить в нашей таблице. От этого будет зависеть выбор хеш-функции.

Чаще всего, объектами (ключами) хеш-таблицы являются строки. Если строки небольшие (< 1000 байт), а характерное количество объектов – большое, то можно выбрать одну из простых хеш-функций. Выдумывать ее точно не стоит. Имеет смысл воспользоваться известными, с доказанными характеристиками. Например Rot13:

59a  $\langle \text{Реализация хеша Rot13 59a} \rangle \equiv$  (63)

```

unsigned int rot13(char* c) {
    unsigned int h = 0;
    while (*c) {
        h += (unsigned char)(*c);
        h -= (h>>19) | (h<<13);
        c++;
    }
    return h;
}

```

Defines:

rot13, used in chunk 63.

Опять будем строить абстрактную Хеш-таблицу, которая могла бы хранить объекты произвольной природы. Для этого нам потребуется массив списков и интерфейсные функции.

59b  $\langle \text{Определение списка 59b} \rangle \equiv$  (62b)

```

typedef struct _HashList {
    void* data;
    struct _HashList* next;
} HashList;

```

Defines:

HashList, used in chunks 59–62.

Интерфейсные функции — это функции копирования данных `cpy`, освобождения `free`, сравнения `eq`, и вычисления хеша `hash`. Также структура интерфейса, для простоты, будет хранить в себе указатель на массив списков, образующих таблицу `table` и размер таблицы `sz`.

59c  $\langle \text{Интерфейс хеш-таблицы 59c} \rangle \equiv$  (62c)

```

typedef struct _HashTable {
    void* (*cpy)(const void*);
    void (*free)(void*);
    int (*eq)(const void*, const void*);
    unsigned int (*hash)(const void*);
    struct _HashList** table;
    size_t sz;
} HashTable;

```

Defines:

HashTable, used in chunks 59–63.

Хеш таблица всегда должна быть инициализирована с достаточным начальным размером. Ячейка массива со значением `NULL` будет означать пустую ячейку, поэтому таблицу заполняем нулями.

59d  $\langle \text{Инициализация хеш-таблицы 59d} \rangle \equiv$  (62b)

```

int ht_init(HashTable* t, size_t sz) {
    t->table = calloc(sizeof(HashList*), sz);
    if ( !t->table ) return -1;
    t->sz = sz;
    return 0;
}

```

Defines:

ht\_init, used in chunks 62c and 63.

Uses HashList 59b and HashTable 59c.

Забегаая вперед, сразу определим для себя стратегию поиска места в таблице, которая будет использоваться как для вставки, так для поиска и удаления. Она будет нам возвращать элемент списка. Сначала ищем нужный нам список в таблице, и, если находим, итерируемся по нему. Если он будет пустой, т.е. указывать на NULL, то это будет означать, что объект `data` не найден.

60a *⟨Поиск места 60a⟩* ≡ (62b)

```

  HashList** _ht_lookup(HashTable* t, const void* data) {
      int h;
      HashList** lst;
      ⟨Вычисляем значение хеш-функции 61b⟩
      ⟨Встаем на нужную ячейку h 61c⟩
      ⟨Бежим по списку, сравнивая элементы 61d⟩
  }

```

Defines:

`_ht_lookup`, used in chunks 60 and 61a.

Uses `HashList` 59b and `HashTable` 59c.

Соответственно, функция вставки объекта в таблицу будет искать место для вставки, если объект уже найден, то вставлять ничего не нужно – возвращаем 1. Иначе – создаем элемент списка и помещаем в него копию объекта `data`.

60b *⟨Функция вставки 60b⟩* ≡ (62b)

```

  int ht_insert(HashTable* t, const void* data) {
      HashList** lst;
      lst = _ht_lookup(t, data);
      if ( *lst ) return 1;
      *lst = malloc(sizeof(HashList));
      **lst = (HashList){t->cpy(data), NULL};
      return 0;
  }

```

Defines:

`ht_insert`, used in chunks 62c and 63.

Uses `_ht_lookup` 60a, `HashList` 59b, and `HashTable` 59c.

Функция поиска, будет строиться аналогично. Только в случае, если ничего не найдено – ничего не делаем.

60c *⟨Функция поиска 60c⟩* ≡ (62b)

```

  int ht_search(HashTable* t, const void* key, void** data) {
      HashList** lst = _ht_lookup(t, key);
      if ( !*lst ) return 0;
      *data = (*lst)->data;
      return 1;
  }

```

Defines:

`ht_search`, used in chunk 62c.

Uses `_ht_lookup` 60a, `HashList` 59b, and `HashTable` 59c.

Удаление – это, фактически, удаление из односвязного списка

61a *⟨Функция удаления 61a⟩≡* (62b)

```
int ht_remove(HashTable* t, const void* key) {
    HashList** lst = _ht_lookup(t, key);
    HashList* l = *lst;
    if ( !l) return -1;
    free(*lst);
    t->free(l->data);
    *lst = l->next;
    return 0;
}
```

Defines:

ht\_remove, used in chunk 62c.

Uses \_ht\_lookup 60a, HashList 59b, and HashTable 59c.

Теперь пришло время разобраться с тем, как работает вспомогательная функция \_ht\_lookup.

61b *⟨Вычисляем значение хеш-функции 61b⟩≡* (60a)

```
h = t->hash(data) % t->sz;
```

Здесь мы приводим диапазон значений хеша к размеру таблицы. Далее, смотрим, что у нас в ячейке и индексом h.

61c *⟨Встаем на нужную ячейку h 61c⟩≡* (60a)

```
lst = &t->table[h];
```

Если в этой ячейке что-то ненулевое, значит, это голова списка — пытаемся сравнить объекты. Если они не равны – продолжаем, если равны – прекращаем бег. Если мы дошли до конца списка и объект не найден, возвращаем указатель на следующий элемент списка, пригодный для добавления в хвост.

61d *⟨Бежим по списку, сравнивая элементы 61d⟩≡* (60a)

```
while(*lst) {
    if ( t->eq((*lst)->data, data) ) break;
    lst = &(*lst)->next;
}
return lst;
```

Отдельно нам еще будет нужна функция, которая итерируется по всем объектам, хранящимся в хеш-таблице:

61e *⟨Функция итерации 61e⟩≡* (62b)

```
void ht_iter(HashTable* t, int (*cb)(const void*)) {
    size_t i;
    HashList* lst;
    for (i = 0; i < t->sz; i++) {
        lst = t->table[i];
        while (lst) {
            if ( cb(lst->data) ) return;
            lst = lst->next;
        }
    }
}
```

Defines:

ht\_iter, used in chunks 62c and 63.

