

Строки

М. А. Ложников

13 февраля 2020 г.

Ревизия: 1

Это предварительная невыверенная версия. Читайте на свой страх и риск.

1 Символы в языке C

В языке C помимо типов данных для целых чисел и чисел с плавающей точкой есть тип данных для хранения символов. Этот тип данных называется `char`. Не смотря на то, что этот тип данных предназначен для хранения отдельных символов, он является обычным целочисленным типом данных, то есть ему можно присваивать целые числа, конвертировать его в целые числа, а также выполнять с переменными этого типа арифметические операции.

Тип данных `char` занимает в памяти 1 байт [2] (то есть `sizeof(char) = 1`) и обычно (но не всегда) принимает целые значения от -128 до 127. Дело в том, что на некоторых архитектурах размер одного байта может отличаться от 8 бит. В языке C количество бит в байте хранится в специальной константе `CHAR_BIT` [3], определённой в заголовочном файле `<limits.h>`.

Таким образом, для хранения символов используются обычные целые числа. Соответствие между числами и символами определяется кодировкой. Как правило числа, которые соответствуют латинским символам, знакам препинания, а также некоторым служебным символам (например, перевод строки) не зависят от кодировки. Это выполняется для символов ASCII. Однако, символы других алфавитов (не ASCII) в различных кодировках могут иметь разные значения или вообще отсутствовать.

2 Символьные константы

Для инициализации символьных типов данных можно использовать специальные символьные константы. Символьная константа записывается в одинарных кавычках. Это может быть либо символ алфавита, либо знак препинания, либо служебный символ. Например,

```
1 char sym = 'a'; // Инициализируем переменную sym строчной латинской буквой a.
2
3 sym = 'D'; // Заглавная латинская буква D.
4 sym = ' '; // Пробел.
```

Может показаться, что компилятор языка C конвертирует символьную константу в числовое значение, которое ей соответствует. На самом деле, всё немного не так. В действительности имеет место обратный процесс. Компилятор всего-навсего использует то значение, при помощи которого символьная константа была представлена в исходном коде программы. То, как эта константа выглядит при редактировании исходного кода в текстовом редакторе определяется кодировкой, выбранной

в текстовом редакторе. А вид этой символьной константы в выводе программы на экран определяется кодировкой, выбранной в терминале.

Таким образом, при использовании символов, отличных от ASCII кодировка файлов исходного кода программы должна совпадать с кодировкой терминала, в котором эта программа запускается.

Помимо символов алфавита 'a', 'b', 'c', ..., 'z' и 'A', 'B', 'C', ..., 'Z' доступны знаки препинания ',', '.', ':', ';', '!', '?', '%', '\$', '#' и т. д., а также различные специальные символы, среди которых важно упомянуть следующие:

- '\n' — символ перевода строки;
- '\t' — символ табуляции;
- '\r' — символ возврата каретки;
- '\"' — символ двойной кавычки;
- '\'' — символ одинарной кавычки;
- '\\ ' — символ \;
- '\0' — символ конца строки, значение которого в любой кодировке равно 0.

Таким образом, специальные символы записываются при помощи символа экранирования \. Отметим, что строковый литерал служит для представления ровно одного символа.

Замечание 2.1 Не смотря на то, что для представления символов в языке C есть специальный тип *char*, типом данных символьной константы в языке C является тип данных *int*.

3 Ввод/вывод символов

3.1 Первый способ

Для ввода/вывода отдельных символов можно использовать функции `scanf()`/`printf()`, указав в формате "%c".

```
1 char sym;
2
3 /* Пример считывания символа с клавиатуры. */
4 if (scanf("%c", &sym) != 1) {
5     printf("Can't read symbol!\n");
6     return 0;
7 }
8
9 /* Пример вывода символа на экран. */
10 if (sym == 'A')
11     printf("sym = A\n");
12 else
13     printf("sym = %c\n", sym);
```

Отметим, что символы считываются последовательно, разделители считываются точно так же, как и все остальные символы и записываются в выходные переменные. Иными словами, никакие символы не игнорируются.

Аналогично, можно использовать функции `fscanf()`/`fprintf()` при работе с файлами.

