

Строки

М. А. Ложников

25 марта 2019 г.

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

1 Класс `std::string`

Для работы со строками в C++ используется шаблон `std::basic_string` с единственным обязательным аргументом — типом данных символа. Этот шаблон определён в заголовочном файле `<string>`. Он представляет из себя оболочку над обычным динамическим массивом, то есть хранит символы в памяти непрерывно. Для удобства в `<string>` определены конкретизации шаблона часто используемыми типами символов.

Так класс `std::string` является конкретизацией базового шаблона `std::basic_string` шаблонным параметром `char`. То есть, он совпадает с `std::basic_string<char>`. Класс подходит для работы со строками в однобайтных кодировках, если нужно работать с Unicode, то следует использовать класс `std::wstring` (равный `std::basic_string<wchar_t>`). Далее везде будем рассматривать лишь класс `std::string`.

1.1 Операции, специфичные для контейнера

Класс `std::string` поддерживает все те же операции, что и класс `std::vector<char>`. В самом деле

```
1 #include <iostream>
2 #include <string>
3 #include <algorithm>
4
5 int main() {
6     std::string str;
7
8     /* Метод push_back() кладёт элементы в конец строки. Ведёт себя также, как и
9        в векторе. */
10    str.push_back('a');
11    str.push_back('p');
12    str.push_back('p');
13    str.push_back('l');
14    str.push_back('e');
15
16    /* Методы size() и length() позволяют узнать длину строки в символах. Методы работают
17       одинаково. */
18    std::cout << "String '" << str << "' has " << str.size() << " symbols." << std::endl;
19
20    // Метод capacity() позволяет узнать размер выделенной памяти (в символах).
```

```
21 std::cout << "Allocated memory for " << str.capacity() << " symbols." << std::endl;
22
23 // Есть доступ по индексу как на чтение
24 for (std::size_t i = 0; i < str.size(); i++)
25     std::cout << "i = " << i << " sym = " << str[i] << std::endl;
26
27 // так и на запись
28 str[0] = 'A';
29
30 /* У класса std::string есть прямые/обратные константные/неконстантные итераторы.
31    Соответственно к строкам можно применять функции из библиотеки алгоритмов. */
32 // Например, найдём первый символ 'l'.
33 std::string::iterator it = std::find(str.begin(), str.end(), 'l');
34
35 /* Удалим этот символ. Функция вернёт итератор на символ, следующий за удалённым. */
36 it = str.erase(it);
37
38 /* Вставим обратно заглавный символ. Вместо одного символа можно передать
39    полуинтервал, который нужно вставить. Полуинтервал задаётся парой итераторов.
40    Метод возвращает итератор на первый вставленный элемент. */
41 str.insert(it, 'L');
42
43 /* Метод empty() позволяет узнать, пуста ли строка. */
44 std::cout << str.empty() << std::endl;
45
46 /* Метод clear() очищает строку. */
47 str.clear();
48
49 std::cout << str.empty() << std::endl;
50 return 0;
51 }
```

1.2 Ввод/вывод строк

В отличие от контейнеров `std::string` поддерживает операции ввода/вывода.

```
1 #include <iostream>
2 #include <string>
3 #include <sstream>
4 #include <cstdio> // Это для printf()
5
6 int main() {
7     /* Класс поддерживает неявную инициализация строкой в кавычках, то есть строкой
8        типа const char*. Если, например, функция принимает const std::string& или
9        просто std::string, то в неё можно передать строку в кавычках без явного
10       преобразования типа. */
11     std::string str = "some string";
12
13     /* Для класса std::string определены операции сравнения, сравнивающие строки в
14        лексикографическом порядке. В примере один из аргументов операции > имеет
```

```
15     min std::string, а второй к нему преобразуется неявным образом. */
16     if (str > "abc")
17         std::cout << "" << str << " is greater than 'abc'" << std::endl;
18
19     // Строки можно выводить в поток
20     std::cout << str << std::endl;
21
22     // Можно читать из потока
23     stringstream ss;
24     ss << "word" << std::endl << "some line 1" << std::endl << "some line 2";
25
26     // Можно читать отдельные слова
27     ss >> str;
28
29     /* Можно считывать целые строки до перевода на новую строку. Если указать третьим
30        аргументом разделитель, то функция будет считывать до разделителя. При этом
31        разделитель вытаскивается из потока, но не добавляется в строку. По умолчанию
32        разделитель равен '\n'. Функция возвращает поток. */
33     while (std::getline(ss, str))
34         std::cout << "Line: " << str << "" << std::endl;
35
36     /* Для тех, кому нравится работать со строками в стиле C можно получить указатель
37        на массив символов. */
38     str = "some string";
39     const char* data = str.data();
40     std::printf("%s\n", data);
41 }
```