Uses HashList 59b and HashTable 59c.

Ну и напоследок, разрушение таблицы.

62a  $\langle$ Функция разрушения 62a $\rangle \equiv$  (62b)

```
void ht_free(HashTable* t) {
    HashList* l, *n;
    size_t i;
    for ( i = 0; i < t->sz; i++ ) {
        l = t->table[i];
        while(l) {
            n = l->next;
            t->free(l->data);
            free(l);
            l = n;
        }
    }
    free(t->table);
}
```

Defines:

ht\_free, used in chunks 62c and 63.

Uses HashList 59b and HashTable 59c.

### 8.3. Библиотека

Соберем все функции в библиотеку:

62b  $\langle$ hash.c 62b $\rangle \equiv$

```
#include <stdio.h>
#include <stdlib.h>
#include <string.h>
#include "hash.h"
 $\langle$ Определение списка 59b $\rangle$ 
 $\langle$ Инициализация хеш-таблицы 59d $\rangle$ 
 $\langle$ Поиск места 60a $\rangle$ 
 $\langle$ Функция вставки 60b $\rangle$ 
 $\langle$ Функция поиска 60c $\rangle$ 
 $\langle$ Функция удаления 61a $\rangle$ 
 $\langle$ Функция итерации 61e $\rangle$ 
 $\langle$ Функция разрушения 62a $\rangle$ 
```

И вынесем их определения в заголовочный файл:

62c  $\langle$ hash.h 62c $\rangle \equiv$

```
#ifndef _HASH_H_
#define _HASH_H_

 $\langle$ Интерфейс хеш-таблицы 59c $\rangle$ 
int ht_init(HashTable* ht, size_t sz);
int ht_insert(HashTable* ht, const void* data);
int ht_search(HashTable* ht, const void* key, void** data);
int ht_remove(HashTable* ht, const void* key);
void ht_iter(HashTable* ht, int (*cb)(const void* data));
void ht_free(HashTable* ht);

#endif /* _HASH_H_ */
```

Defines:

\_HASH\_H\_, never used.

Uses HashTable 59c, ht\_free 62a, ht\_init 59d, ht\_insert 60b, ht\_iter 61e, ht\_remove 61a, and ht\_search 60c.

## 8.4. Тестирование

Теперь немного тестирования. Будем считать количество уникальных строк, поступающих на вход `stdin`, предполагая, что число строк будет не более 20-50 тыс.

```
63 <hashtest.c 63>≡
    #include <stdio.h>
    #include <string.h>
    #include <stdlib.h>
    #include "hash.h"
    <Реализация хеша Rot13 59a>
    int streq(const void* a, const void* b) {
        return !strcmp((const char*)a, (const char*)b);
    }

    int main() {
        int uniqs = 0;
        char buf[4096];
        int count(const void* data) {
            uniqs += 1;
            return 0;
        }

        HashTable tbl = {(void*)(*) (const void*)}strdup,
                        free, streq,
                        (unsigned int)(*)(const void*)}rot13, NULL, 0};

        ht_init(&tbl, 10093);

        while( fgets(buf, sizeof(buf), stdin) ) {
            ht_insert(&tbl, buf);
        }

        ht_iter(&tbl, count);
        ht_free(&tbl);
        printf("Unique strings: %d\n", uniqs);
        return 0;
    }
```

Defines:

`main`, never used.  
`streq`, never used.

Uses HashTable 59c, `ht_free` 62a, `ht_init` 59d, `ht_insert` 60b, `ht_iter` 61e, and `rot13` 59a.

## 8.5. Задания

1. Доработайте тестовый пример, чтобы он выводил уникальные строки
2. Ускорьте работу хеш таблицы – храните в таблице сразу головы списков чтобы избежать лишней аллокации при добавлении в пустую ячейку.
3. Попробуйте различные алгоритмы хеширования и хеш-функции.

## Глава 9

# Элементарная база данных

В этой главе мы попробуем создать небольшую (в плане функциональности) базу данных, основанную на хеш-таблицах и заодно познакомимся с файловым вводом-выводом. База данных будет представлять собой персистентный (хранящийся на диске, переживающий выключение) ассоциативный массив ключ-значение, хранящийся в одном файле, где в качестве ключа и значения могут выступать строки произвольной длины. Приблизительно представляя, что такая база данных будет иметь интерфейс типа CRUD (от англ. названия четырех базовых операций Create-Retrieve-Update-Delete), попробуем, для начала, понять, как хеш-таблицу уложить в файл с возможностью изменения и дополнения.

### 9.1. Файловая грамматика

Даже бинарные данные, при упаковке в файл имеют иерархическую структуру. Эта структура описывается *грамматикой*. С грамматиками и синтаксисом мы уже знакомились в задаче про скобки. Грамматика, описывающая структуру хранения базы данных внутри файла имеет некоторые особенности, которые мы сейчас разберем. Итак, файл базы данных представляет собой заголовок и список таблиц (будем использовать способ расширения таблицы за счет добавления связанных таблиц большего размера).

$DB$	$\leftarrow$	$DHeaderTable$
$DHeader$	$\leftarrow$	$DMagicVersionStat$
$Table$	$\leftarrow$	$THeaderNode +$
$THeader$	$\leftarrow$	$TMagicCapacitySizeNext$
$Node$	$\leftarrow$	$KeyOffKeySzValueOffValueSzNextEl$

Начнем по порядку строить структуру файла нашей БД.

$$DB \leftarrow DHeaderTable$$

Файл базы данных  $DB$  начинается с заголовка  $DHeader$ , за которым сразу следует таблица  $Table$ .

$$DHeader \leftarrow DMagicVersionStat$$

Заголовок  $DHeader$  будет содержать в себе некоторую "магическую" последовательность байт  $DMagic$ , которая будет указывать на тип файла, информацию о версии формата файла  $Version$  – пригодится нам в будущем, если формат будет меняться, а также статистику нашей БД  $Stat$  для быстрой информации о том, сколько ключей хранится в базе, сколько таблиц создано, средний размер ключа, средний размер значений и т.д.

Заголовку соответствует структура



65a  $\langle \text{Структура заголовка БД 65a} \rangle \equiv$  (80)

```
typedef struct _DHeader {
    char magic[4];
    uint32_t version;
    struct _Stat stat;
} DHeader;
```

Defines:

DHeader, used in chunks 71a, 77, and 78d.