3.2 Второй способ

Функции `scanf()/printf()` слишком “тяжёлые” для считывания/вывода одного символа. Есть функции полегче, предназначенные только для этого, а именно `getchar()/putchar()` и `fgetc()/fputc()`.

```
1 int sym;
2
3 /* Пример считывания символа с клавиатуры. В случае ошибки или конца файла
4    возвращается константа EOF, определённая в заголовочном файле <stdio.h>. */
5 if ((sym = getchar()) == EOF) {
6     /* Функция perror() выводит сообщение в поток вывода ошибок stderr. */
7     perror("Can't read symbol!\n");
8     return 0;
9 }
10
11 printf("Symbol: ");
12
13 /* Вывод символа на экран. Функция возвращает EOF в случае ошибки. */
14 if (putchar(sym) == EOF) {
15     perror("Can't write the symbol!\n");
16 }
17
18 putchar('\n');
```

То же самое, только с файлами:

```
1 FILE* fin;
2 FILE* fout;
3 int sym;
4
5 if ((fin = fopen("input.txt", "rt")) == NULL) {
6     perror("Can't open input file!\n");
7     return 0;
8 }
9
10 /* Пример считывания символа из файла. В случае ошибки возвращается константа EOF. */
11 if ((sym = fgetc(fin)) == EOF) {
12     perror("Can't read symbol!\n");
13     fclose(fin);
14     return 0;
15 }
16
17 if ((fout = fopen("output.txt", "wt")) == NULL) {
18     perror("Can't open output file!\n");
19     fclose(fin);
20     return 0;
21 }
22
23 /* Пример записи символа в файл. */
24 if (fputc(sym, fout) == EOF) {
```

```
25 perror("Can't write the symbol!\n");
26 }
27
28 putchar('\n');
29
30 fclose(fin);
31 fclose(fout);
```

Кроме того, есть функция `int ungetc(int ch, FILE* stream)`, которая кладёт символ `ch` обратно в файл/поток `stream` и возвращает сам символ в случае успеха и `EOF` в случае ошибки. Эта функция может быть полезна для того, чтобы положить обратно в поток по ошибке считанный символ для того, чтобы считать его в дальнейшем из этого потока. Гарантируется, что функция отработает без ошибки хотя бы один раз, однако несколько последовательных вызовов `ungetc()`, не прерываемых операциями чтения из потока, могут не сработать.

Замечание 3.1 При использовании функций `getchar()` и `fgetc()` важно сохранять значение в тип данных `int`, а не `char`, поскольку множество возвращаемых значений функции больше множества возможных значений типа `char`. Функция возвращает значение типа `unsigned char` (беззнаковый `char`), приведённое к типу `int` или константу `EOF`, значение которой равно `-1`. Значения типов `char` и `unsigned char` переводятся друг в друга без потерь, однако константа `EOF` может совпасть со значением считанного символа при конвертации возвращаемого значения в тип `char/unsigned char`.

Замечание 3.2 В заголовочном файле `<stdio.h>` определены стандартные потоки ввода/вывода и поток вывода ошибок, а именно `stdin`, `stdout` и `stderr` соответственно, которые имеют тип данных `FILE*`. Таким образом, вызов `fgetc(stdin)` эквивалентен вызову `getchar()`, а вызов `fputc(sym, stdout)` эквивалентен вызову `putchar(sym)`.

4 Относительный порядок символов

Во всех кодировках символы алфавита расположены последовательно в алфавитном порядке, то есть им соответствуют последовательные целые числа. Символы, обозначающие цифры также расположены последовательно от `'0'` до `'9'`. Кроме того, символы пунктуации тоже расположены последовательно.

```
1 int i;
2
3 /* Выводим строчные латинские символы от 'a' до 'j'. */
4 for (i = 0; i < 10; i++)
5     putchar('a' + i);
6
7 /* Выводим заглавные латинские символы от 'D' до 'H'. */
8 for (i = 0; i < 5; i++)
9     putchar('D' + i);
```