1.3 Операции, специфичные для строк

Кроме того, `std::string` поддерживает операции, специфичные именно для строк. Поясним на примере.

```
1 #include <iostream>
2 #include <string>
3
4 int main() {
5     std::string str = "some string";
6
7     // Строки можно складывать. Операции + и += осуществляют конкатенцию строк.
8     str += " and another string";
9     str = str + " and the third string";
10
11     std::cout << str << std::endl;
12
13     /* Для класса std::string определены операции сравнения, сравнивающие строки в
14        лексикографическом порядке. Доступны все операции сравнения <, >, <=, >=, ==, !=.
15        Отметим, что все эти операции доступны и в векторе. В примере один из аргументов
16        операции > имеет min std::string, а второй к нему преобразуется неявным образом. */
17     if (str > "abc")
```

```
18     std::cout << "\"" << str << "' is greater than 'abc'" << std::endl;
19
20     /* Метод replace() заменяет подстроку новой строкой. Длины не обязаны совпадать.
21     У метода есть довольно много перегрузок, приведём одну из них. Первый аргумент
22     (first) задаёт позицию, второй (count) длину подстроки, которую нужно заменить.
23     Третьим аргументом передаётся подстрока, на которую нужно заменить. Метод
24     возвращает саму строку. */
25     // Заменяем первое вхождение "string" на "STRING".
26     str.replace(5, 6, "STRING");
27
28     std::cout << str << std::endl;
29
30     /* Метод принимает два аргумента: начало промежутка (first) и количество символов
31     (count). Метод копирует полуинтервал [first, first + count) в новую строку.
32     Позиции задаются индексами. По умолчанию, first равен нулю, а count равен
33     статической константе std::string::npos, означающей, что нужно скопировать символы
34     до конца строки. Метод возвращает новую строку. */
35     // Выделим подстроку из 10 символов от 5-го символа.
36     std::string anotherStr = str.substr(5, 10);
37
38     std::cout << anotherStr << std::endl;
39
40     // Выделим подстроку от 5-го символа до конца.
41     std::string anotherStr2 = str.substr(5);
42
43     std::cout << anotherStr2 << std::endl;
44
45     /* Вставка символов на определённую позицию. Отличие от обычных контейнеров
46     заключается в том, что позиция задаётся индексом. Поддерживается вставка
47     отдельных символов, строк типа std::string и const char*. Метод возвращает саму
48     строку, в которую было осуществлена вставка. */
49     // Вставим 3 символа 'a' на позицию 5.
50     std.insert(5, 3, 'a');
51
52     // Вставим строку "abcde" на позицию 10.
53     std.insert(10, "abcde");
54
55     /* У метода точно такие же аргументы, как и у substr(). Метод удаляет полуинтервал
56     [first, first + count). Метод возвращает саму строку. */
57     // Удалить полуинтервал [10, 15).
58     str.erase(10, 15);
59     // Удалить все символы от 5-ого.
60     str.erase(5);
61
62     return 0;
63 }
```

Также есть метод `copy()`, позволяющий скопировать подстроку в заранее выделенный массив элементов типа `char`.

1.4 Преобразование из строки в число и обратно

Для преобразования числа в строку используется функция `std::to_string()`. Поддерживаются как целые числа, так и числа с плавающей точкой.

```
1 std::string str;
2
3 str = std::to_string(10);
4
5 std::cout << str << std::endl;
6
7 str = std::to_string(3.14);
8
9 std::cout << str << std::endl;
```

Для конвертации строки в число используется семейство функций `std::sto*`. Например, для преобразования строки в целое типов `int` или `long` можно использовать `std::stoi()` или `std::stol()` соответственно. Для преобразования строки в числа типа `double` или `float` можно использовать функции `std::stod()` или `std::stof()` соответственно.