Подробное внутреннее устройство статистики *Stat* можно определить позже.

$$Table \leftarrow THeaderNode +$$

Таблица *Table* состоит из заголовка *THeader* и массива элементов *Node*. Количество элементов, следующих за заголовком определяется в самом заголовке:

$$THeader \leftarrow TMagicCapacitySizeNext$$

Заголовок таблицы *THeader* содержит в себе также ”магическую” последовательность байт *TMagic*, которая будет говорить о том, что мы читаем именно нашу таблицу, Указание на общий размер таблицы *Capacity*, заполнение *Size* и указатель на следующую таблицу.

65b  $\langle \text{Структура заголовка таблицы 65b} \rangle \equiv$  (80)

```
typedef struct _THeader {
    char magic[4];
    uint32_t capacity;
    uint32_t size;
    int64_t next;
} THeader;
```

Defines:

THeader, used in chunks 67, 70b, 72a, 73b, and 78c.

$$Node \leftarrow KeyOffKeySzValueOffValueSzNext$$

Элемент содержит в себе информацию о расположении блока данных с ключом *KeyOff* – смещение внутри файла, *KeySz* – размер ключа и блока данных со значением (аналогично *ValueOff*, *ValueSz* – смещение и размер соответственно). Также, элемент содержит в себе ссылку на следующий элемент *Next* (смещение внутри файла) на случай коллизии.

65c  $\langle \text{Структура элемента таблицы 65c} \rangle \equiv$  (80)

```
typedef struct _Node {
    int64_t keyoff;
    uint64_t keysz;
    int64_t valueoff;
    uint64_t valuesz;
    int64_t next;
} Node;
```

Defines:

Node, used in chunks 34–40, 67, 68a, 70, 71, 73, 77, and 78.

Поскольку некоторые элементы могут быть незаполненными (пустыми), мы будем считать таковыми элементы с  $KeyOff = 0$ .

## 9.2. Описание работы

Работа нашей базы данных будет практически повторять работу обычной хеш-таблицы в памяти с тем исключением, данные будут храниться на диске, а следовательно, нам потребуется свой аналог функций `malloc` и `free`, которые бы управляли выделением пространства для блоков данных в файле.

Проще всего выделять память в конце файла. Обычно, все файловые системы оптимизированы по скорости добавления данных в конец файла. Поэтому каждый новый блок можно дописывать в конец файла, сохраняя его начальное смещение (которое будет равно размеру файла). Вот такой простой `malloc`. А `free` может вообще ничего не делать. Старые значения и ключи, которые больше не используются проще оставлять внутри файла. Так можно себе позволить поступать при перезаписи значений и ключей, которые могут иметь произвольные размеры, и перезаписать новое значение на старое место уже невозможно без повреждения соседних данных, или удаление элемента из односвязного списка – старый элемент придется оставить на месте – переиспользовать его уже мы вряд ли сможем. А вот перезаписывать блоки фиксированной длины вполне возможно – это касается заголовков БД, Таблиц, а также самих элементов. Да, конечно, файл при таком расточительном подходе на каждое изменение будет только расти, однако скорость записи в таком случае будет максимальной, поскольку мы не будем тратить время на то, чтобы подыскать свободное местечко для очередного блока внутри файла. Ведь существенную часть файла мы собираемся оставить на месте, а не перезаписывать.

На случай, когда файл окажется слишком большим для реально хранимых значений (например, шла долгая история с перезаписями по одним и тем же ключам и файл будет копить в себе всю историю изменений, которая уже не может быть использована) лучше предусмотреть специальную процедуру "вакуумизации" – перезаписи файла базы данных как-бы заново, оставляя только используемые блоки.

Итак, рассмотрим добавление новой записи ключ-значение в нашу БД. Изначально, предполагаем, что файл БД содержит в себе только заголовок и таблицу с  $Capacity = 10$  (нижние индексы обозначают "версию" блоков)

$$DB_0 = DHeader_0(DMagic, Version, Stat_0), THeader_0(TMagic, 10, 0, 0), Node_0^0(0, 0, 0, 0, 0), \dots, Node_0^9(0, 0, 0, 0, 0)$$

Добавление пары ключ-значение должно повлиять на статистику  $Stat$ ,  $Size$  в  $THeader$  и какой-то элемент, допустим  $i$ . Новые ключ-значение дописываются в конец файла:

$$DB_1 = DHeader_1(DMagic, Verison, Stat_1), THeader_1(TMagic, 10, 1, 0), Node_0^0(0, 0, 0, 0, 0), \dots, Node_1^i([K_0], < X >)$$

Здесь  $[X]$  означает смещение блока  $X$ , а  $< X >$  – его размер. Последовательность модификации файла такова, что сначала дописываются в конец блоки с ключом и значением, их смещения записываются и записываются в элемент  $Node^i$ , Обновляется заголовок таблицы  $THeader$  и заголовок базы, поскольку изменяется статистика  $Stat$ .

Если мы теперь захотим изменить значение по данному ключу  $K_0$  на  $V_1$ , то файл перезапишется следующим образом:

$$DB_2 = DHeader_2(DMagic, Version, Stat_2), THeader_1(TMagic, 10, 1, 0), Node_0^0(0, 0, 0, 0, 0), \dots, Node_2^i([K_0], < V_1 >)$$

Здесь происходит допись в конец только нового значения  $V_1$  и обновление элемента  $Node^i$  и всех заголовков.

Рассмотрим теперь вариант поведения, когда происходит коллизия ключей при добавлении пары  $K_2, V_2$ . Коллизия ключей  $K_0$  и  $K_2$  на значении хеша  $i$  должна приводить к добавлению элемента в односвязный список, место для нового элемента  $Node^{i'}$  также выделяется в конце файла:

$$DB_3 = DHeader_3(DMagic, Version, Stat_3), THeader_2(Magic, 10, 2, 0), Node_0^0(0, 0, 0, 0, 0), \dots, Node_3^i([K_0], < V_1 >), Node_3^{i'}([K_2], < V_2 >)$$

Как видим, тут тоже важна последовательность записи в файл, чтобы не делать несколько перезаписей. Сначала дописываются блоки с парой ключ-значение  $K_2, V_2$ , затем дописывается следующий элемент цепочки  $Node^{i'}$ , который ссылается на них, затем процесс повторяется для элемента таблицы и всех заголовков.

Введем такое понятие как *курсор* – это такая структура данных, которая сохраняет состояние поиска нужного элемента. Поскольку мы будем иметь несколько связанных таблиц, которые будем обходить последовательно, имеет смысл сохранять прочитанный заголовок текущей таблицы *ht*, а заодно и его смещение внутри файла *tableoff*, чтобы, в случае чего, его обновить. Когда будем искать ключ в самой таблице, запомним также прочитанную структуру элемента *node* таблицы, его индекс *idx* в таблице и смещение *nodeoff*, опять же, на случай модификации. Не будем также забывать, что потенциально элементы могут образовывать список и текущий элемент списка также было бы не плохо иметь под рукой, поскольку его придется модифицировать. Назовем эти поля *chain* и *chainoff* соответственно.