Тот факт, что символы алфавита расположены по порядку, можно использовать для проверки того, что символ является символом алфавита, а также, например, для конвертации заглавных символов в строчные.

```
1 /* Функция проверяет, является ли символ цифрой. */
2 int IsDigit(int sym) {
3     if (sym >= '0' && sym <= '9')
4         return 1;
5
6     return 0;
7 }
8
9 /* Функция проверяет, является ли символ символом латинского алфавита. */
10 int IsAlpha(int sym) {
11     if ((sym >= 'a' && sym <= 'z') || (sym >= 'A' && sym <= 'Z'))
12         return 1;
13
14     return 0;
15 }
16
17 /* Функция конвертирует заглавные латинские буквы в строчные, а остальные оставляет
18     без изменения. */
19 int ToLower(int sym) {
20     if (sym >= 'A' && sym <= 'Z')
21         return sym + 'a' - 'A';
22
23     return sym;
24 }
```

В заголовочном файле `<ctype.h>` определены функции для работы с отдельными символами. Например, функции `int isalnum(int c)`, `int isalpha(int c)`, `int isdigit(int c)`, `int ispunct(int c)` и `int isspace(int c)` возвращают число, отличное от нуля, если символ `c` является алфавитным или цифрой; алфавитным; цифрой; символом пунктуации или пробельным символом соответственно.

Замечание 4.1 Функции стандартной библиотеки умеют работать не только с кодировкой ASCII. См. функцию `setlocale()` [4].

5 Строки и строковые литералы в языке C

Строки в языке C представляются в виде массивов элементов типа `char`, причём каждая строка должна оканчиваться символом конца строки `'\0'`. Таким образом, строка, состоящая из N символов должна храниться в массиве длины не менее $N + 1$. Длину строки хранить необязательно.

Строковые литералы (константы) в языке C записываются внутри двойных кавычек и имеют тип `const char*` — массив указателей на `const char`, ключевое слово `const` означает, что элементы этого массива нельзя изменять. Каждый строковый литерал неявно заканчивается символом конца строки. Статический массив элементов типа `char` можно инициализировать при объявлении строковым литералом.

```
1 /* Инициализируем при объявлении статический массив str строковым литералом. Каждый
2     строковый литерал заканчивается символом конца строки '\0', и в массиве str
3     он тоже пропихнется. */
4 char str[256] = "Hello world!";
5
```

```
6 /* А вот следующий код уже вызовет ошибку компиляции. Нельзя преобразовать указатель
7    в статический массив, то есть тип const char* в тип char[]. */
8 /* str = "Something else."; */
9
10 /* Объявляем указатель на const char. */
11 const char* ptr;
12
13 /* Здесь всё в порядке, ошибки нет. */
14 ptr = "Some string.";
```

Замечание 5.1 При работе со строками в языке C важно отличать символ конца строки `'\0'` от символа перевода строки `'\n'`. Символ конца строки предназначен для того, чтобы функции работы со строками могли определить количество символов, которое нужно обработать. Символ перевода строки влияет на способ вывода строк, например, функции вывода строки на экран выводят символы с новой строки как только встретят символ перевода строки. Таким образом, одна строка в языке C может содержать несколько символов перевода строки. Отметим, что в английском языке строка в языке C и строка в текстовом редакторе или в выводе на экране называются разными словами: *string* и *line* соответственно.

Замечание 5.2 Обычно текстовые редакторы в Windows вставляют два символа `"\r\n"` для обозначения перевода строки, в то время как текстовые редакторы в Linux обычно вставляют только символ `'\n'`.

5.1 Пример перевода строки в число

```
1 /* Функция для перевода строки в целое число. */
2 int Atoi(const char* str) {
3     int number = 0;
4     int sign = 1;
5
6     /* Если строка пустая. */
7     if (*str == 0)
8         return 0;
9
10    /* Если первый символ --- знак "минус". */
11    if (*str == '-') {
12        sign = -1;
13
14        /* Перемещаем указатель на следующий символ. */
15        str++;
16    }
17
18    /* Цикл до тех пор пока указатель указывает на символ, обозначающий цифру. */
19    for (; IsDigit(*str); str++) {
20        number *= 10;
21
22        /* Символы, обозначающие цифры, расположены по порядку. Значение выражения
23           *str - '0' равно текущей цифре. */
24        number += *str - '0';
25    }
```

```
26
27 return number * sign;
28 }
```

В заголовочном файле `<stdlib.h>` объявлена функция `int atoi(const char* str)`, которая переводит строку в целое число, а также функция `double atof(const char* str)`, которая переводит строку в число с плавающей точкой.