Все эти функции работают одинаково. Принимают строку первым аргументом и возвращают число. У каждой функции есть необязательный второй параметр типа `std::size_t*`, в который может быть записано число обработанных символов, то есть позиция первого необработанного символа. Функции для конвертации в целые числа имеют необязательный третий аргумент — основание системы исчисления, равный десяти по умолчанию. Функции игнорируют пробелы в начале строки. Функции могут выбросить исключение `std::invalid_argument`, если не было осуществлено преобразования или `std::out_of_range` если значение оказалось за пределами соответствующего типа данных. Все функции являются обёртками надо соответствующими функциями из `<cstdlib>` (например, `std::strtol()`, `std::strtod()`, `std::strtof()`).

```
1 std::string str = "1 2 3 4 abcde";
2
3 // Самый простой вариант
4 int num = std::stoi(str);
5
6 std::size_t pos = 0;
7
8 while(1) {
9     try {
10         // Считаем следующее число
11         num = std::stoi(str, &pos);
12
13         // Обрежем строку
14         str = str.substr(pos);
15
16         std::cout << "Read " << num << ". The remaining string equals " << str << std::endl;
17     } catch(std::invalid_argument&) {
18         std::cout << "Can't convert any more." << std::endl;
19         break; // Выходим из цикла
20     }
21 }
```

```
22
23 // Для чисел с плавающей точкой всё аналогично
24 str = "3.14";
25 double = std::stod(str);
```

1.5 Методы поиска

Метод `find()` предназначен для нахождения первого вхождения символа или подстроки в строке. Функция принимает символ (`char`) или подстроку (`const std::string&` или `const char*`) для поиска, а также необязательный параметр — позицию, с которой начинать поиск, равную по умолчанию нулю. Функция возвращает индекс первого вхождения или константу `std::string::npos` если вхождения не обнаружено.

Метод `rfind()` ищет последнее вхождение (с конца строки). Позиция начала поиска по умолчанию равна `std::string::npos`.

```
1 /* Разобьём строку на части, используя в качестве разделителя слово "and". */
2 std::string str = "123 and abcde and 456 and defg";
3 std::size_t next = str.find("and"); // Ищем первое вхождение
4 const std::size_t len = 3; // Длина искомой строки "and".
5 std::size_t prev = 0;
6
7 while (next != std::string::npos) {
8     // Проверяем, что токен не пуст
9     if (next > prev) {
10         std::string token = str.substr(prev, next - prev);
11         std::cout << "Found token '" << token << "'!" << std::endl;
12     }
13
14     prev = next + len;
15     next = str.find("and", prev); // Ищем следующее вхождение
16 }
17
18 if (prev != str.size()) { // Печатаем последнее слово
19     std::string token = str.substr(prev);
20     std::cout << "Found token '" << token << "'!" << std::endl;
21 }
```

Кроме того, есть методы для поиска в строке первого элемента из набора (`find_first_of()`), не из набора (`find_first_not_of()`). Параметры и возвращаемое значение у этих методов те же, что и у метода `find()`. Есть также методы, выполняющие поиск с конца.

```
1 /* Разобьём строку на части по набору разделителей. */
2 std::string str = "123 and abcde and 456 and defg";
3 std::size_t next = str.find_first_of(" \\t\\n\\r"); // Ищем первый разделитель
4 std::size_t prev = 0;
5
6 while (next != std::string::npos) {
7     // Проверяем, что токен не пуст
8     if (next > prev) {
```

```
9     std::string token = str.substr(prev, next - prev);
10     std::cout << "Found token '" << token << "'" << std::endl;
11 }
12
13 prev = next + 1;
14 next = str.find_first_of(" \t\n\r", prev); // Ищем следующий разделитель
15 }
16
17 if (prev != str.size()) { // Печатаем последнее слово
18     std::string token = str.substr(prev);
19     std::cout << "Found token '" << token << "'" << std::endl;
20 }
```

2 Класс `std::string_view` (C++17)

У класса `std::string` есть одна проблема: он не позволяет получить подстроку не скопировав её в отдельный экземпляр класса `string`. Таким образом при извлечении подстроки происходит выделение памяти и копирование символов.

Во многих задачах требуется получить некоторую подстроку только для чтения. Для того, чтобы сделать это без копирования элементов, в C++17 был введён невладеющий контейнер шаблон `std::basic_string_view` с одним обязательным шаблонным параметром — типом данных символа. Шаблон определён в заголовочном файле `<string_view>`. Невладеющий означает, что он не хранит данных внутри себя, он хранит лишь указатель на начало промежутка и указатель на конец промежутка. Таким образом, `std::basic_string_view` может ссылаться на непрерывный промежуток элементов в памяти.