67  $\langle \text{Курсор } 67 \rangle \equiv$  (80)

```
typedef struct _Cursor {
    FILE* fh;
    Stat* stat;
    THeader th;
    off_t tableoff;
    uint64_t hash;
    int idx;
    Node node;
    off_t nodeoff;
    Node chain;
    off_t chainoff;
    Node prev;
    off_t prevoff;
}Cursor;
// node    chain
// v        v
// N1->N3->N4->N5
// N2
// N5->N6->N7->N8 <<< cursor
// x
```

Defines:

Cursor, used in chunks 68–72, 74, 75a, 77, and 78.  
Uses Node 35b 65c, Stat 79c, and THeader 65b.

Теперь рассмотрим алгоритм "дописывания" в конец цепочки элементов нового элемента и пары ключ-значение. Будем предполагать, что мы уже встали на последний элемент цепочки и это состояние будет передано структуре `Cursor`. Выделим некоторые "низкоуровневые" операции с файлов в виде функций с префиксом `_file_`. Тогда алгоритм дописывания будет выглядеть следующим образом:

68a *⟨Запись в цепочку 68a⟩* ≡ (77c)

```
int _cur_write_chain(Cursor* cur, const char* key, const char* v) {
    Node new;
    _file_append_block(cur->fh, key, &new.keyoff, &new.keysz);
    _file_append_block(cur->fh, v, &new.valueoff, &new.valuesz);
    _file_append_node(cur->fh, &new, &cur->chain.next);
    _cur_update_chain(cur);
    cur->th.size ++;
    _cur_update_table(cur);
    cur->stat->keys += new.keysz;
    cur->stat->valuesz += new.valuesz;
    cur->stat->keys ++;
    _cur_update_stat(cur);
    return 0;
}
```

Defines:

`_cur_write_chain`, used in chunk 74a.

Uses `_cur_update_chain` 78b, `_cur_update_stat` 78d, `_cur_update_table` 78c, `_file_append_block` 72d, `_file_append_node` 73a, `Cursor` 67, and `Node` 35b 65c.

Здесь мы используем предполагаемые низкоуровневые операции добавления блока и элемента в конец файла (`_file_append_block` и `_file_append_node` соответственно). А также высокоуровневые операции с курсором, которые позволяют на месте вносить изменения в файл. Функция `_cur_update_chain` обновляет элемент `chain` курсора после обновления поля `next`. Функция `_cur_update_table` обновляет заголовок текущей таблицы после внесения туда изменений, связанных со статистикой. Функция `_cur_update_stat` обновляет структуру со статистикой всей БД.

Алгоритм записи в таблицу (безконфликтной) выглядит похожим образом, но обновляется табличный элемент `node`, а не элемент цепочки `chain`:

68b *⟨Запись в таблицу 68b⟩* ≡ (77c)

```
int _cur_write_node(Cursor* cur, const char* key, const char* v) {
    _file_append_block(cur->fh, key, &cur->node.keyoff, &cur->node.keysz);
    _file_append_block(cur->fh, v, &cur->node.valueoff, &cur->node.valuesz);
    _cur_update_node(cur);
    cur->th.size ++;
    _cur_update_table(cur);
    cur->stat->keys += cur->node.keysz;
    cur->stat->valuesz += cur->node.valuesz;
    cur->stat->keys ++;
    _cur_update_stat(cur);
    return 0;
}
```

Defines:

`_cur_write_node`, used in chunk 74a.

Uses `_cur_update_node` 77d, `_cur_update_stat` 78d, `_cur_update_table` 78c, `_file_append_block` 72d, and `Cursor` 67.

Поскольку мы имеем дело с мутабельным ассоциативным массивом, значение *v*, которое ассоциировано с ключом может быть изменено независимо от ключа. Опять же, предполагаем, что курсор стоит на нужном элементе. В этом случае мы только добавляем в конец файла блок с новым значением *v*, обновляем элемент, чтобы он указывал на новый блок и статистику:

69a *⟨Обновление элемента 69a⟩*≡ (77c)

```
int _cur_set_value(Cursor* cur, const char* value) {
    size_t oldsz = cur->chain.valuesz;
    _file_append_block(cur->fh, value, &cur->chain.valueoff, &cur->chain.valuesz);
    _cur_update_chain(cur);
    cur->stat->valuesz += cur->chain.valuesz - oldsz;
    _cur_update_stat(cur);
}
```

Defines:

`_cur_set_value`, used in chunk 74a.

Uses `_cur_update_chain` 78b, `_cur_update_stat` 78d, `_file_append_block` 72d, and `Cursor` 67.

Удаление из цепочки и удаление элемента из самой таблицы выглядят по-разному, это их можно объединить в одной функции. Если курсор находится на элементе внутри таблицы, у него смещения `nodeoff` и `chainoff` должны совпадать. Физически удалить или забыть элемент в этом случае мы не можем, поэтому просто пометим его как удаленный при помощи указания нулевого смещения для значения `valueoff`. В противном случае, это просто удаление из односвязного списка. Аналогично, при этом, модифицируем статистику текущей таблицы и статистику БД.

69b *⟨Удаление элемента 69b⟩*≡ (77c)

```
int _cur_del_node(Cursor* cur) {
    if ( cur->nodeoff != cur->chainoff ) {
        cur->prev.next = cur->chain.next;
        _cur_update_prev(cur);
    } else {
        cur->node.valueoff = 0;
        _cur_update_node(cur);
    }
    cur->th.size --;
    _cur_update_table(cur);
    cur->stat->valuesz -= cur->node.valuesz;
    cur->stat->keysz -= cur->node.keysz;
    _cur_update_stat(cur);
}
```

Defines:

`_cur_del_node`, used in chunk 75a.

Uses `_cur_update_node` 77d, `_cur_update_prev` 78a, `_cur_update_stat` 78d, `_cur_update_table` 78c, and `Cursor` 67.

До этого мы предполагали, что курсор встал на нужную нам позицию для записи. Теперь посмотрим, как может происходить поиск элемента. Забегая вперед, можно сказать, что поиск может иметь две различных цели. Первая – это чтение – мы ищем элемент по ключу, чтобы прочитать его или удостовериться, что он есть. Если ключ не найден – ничего в этом случае не делаем. Вторая цель – это модификация, если элемент не найден, мы пытаемся найти место, чтобы записать новый. При этом вполне может быть создана новая таблица, если предыдущая достигла критического значения фактора заполнения ( $\approx 0.7$ ). Цель поиска можно задать в самом курсоре, для этого предусмотрено поле `mod`.