6 Ввод/вывод строк

Как и в случае с отдельными символами строки можно считывать и выводить различными способами.

6.1 Первый способ

Для ввода/вывода строк можно использовать функции `scanf()/printf()` с форматом `"%s"`, а также их аналоги `fscanf()/fprintf()` для чтения/записи из файла/в файл. Стоит отметить, что функции `scanf()/fscanf()` считывают строку до первого разделителя.

```
1 char buffer[256]; /* Заведём массив достаточного размера. */
2
3 /* В функции scanf() перед buffer не нужно ставить &, поскольку массивы передаются
4  в функции по указателям. */
5 if (scanf("%s", buffer) != 1) {
6     printf("Can't read string!\n");
7     return 0;
8 }
9
10 printf("String: \"%s\"\n", buffer);
```

Замечание 6.1 Отметим, что в `printf()` не следует передавать строку напрямую вместо формата, поскольку строка может содержать символы форматирования `%`.

6.2 Второй способ

Функция `char* fgets(char* str, int count, FILE* stream)` считывает из потока/файла `stream` символы в массив `str` до символа перевода строки `'\n'`, но не более `count - 1` символов. Если функция `fgets()` дошла до символа перевода строки, то она запишет его в массив. Кроме того, она всегда записывает символ конца строки `'\0'`. Функция возвращает указатель `str` в случае успеха и `NULL` в случае ошибки. В массиве `str` должна быть выделена память под не менее чем `count` элементов типа `char`.

Функция `int fputs(const char* str, FILE* stream)` выводит строку `str` в файл/поток `stream`. В случае успеха она возвращает неотрицательное значение, а в случае ошибки константу `EOF`.

```
1 char buffer[256];
2
3 /* Считываем строки с клавиатуры до тех пор, пока считываются. */
4 while (fgets(buffer, 256, stdin) != NULL) {
5     printf("New line: %s", buffer);
6 }
```

6.3 Считывание чисел из строк и форматирование строк

Точно так же как и в случае с форматированным вводом/выводом в файл/из файла в стандартной библиотеке языка C существуют функции для форматирования строк, а также для считывания форматированных данных из строки. Функции называются `sscanf()/sprintf()`. Они работают точно так же, как и `fscanf()/fprintf()`, но первым параметром у них является строка.

```
1 char buffer[256] = "3.14 2.71";
2 double a, b;
3
4 if (sscanf(buffer, "%lf%lf", &a, &b) != 2) {
5     /* Этого не должно произойти. */
6     printf("Can't read numbers!\n");
7     return 0;
8 }
9
10 /* Функция sprintf() затрёт данные, которые были в строке buffer до вызова функции. */
11 sprintf(buffer, "Sum is equal to %f", a + b);
12
13 printf("%s\n", buffer);
```

6.4 Считывание чисел из строк

Следующие функции, определённые в `<stdlib.h>`, позволяют считать число из строки, а также узнать, на какой позиции считывание остановилось.

```
1 char firstString[256] = "12345abcde";
2 long firstValue;
3 double secondValue;
4 char* pos;
5
6 /* Первый параметр --- строка, откуда считываются данные, второй параметр ---
7    указатель на указатель на первый символ, который не будет считан, то есть первый
8    символ, не входящий в число. Третий параметр --- основание системы исчисления,
9    в которой записано число. Если вторым аргументом передать NULL, то аргумент будет
10   проигнорирован функцией. */
11 firstValue = strtol(firstString, &pos, 10);
12
13 /* Выведет Value: 12345, position: abcde */
14 printf("Value: %ld, position: %s\n", firstValue, pos);
15
16 secondValue = strtod("1.2345e-3abcd", &pos);
17
18 /* Выведет Value: 0.001234, position: abcd
19    Числа с плавающей точкой по-умолчанию выводятся в фиксированном формате
20    с точностью 6 знаков после запятой. */
21 printf("Value: %f, position: %s\n", secondValue, pos);
```

7 Другие функции работы со строками

Все функции данного раздела объявлены в заголовочном файле `<string.h>`.