Для работами со строками, состоящими из символов типа `char` в заголовочном файле `<string_view>` определён класс `std::string_view`, являющийся конкретизацией шаблона `std::basic_string_view` типом `char` (то есть `std::basic_string_view<char>`).

```
1 std::string str = "string";
2 const char* arr = "array of const chars";
3
4 /* std::string_view можно создать из класса std::string или из массива char'ов.
5  Объекты типа std::string_view можно неявно создавать из строк. */
6 std::string_view strView = str;
7
8 /* Сложность этого конструктора линейная, поскольку требуется подсчитать длину
9  строки. */
10 std::string_view arrView = arr;
11
12 /* std::string_view поддерживает все те же методы, предназначенные для чтения
13  элементов, что и std::string. */
14 std::cout << "View size equals " << strView.size() << std::endl;
15 std::cout << "The third symbol equals " << strView[2] << std::endl;
16 std::cout << strView.empty() << std::endl;
17
18 // Также можно получить указатель на const char
19 const char* data = strView.data();
```

Отметим, что `std::string_view` не позволяет менять символы в строке, на которую он ссылается. Кроме того, `std::string_view` может инвалидироваться в том случае, если освободится память, на которую он ссылается. Например, если добавить в строку типа `std::string`, на которую ссылается `std::string_view`, несколько символов, то это может привести к перевыделению памяти внутри контейнера `std::string` и следовательно к инвалидации `std::string_view`.

У класса `std::string_view` также есть итераторы, предоставляющие произвольный прямой/обратный доступ на чтение.

2.1 Извлечение подстрок

```
1 std::string str = "There were four of us --- George, and William Samuel Harris, "  
2                 "and myself, and Montmorency.";  
3 std::string_view view = str;  
4  
5 /* Метод substr() принимает два аргумента: первый (pos) --- начало подстроки  
6    и второй (count) количество символов в подстроке. Оба аргумента необязательный,  
7    pos по умолчанию равен 0, а count статической константе std::string_view::npos,  
8    означающей, что нужно брать подстроку до конца. Метод возвращает новый string_view  
9    на требуемую подстроку. */  
10 std::string_view george = view.substr(26, 6);  
11  
12 std::cout << george << std::endl;  
13  
14 std::string_view persons = view.substr(26);  
15  
16 std::cout << persons << std::endl;  
17  
18 /* В тех случаях, когда требуется изменить текущий string_view, можно воспользоваться  
19    методами, убирающими из начала или из конца N символов. Методы называются  
20    remove_prefix() и remove_suffix() соответственно. Их поведение не определено,  
21    если передать в них число, большее количества символов в string_view. */  
22  
23 str.remove_prefix(26);  
24  
25 std::cout << str << std::endl; // Выведет всё, начиная с "George".  
26  
27 str.remove_suffix(18);  
28  
29 std::cout << str << std::endl; // Выведет от "George" до "myself".
```

2.2 Методы поиска

В классе `std::string_view` есть все те же методы поиска, которые есть и в классе `string`. Эти методы принимают символ (типа `char`) или набор символов (либо `const char*` либо `std::string_view`, а следовательно и `string`) первым аргументом, кроме того, вторым необязательным аргументом они принимают позицию, с которой нужно начинать поиск. Методы возвращают позицию первого (или последнего) вхождения или константу `std::string_view::npos` если символ/подстрока не были найдены.

Напишем функцию, извлекающую токены по заданному набору разделителей:


```
1 std::string_view GetNextToken(std::string_view& str, std::string_view delimiters) {
2     std::string_view retval; // Изначально view пуст.
3
4     while (retval.size() == 0u) { // Избегаем подряд идущих разделителей
5         std::size_t pos = str.find_first_of(delimiters);
6
7         if (pos == std::string_view::npos)
8             return retval;
9
10        retval = str.substr(0, pos);
11
12        str.remove_prefix(pos + 1);
13    }
14
15    return retval;
16 }
```

Обратите, что объекты `std::string_view` достаточно легковесные, их можно передавать в функцию без ссылки.

Разобьём строку на слова при помощи предложенной функции.

```
1 std::string str = "There were four of us --- George, and William Samuel Harris, "
2                 "and myself, and Montmorency.";
3 std::string_view view = str;
4
5 std::string_view token = GetNextToken(view, " -,.");
6
7 while (token.size() != 0) {
8     std::cout << "Found token '" << token << "'!" << std::endl;
9     token = GetNextToken(view, " -,.");
10 }
```

Список литературы

- [1] https://en.cppreference.com/w/cpp/string/basic_string
- [2] https://en.cppreference.com/w/cpp/string/basic_string_view
- [3] <https://en.cppreference.com/w/cpp/regex>