Для начала опишем проход по цепочке:

70a *⟨Проход по цепочке 70a⟩* ≡ (77c)

```
int _cur_find(Cursor* cur, const char* key) {
    if ( cur->chain.valueoff && _cur_cmp_chain(cur, key) ) return 1;
    if ( cur->chain.next ) {
        printf("Reading next node at %ld\n", cur->chain.next);
        _cur_read_chain(cur, cur->chain.next);
        return _cur_find(cur, key);
    }
    return 0;
}
```

Defines:

`_cur_find`, used in chunk 70b.

Uses `_cur_cmp_chain` 71b, `_cur_read_chain` 71d, and `Cursor` 67.

Здесь мы пропускаем элементы с нулевым `valueoff`. Функция `_cur_cmp_chain` сравнивает побайтно ключ, а функция `_cur_load_chain` загружает поле `chain` курсора следующим элементом цепочки.

Собственно, сам поиск элемента также является рекурсивной функцией обхода таблиц и перемещения на позицию с индексом, соответствующим хеш-функции.

70b *⟨Поиск элемента 70b⟩* ≡ (77c)

```
int _cur_search(Cursor* cur, const char* key) {
    cur->idx = cur->hash % cur->th.capacity;
    printf("Searching for key %s at table %ld idx %u\n", key, cur->tableoff, cur->idx);
    cur->nodeoff = cur->tableoff + sizeof(THeader) + sizeof(Node)*cur->idx;
    _cur_load_node(cur);
    cur->chainoff = cur->nodeoff;
    cur->chain = cur->node;
    if(_cur_find(cur, key))
        return 1;
    if ( cur->th.next ) {
        _cur_read_table(cur, cur->th.next);
        return _cur_search(cur, key);
    }
    return 0;
}
```

Defines:

`_cur_search`, used in chunks 71a and 74a.

Uses `_cur_find` 70a, `_cur_load_node` 71c, `_cur_read_table` 72a, `Cursor` 67, `Node` 35b 65c, and `THeader` 65b.

Чтобы начать поиск, сначала вычислим значение хеш-функции, загрузим первую таблицу, инициализируем поля курсора и стартуем поиск.

71a *⟨Инициализация курсора 71a⟩*≡ (80)

```
int _ht_search(DB* db, Cursor* cur, const char* key) {
    memset(cur, 0, sizeof(Cursor));
    cur->stat = &db->stat;
    cur->fh = db->fh;
    cur->hash = db->hash(key);
    _cur_read_table(cur, (off_t)sizeof(DHeader));
    return _cur_search(cur, key);
}
```

Defines:

\_ht\_search, used in chunks 74 and 75a.

Uses \_cur\_read\_table 72a, \_cur\_search 70b, Cursor 67, and DHeader 65a.

71b *⟨Сравнение ключа 71b⟩*≡ (77c)

```
int _cur_cmp_chain(Cursor* cur, const char* key) {
    size_t ks = strlen(key) + 1;
    printf("Comparing key '%s' with cursor\n", key);
    if ( cur->chain.keysz != ks ) return 0;
    printf("Sizes %d are equal\n", ks);
    fseek(cur->fh, cur->chain.keyoff, SEEK_SET);
    return _file_cmp_block(cur->fh, key, ks);
}
```

Defines:

\_cur\_cmp\_chain, used in chunk 70a.

Uses \_file\_cmp\_block 72b and Cursor 67.

71c *⟨Загрузка элемента таблицы 71c⟩*≡ (77c)

```
int _cur_load_node(Cursor* cur) {
    fseek(cur->fh, cur->nodeoff, SEEK_SET);
    if (fread(&cur->node, sizeof(Node), 1, cur->fh) != 1 ) {
        error("Cannot read node at %ld\n", cur->nodeoff);
        return 1;
    }
    printf("Loaded node at %d: keyoff=%ld, keysz=%lu, valueoff=%ld, valuesz=%lu\n",
        cur->idx, cur->node.keyoff, cur->node.keysz,
        cur->node.valueoff, cur->node.valuesz);
}
```

Defines:

\_cur\_load\_node, used in chunk 70b.

Uses Cursor 67, error 80, and Node 35b 65c.

71d *⟨Чтение элемента цепочки 71d⟩*≡ (77c)

```
int _cur_read_chain(Cursor* cur, off_t offset) {
    cur->prev = cur->chain;
    cur->prevoff = cur->chainoff;
    cur->chainoff = offset;
    fseek(cur->fh, cur->chainoff, SEEK_SET);
    return fread(&cur->chain, sizeof(Node), 1, cur->fh) != 1;
}
```

Defines:

\_cur\_read\_chain, used in chunk 70a.

Uses Cursor 67 and Node 35b 65c.

- 72a  $\langle \text{Чтение таблицы 72a} \rangle \equiv$  (77c)
- ```

int _cur_read_table(Cursor* cur, off_t offset) {
    cur->tableoff = offset;
    fseek(cur->fh, cur->tableoff, SEEK_SET);
    return fread(&cur->th, sizeof(THeader), 1, cur->fh) != 1;
}

```
- Defines:  
 \_cur\_read\_table, used in chunks 70b, 71a, and 74a.  
 Uses Cursor 67 and THeader 65b.
- 72b  $\langle \text{Низкоуровневое побайтовое сравнение 72b} \rangle \equiv$  (72c)
- ```

int _file_cmp_block(FILE* fh, const char* key, size_t ks) {
    char buf[4096];
    size_t bytes;
    if ( ks == 0 ) {
        printf("Keys are equal\n");
        return 1;
    }
    bytes = MIN(4096, ks);
    printf("Reading %lu bytes\n", bytes);
    bytes = fread(buf, 1, MIN(4096, ks), fh);
    printf("Comparing %lu bytes of key blocks '%.4s' and '%.4s'...\n", bytes, buf, key);
    if ( memcmp(buf, key, MIN(bytes, ks)) == 0 ) {
        printf("Chunks are equal\n");
        return _file_cmp_block(fh, key+ks, ks-bytes);
    }
    printf("Chunks are not equal\n");
    return 0;
}

```
- Defines:  
 \_file\_cmp\_block, used in chunk 71b.  
 Uses MIN 41b 80.
- Соберем все низкоуровневые операции, которым нужен только файл:
- 72c  $\langle \text{Низкоуровневые операции 72c} \rangle \equiv$  (80) 76d▷
- $\langle \text{Низкоуровневое побайтовое сравнение 72b} \rangle$   
 $\langle \text{Добавление блока 72d} \rangle$   
 $\langle \text{Добавление элемента 73a} \rangle$   
 $\langle \text{Добавление таблицы 73b} \rangle$
- Доопределим недостающие функции для добавления в конец новых записей.
- 72d  $\langle \text{Добавление блока 72d} \rangle \equiv$  (72c)
- ```

int _file_append_block(FILE* fh, const char* value, off_t* offset, size_t* size) {
    *size = strlen(value) + 1;
    fseek(fh, 0, SEEK_END);
    *offset = ftell(fh);
    return fwrite(value, *size, 1, fh) != 1;
}

```
- Defines:  
 \_file\_append\_block, used in chunks 68 and 69a.