7.1 Реализация некоторых функций

Рассмотрим реализации некоторых функций, аналогичные которым есть в стандартной библиотеке языка C.

```
1  /* Функция возвращает длину строки. */
2  size_t StrLen(const char* str) {
3      size_t length = 0;
4
5      for (; *str != 0; str++)
6          length++;
7
8      return length;
9  }
10
11 /* Копирует строку src на место строки dest. Данные, которые были в строке dest
12 затираются. */
13 char* StrCpy(char* dest, const char* src) {
14     char* ptrToDest = dest;
15
16     for (; *src != 0; src++) {
17         *dest = *src;
18         dest++;
19     }
20
21     /* Прописываем символ конца строки. */
22     *dest = 0;
23
24     return ptrToDest;
25 }
26
27 /* Добавляет строку src в конец строки dest. */
28 char* StrCat(char* dest, const char* src) {
29     char* ptrToDest = dest;
30
31     /* Пропускаем символы, которые записаны в dest. У этого цикла нет тела, поскольку
32 после закрывающейся круглой скобки стоит точка с запятой. */
33     for (; *dest != 0; dest++);
34
35     /* Добавляем в конец dest символы из src. */
36     for (; *src != 0; src++) {
37         *dest = *src;
38         dest++;
39     }
40
41     /* Прописываем символ конца строки. */
```

```
42 *dest = 0;
43
44 return ptrToDest;
45 }
```

В стандартной библиотеке языка C есть функции, аналогичные рассмотренным — `strlen()`, `strcpy()`, `strcat()`, а также функция `strcmp()`, предназначенная для сравнения двух строк в лексикографическом порядке.

Замечание 7.1 Отметим, что функции стандартной библиотеки `strlen()`, `strcpy()`, `strcat()`, `strcmp()` могут привести к неопределённому поведению в случае, если в их аргументах не прописан символ конца строки. Для того, чтобы предостеречь себя от подобных ошибок рекомендуется использовать функции `strlen()`, `strncpy()`, `strncat()`, `strncmp()`, которые принимают ещё один дополнительный параметр, ограничивающий сверху число обрабатываемых символов.

7.2 Поиск подстроки в строке

Рассмотрим пример использования некоторых других функций.

```
1 const char* str = "Hello world!";
2 char* ptr;
3
4 /* Возвращает указатель на первое вхождение подстроки (второй параметр)
5    в строку (первый параметр) или NULL, если подстрока в строке не найдена. */
6 ptr = strstr(str, "wor");
7
8 if (ptr != NULL) {
9     printf("Found substring: %s\n", ptr);
10 }
```