73a *⟨Добавление элемента 73a⟩*≡ (72c)

```
int _file_append_node(FILE* fh, Node* node, off_t* offset) {
    fseek(fh, 0, SEEK_END);
    *offset = ftell(fh);
    return fwrite((void*)node, sizeof(Node), 1, fh) != 1;
}
```

Defines:

`_file_append_node`, used in chunk 68a.

Uses Node 35b 65c.

73b *⟨Добавление таблицы 73b⟩*≡ (72c)

```
int _file_append_table(FILE* fh, size_t capacity, off_t* offset) {
    THeader th = (THeader){TABLE_MAGIC, capacity, 0, 0};
    Node empty = (Node){0, 0, 0, 0, 0};
    int i;
    fseek(fh, 0, SEEK_END);
    *offset = ftell(fh);
    if (fwrite(&th, sizeof(th), 1, fh) != 1) return -1;
    for (i = 0; i < capacity; i++) {
        if (fwrite(&empty, sizeof(Node), 1, fh) != 1) return -2;
    }
    return 0;
}
```

Defines:

`_file_append_table`, used in chunks 74a and 75b.

Uses Node 35b 65c, TABLE\_MAGIC 78e, and THeader 65b.

Пришло время определить высокоуровневые операции, которые, собственно, и образуют интерфейсные функции:

73c *⟨Интерфейсные функции 73c⟩*≡ (80)

*⟨Открытие файла базы данных 75b⟩*  
*⟨Закрытие файла базы данных 76c⟩*  
*⟨Установка ассоциации 74a⟩*  
*⟨Получение значения по ключу 74b⟩*  
*⟨Удаление значения по ключу 75a⟩*  
*⟨Получение статистики 76a⟩*

Самая интересная операция – это операция помещения пары ключ-значение в базу. Сделаем это с учетом расширения таблиц.

74a *⟨Установка ассоциации 74a⟩* ≡ (73c)

```
int ht_set(DB* db, const char* key, const char* value) {
    Cursor cur;
    if ( _ht_search(db, &cur, key) )
        // элемент с нужным ключом найден -- обновляем значение
        return _cur_set_value(&cur, value);
    else {
        // элемент не найден
        if ( cur.th.size*10 >= cur.th.capacity*7 ) {
            // создаем новую таблицу, если текущая таблица переполнена
            _file_append_table(cur.fh, cur.th.capacity*2, &cur.th.next);
            cur.stat->tables ++;
            cur.stat->capacity += cur.th.capacity*2;
            _cur_update_stat(&cur);
            _cur_update_table(&cur);
            _cur_read_table(&cur, cur.th.next);
            _cur_search(&cur, key);
        }
        if ( cur.nodeoff == cur.chainoff && !cur.node.keyoff )
            return _cur_write_node(&cur, key, value);
        else
            return _cur_write_chain(&cur, key, value);
    }
}
```

Defines:

ht\_set, used in chunks 79a and 81.

Uses \_cur\_read\_table 72a, \_cur\_search 70b, \_cur\_set\_value 69a, \_cur\_update\_stat 78d, \_cur\_update\_table 78c, \_cur\_write\_chain 68a, \_cur\_write\_node 68b, \_file\_append\_table 73b, \_ht\_search 71a, and Cursor 67.

74b *⟨Получение значения по ключу 74b⟩* ≡ (73c)

```
int ht_get(DB* db, const char* key, char** value) {
    Cursor cur;
    if ( _ht_search(db, &cur, key) ) {
        *value = malloc(cur.chain.valuesz);
        fseek(cur.fh, cur.chain.valueoff, SEEK_SET);
        if ( fread(*value, 1, cur.chain.valuesz, cur.fh) == cur.chain.valuesz ) {
            return 1;
        } else {
            error("Cannod read %lu bytes for value at %ld\n", cur.chain.valuesz,
                cur.chain.valueoff);
            return 0;
        }
    }
    return 0;
}
```

Defines:

ht\_get, used in chunks 79a and 81.

Uses \_ht\_search 71a, Cursor 67, and error 80.

75a *⟨Удаление значения по ключу 75a⟩*≡ (73c)

```

int ht_del(DB* db, const char* key) {
    Cursor cur;
    if ( _ht_search(db, &cur, key) ) {
        _cur_del_node(&cur);
        return 1;
    }
    return 0;
}

```

Defines:

ht\_del, used in chunks 79a and 81.

Uses \_cur\_del\_node 69b, \_ht\_search 71a, and Cursor 67.

Открытие Базы данных – достаточно сложная операция, поскольку сначала нужно удостовериться, если файл существует, то он - наш. Для этого мы сначала его открываем в режиме чтения и вызываем функцию `_file_check_magic`, после проверки – переоткрываем файл на запись, а если файла не существует – создаем файл и записываем заголовок `_file_write_header`, и сразу после заголовка добавляем таблицу. Как только файл успешно открыт, зачитываем из него статистику.

75b *⟨Открытие файла базы данных 75b⟩*≡ (73c)

```

DB* ht_open(const char* filename, size_t initial_capacity) {
    DB* dbh;
    FILE* f = fopen(filename, "r");
    if ( f ) {
        if ( _file_check_magic(f) ) {
            f = freopen(NULL, "r+", f);
        } else {
            error("Wrong magic file %s\n", filename);
            return NULL;
        }
    } else {
        off_t first_table_offset;
        f = fopen(filename, "w+");
        if ( !f ) {
            error("Cannot open file %s for writing\n", filename);
            return NULL;
        }
        _file_write_header(f, initial_capacity);
        _file_append_table(f, initial_capacity, &first_table_offset);
    }
    if ( !f ) {
        error("Cannot open file %s\n", filename);
        return NULL;
    }
    dbh = malloc(sizeof(DB));
    dbh->fh = f;
    dbh->hash = rot1333;
    _file_load_stat(f, &dbh->stat);
    return dbh;
}

```

Defines:

ht\_open, used in chunk 81.

Uses \_file\_append\_table 73b, \_file\_check\_magic 76e, \_file\_load\_stat 77b, \_file\_write\_header 77a, and error 80.

76a *⟨Получение статистики 76a⟩*≡ (73c)

```
int ht_get_stat(DB* dbh, Stat* stat) {
    *stat = dbh->stat;
    return 0;
}
```

Defines:

ht\_get\_stat, used in chunks 79a and 81.

Uses Stat 79c.

Соответственно, мы предполагаем, что в самом объекте БД DB достаточно хранить указатель на открытый файл, функцию хеширования и статистику:

76b *⟨Определение структуры БД 76b⟩*≡ (80)

```
struct _DB {
    FILE* fh;
    uint64_t (*hash)(const char*);
    Stat stat;
};
```

Uses Stat 79c.

76c *⟨Закрытие файла базы данных 76c⟩*≡ (73c)

```
int ht_close(DB* db) {
    if( db ) {
        if (db->fh) fclose(db->fh);
        free(db);
    }
    return 0;
}
```

Defines:

ht\_close, used in chunks 79a and 81.

Доопределим наши низкоуровневые функции

76d *⟨Низкоуровневые операции 72c⟩*+≡ (80) <72c

*⟨Проверка файла 76e⟩*  
*⟨Добавление заголовка 77a⟩*  
*⟨Чтение статистики 77b⟩*

76e *⟨Проверка файла 76e⟩*≡ (76d)

```
int _file_check_magic(FILE* fh) {
    char magic[] = DB_MAGIC;
    char buf[4];
    fseek(fh, 0, SEEK_SET);
    if (fread(buf, 4, 1, fh) != 1 ) {
        error("Cannot read magic\n");
        return 0;
    }
    if ( memcmp(magic, buf, 4) ) {
        error("Wrong magic\n");
        return 0;
    }
    return 1;
}
```

Defines:

\_file\_check\_magic, used in chunk 75b.

Uses DB\_MAGIC 78e and error 80.

77a  $\langle$ Добавление заголовка 77a $\rangle \equiv$  (76d)

```

int _file_write_header(FILE* fh, size_t initial_capacity) {
    Stat stat = {0};
    stat.capacity = initial_capacity;
    stat.tables = 1;
    DHeader dh = (DHeader){DB_MAGIC, CURRENT_VERSION, stat};
    fseek(fh, 0, SEEK_SET);
    if ( fwrite(&dh, sizeof(dh), 1, fh) != 1 ) return -1;
    return 0;
}

```

Defines:

`_file_write_header`, used in chunk 75b.

Uses `CURRENT_VERSION` 78e, `DB_MAGIC` 78e, `DHeader` 65a, and `Stat` 79c.

77b  $\langle$ Чтение статистики 77b $\rangle \equiv$  (76d)

```

int _file_load_stat(FILE* fh, Stat* stat) {
    fseek(fh, offsetof(DHeader, stat), SEEK_SET);
    return fread(stat, sizeof(Stat), 1, fh) != 1;
}

```

Defines:

`_file_load_stat`, used in chunk 75b.

Uses `DHeader` 65a and `Stat` 79c.

А также объединим все операции с курсором.

77c  $\langle$ Операции с курсором 77c $\rangle \equiv$  (80)

$\langle$ Загрузка элемента таблицы 71c $\rangle$   
 $\langle$ Чтение элемента цепочки 71d $\rangle$   
 $\langle$ Чтение таблицы 72a $\rangle$   
 $\langle$ Запись элемента таблицы 77d $\rangle$   
 $\langle$ Запись элемента цепочки 78b $\rangle$   
 $\langle$ Запись предыдущего элемента 78a $\rangle$   
 $\langle$ Запись таблицы 78c $\rangle$   
 $\langle$ Запись статистики 78d $\rangle$   
 $\langle$ Сравнение ключа 71b $\rangle$   
 $\langle$ Проход по цепочке 70a $\rangle$   
 $\langle$ Поиск элемента 70b $\rangle$   
 $\langle$ Запись в цепочку 68a $\rangle$   
 $\langle$ Запись в таблицу 68b $\rangle$   
 $\langle$ Обновление элемента 69a $\rangle$   
 $\langle$ Удаление элемента 69b $\rangle$

Определим последний набор операций с курсором, которые записывают (обновляют) соответствующие поля курсора.

77d  $\langle$ Запись элемента таблицы 77d $\rangle \equiv$  (77c)

```

int _cur_update_node(Cursor* cur) {
    fseek(cur->fh, cur->nodeoff, SEEK_SET);
    return fwrite(&cur->node, sizeof(Node), 1, cur->fh) != 1;
}

```

Defines:

`_cur_update_node`, used in chunks 68b and 69b.

Uses `Cursor` 67 and `Node` 35b 65c.

78a  $\langle \text{Запись предыдущего элемента 78a} \rangle \equiv$  (77c)

```
int _cur_update_prev(Cursor* cur) {
    fseek(cur->fh, cur->prevoff, SEEK_SET);
    return fwrite(&cur->prev, sizeof(Node), 1, cur->fh) != 1;
}
```

Defines:

\_cur\_update\_prev, used in chunk 69b.

Uses Cursor 67 and Node 35b 65c.

78b  $\langle \text{Запись элемента цепочки 78b} \rangle \equiv$  (77c)

```
int _cur_update_chain(Cursor* cur) {
    fseek(cur->fh, cur->chainoff, SEEK_SET);
    return fwrite(&cur->chain, sizeof(Node), 1, cur->fh) != 1;
}
```

Defines:

\_cur\_update\_chain, used in chunks 68a and 69a.

Uses Cursor 67 and Node 35b 65c.

78c  $\langle \text{Запись таблицы 78c} \rangle \equiv$  (77c)

```
int _cur_update_table(Cursor* cur) {
    fseek(cur->fh, cur->tableoff, SEEK_SET);
    return fwrite(&cur->th, sizeof(THeader), 1, cur->fh) != 1;
}
```

Defines:

\_cur\_update\_table, used in chunks 68, 69b, and 74a.

Uses Cursor 67 and THeader 65b.

78d  $\langle \text{Запись статистики 78d} \rangle \equiv$  (77c)

```
int _cur_update_stat(Cursor* cur) {
    fseek(cur->fh, offsetof(DHeader, stat), SEEK_SET);
    return fwrite(cur->stat, sizeof(Stat), 1, cur->fh) != 1;
}
```

Defines:

\_cur\_update\_stat, used in chunks 68, 69, and 74a.