7.3 Сортировка строк по возрастанию

Рассмотрим пример простейшей сортировки строк методом “пузырька”.

```
1 #include <stdio.h>
2 #include <stdlib.h>
3 #include <string.h>
4
5 #define MAX_LINES 10
6 #define MAX_LINE_LENGTH 256
7
8 /* Функция для освобождения памяти. Она принимает массив указателей и освобождает
9    память по каждому указателю. */
10 void FreeMemory(char* lines[], int numLines) {
11     int i;
12
13     for (i = 0; i < numLines; i++)
14         free(lines[i]);
15 }
```

```
16
17 int main(void) {
18     /* Статический массив указателей (MAX_LINES элементов, тип элемента --- char*).
19     Инициализируем все элементы нулём. */
20     char* lines[MAX_LINES] = { 0 };
21     char buffer[MAX_LINE_LENGTH];
22     int numLines, i, k;
23
24     for (numLines = 0; numLines < MAX_LINES; numLines++) {
25         /* Считываем новую строку в массиве buffer. */
26         if (fgets(buffer, MAX_LINE_LENGTH, stdin) == NULL)
27             break;
28
29         /* Выделяем память под новую строку. */
30         lines[numLines] = (char*) malloc(sizeof(char) * (strlen(buffer) + 1));
31
32         if (lines[numLines] == NULL) {
33             /* Не удалось выделить память. */
34             printf("Can't allocate memory!\n");
35             FreeMemory(lines, numLines);
36             return 0;
37         }
38
39         /* Копируем строку в массиве lines. */
40         strcpy(lines[numLines], buffer);
41     }
42
43     /* Сортировка по возрастанию в лексикографическом порядке пузырьком. */
44     for (i = 0; i < numLines - 1; i++) {
45         for (k = 0; k < numLines - i - 1; k++) {
46             /* Функция strcmp() возвращает положительное число, если первая строка
47             больше второй в лексикографическом порядке; отрицательное число, если первая
48             строка меньше второй в лексикографическом порядке и 0, если строки равны. */
49             if (strcmp(lines[k], lines[k + 1]) > 0) {
50                 char* tmp;
51
52                 /* В случае неправильного прядка просто переставляем местами указатели.
53                 Никакого копирования не происходит. */
54                 tmp = lines[k];
55                 lines[k] = lines[k + 1];
56                 lines[k + 1] = tmp;
57             }
58         }
59     }
60
61     /* Выводим результат на экран. */
62     printf("Result:\n");
63     for (i = 0; i < numLines; i++)
64         printf("%d. %s", i + 1, lines[i]);
65
66     /* Освобождаем память. */
```

```
67 FreeMemory(lines, numLines);
68
69 return 0;
70 }
```

7.4 Разбиение строки, используя набор разделителей

Следующий пример демонстрирует разбиение строки при помощи набора разделителей.

```
1 /* Следующие друг за другом строковые литералы объединяются в одну строку. */
2 char buffer[256] = "There were four of us --- George, and William Samuel Harris, "
3                  "and myself, and Montmorency.";
4 const char* delimiters = "\r\n\t,?!?- ";
5 char* token;
6
7 token = strtok(buffer, delimiters);
8
9 while (token != NULL) {
10  printf("Token: %s\n", token);
11
12  /* В последующие разы строку передавать не нужно, функция strtok() её запомнит. */
13  token = strtok(NULL, delimiters);
14 }
```

Замечание 7.2 Функция `strtok()` не выделяет памяти. В переменную `token` записывается указатель на часть массива `buffer`. Таким образом, функция меняет строку, которая в неё передаётся, вставляя на место разделителей символы конца строки.

8 Несколько слов о Unicode

В этом разделе приведена лишь минимальная информация для ознакомления.

При работе с символами, отличными от ASCII, приходится уделять внимание выбору кодировки, поскольку тип `char` не может вместить в себя все возможные символы всех алфавитов. Это создаёт некоторые трудности, особенно в программах, которые должны работать с несколькими различными языками, поскольку каждый язык требует своей кодировки.

Для того, чтобы избежать лишних проблем с кодировками был введён стандарт Unicode. Однако и этот стандарт имеет несколько различных вариантов, например, UTF-16, в котором каждый символ кодируется 16 битами или UTF-32, в котором каждый символ кодируется 32 битами.

В языке C есть специальный тип данных `wchar_t`, предназначенный для хранения символов Unicode. Однако и здесь отсутствует какой-либо стандарт, на системах, поддерживающих UTF-32, размер `wchar_t` равен 4 байтам (например, в Linux), исключением является ОС Windows, в которой размер `wchar_t` равен 2 байтам. В стандартной библиотеке языка C есть специальные функции, работающие со строками из символов `wchar_t`, а также функции для того, чтобы преобразовывать обычные строки в “широкие” строки (массивы `wchar_t`) и обратно, используя различные кодировки.

В настоящее время часто используется кодировка UTF-8, в которой символы кодируются переменным числом байт. Например, ASCII символы в этой кодировке кодируются 1 байтом, а символы русского алфавита — двумя байтами. Во многих дистрибутивах Linux эта кодировка выбрана по умолчанию, например, в файловой системе или в текстовых редакторах.

Список литературы

- [1] *Брайан Керниган, Деннис Ритчи* Язык программирования Си.
- [2] <https://en.cppreference.com/w/c/language/sizeof>
- [3] <https://en.cppreference.com/w/c/types/limits>
- [4] <https://en.cppreference.com/w/c/locale/setlocale>