Uses Cursor 67, DHeader 65a, and Stat 79c.

Все магические последовательности отнесем в Константы

78e  $\langle \text{Константы 78e} \rangle \equiv$  (80)

```
#define DB_MAGIC {'H', 'T', 'd', 'b'}
#define TABLE_MAGIC {'H', 'T', 'T', 'b'}
#define CURRENT_VERSION 1
```

Defines:

CURRENT\_VERSION, used in chunk 77a.

DB\_MAGIC, used in chunks 76e and 77a.

TABLE\_MAGIC, used in chunk 73b.

Заголовочный файл с интерфейсом

```
79a  <hashdb.h 79a>≡
    #ifndef _HASHDB_H_
    #define _HASHDB_H_
    #include <stdint.h>

    typedef struct _DB DB;

    <Определение структуры статистики 79c>

    DB* ht_open(const char* filename, size_t initial_capacity);

    int ht_set(DB* dbh, const char* key, const char* value);
    int ht_get(DB* dbh, const char* key, char** value);
    int ht_del(DB* dbh, const char* key);

    int ht_close(DB* dbh);

    int ht_get_stat(DB* dbh, Stat* stat);

    #endif /* _HASHDB_H_ */
```

Defines:

\_HASHDB\_H\_, never used.

ht\_open, used in chunk 81.

Uses ht\_close 76c, ht\_del 75a, ht\_get 74b, ht\_get\_stat 76a, ht\_set 74a, and Stat 79c.

Добавим определение хеш-функции. Возьмем пока простую функцию типа rot13, но расширенную до 64 бит

```
79b  <Хеш-функция 79b>≡ (80)
    uint64_t rot1333(const char* v) {
        uint64_t h = 0x87f3e9c3b1738bf1;
        while(*v) {
            h += 31*(*v);
            h = (h >> 33 ) | (h << 31);
            h = (h >> 35 ) |(h << 29);
            v++;
        }
        return h;
    }
```

```
79c  <Определение структуры статистики 79c>≡ (79a)
    typedef struct _Stat {
        uint64_t keys;
        uint64_t tables;
        uint64_t capacity;
        uint64_t keysz;
        uint64_t valuesz;
    } Stat;
```

Defines:

Stat, used in chunks 67, 76–79, and 81.

### 9.3. Реализация

```

80  <hashdb.c 80>≡
    #include <stdio.h>
    #include <stdlib.h>
    #include <stdint.h>
    #include <stddef.h>
    #include <string.h>
    #include "hashdb.h"

    #define error(...) fprintf(stderr, __VA_ARGS__)
    #define MIN(a, b) ((a)<(b)? (a) : (b))

```

<Константы 78e>

<Определение структуры БД 76b>

<Структура заголовка БД 65a>

<Структура заголовка таблицы 65b>

<Структура элемента таблицы 65c>

<Хеш-функция 79b>

<Курсор 67>

<Низкоуровневые операции 72c>

<Операции с курсором 77c>

<Инициализация курсора 71a>

<Интерфейсные функции 73c>

Defines:

error, used in chunks 71c and 74–76.

MIN, used in chunks 35a and 72b.



Компиляция реализации выглядит просто:

```
gcc -c -O3 -o hashdb.o hashdb.c
```

При этом мы получим библиотеку `hashdb.o`, которую затем можно будет включать в свои модули.

## 9.4. Тестирование

Напишем простую программу для тестирования нашей базы

```
81 <hashdbtest.c 81>≡
#include <stdio.h>
#include <string.h>
#include <stdlib.h>
#include "hashdb.h"

int main(int argc, char* argv[]) {
    DB* dbh;
    char buf[4096];
    if (argc < 2) {
        printf("usage: %s <dbfilename>\n", argv[0]);
        exit(1);
    }
    dbh = ht_open(argv[1], 1024);
    if ( !dbh ) {
        printf("Cannot open database %s\n", argv[1]);
        exit(1);
    }
    while(fgets(buf, 4096, stdin)) {
        char*s = buf;
        char* cmd = strsep(&s, " \t\n");
        if ( cmd ) {
            if (!strcmp(cmd, "SET")) {
                char* key = strsep(&s, " \t\n");
                char* val = strsep(&s, "\n");
                if (key && val) {
                    printf("Setting '%s' = '%s'\n", key, val);
                    ht_set(dbh, key, val);
                    printf("OK\n");
                } else
                    printf("No key or val given\n");
            } else if ( !strcmp(cmd, "GET")) {
                char* key = strsep(&s, " \t\n");
                if ( key ) {
                    char* val;
                    if ( ht_get(dbh, key, &val) ) {
                        printf("%s\n", val);
                        free(val);
                    } else {
                        printf("No such key\n");
                    }
                }
            } else if ( !strcmp(cmd, "DEL")) {
                char* key = strsep(&s, " \t\n");
                if ( key ) {
```

```

        ht_del(dbh, key);
        printf("OK\n");
    }
} else if ( !strcmp(cmd, "STAT")) {
    Stat stat;
    ht_get_stat(dbh, &stat);
    printf("Keys: %lu\n tables: %lu\nAvg. key size: %f\n"
        "Avg. value size: %f\n"
        "Total capacity: %lu\n"
        "Avg. Fill factor: %f\n",
        stat.keys, stat.tables,
        stat.keysz/(float)stat.keys,
        stat.valuesz/(float)stat.keys,
        stat.capacity,
        stat.keys/(float)stat.capacity);
} else if ( !strcmp(cmd, "QUIT") || !strcmp(cmd, "EXIT")) {
    break;
} else {
    printf("unknown command\n");
}
}
}
ht_close(dbh);
return 0;
}

```

Defines:

`main`, never used.

Uses `ht_close` 76c, `ht_del` 75a, `ht_get` 74b, `ht_get_stat` 76a, `ht_open` 75b 79a, `ht_set` 74a, and `Stat` 79c.

## Задание

- Добавить в статистику максимальную длину цепочки.
- Протестировать таблицу на скорость добавления нового элемента. Нарисовать график зависимости времени вставки от количества хранимых элементов. Все делать с отключенным отладочным выводом.
- При большом количестве ключей посчитать среднее отношение добавочной информации к полезной (Отношение размера таблицы к сумме размеров всех ключей и значений). Нарисовать график от количества элементов и их среднего размера.
- проделать все предыдущие задачи с различными хеш-функциями. Выбрать лучшую и доложить.