

Methods Necessary for String Classes in Object-Oriented Programming Languages

1 INTRODUCTION

In programming, a string is a primitive type designed to provide an interface for an ordered set of characters. However, standard library implementations of this interface vary in different languages. Thus, it is clear that the programming community has not decided on the optimal implementation of a string class. Therefore, this paper aims to determine a minimalistic set of string functions sufficient for a string class.

2 SET BUILDING

2.1 Necessary and Sufficient Functions

The memory tape of the Turing machine can be represented as three strings: parts of the memory tape right and left from the pointer and the memory cell with the pointer. Three operations are necessary for a Turing machine: move the head to the left, move the head to the right, and write a symbol to the cell. Therefore, the following methods are necessary:

- (1) `concatenate`—concatenate two strings
- (2) `slice`—returns the substring of requested length, starting at a given position as a string
- (3) `length`—returns the length of the string

It is easy to see that these methods are sufficient to implement Turing machine commands. Therefore, all computable algorithms, including string-related algorithms, can be described through strings. Because standard libraries include regular expressions, function `regex-split`—split string into groups according to regular expression pattern.

2.2 Populating the Set

However, such an approach is not practical, and other methods need to be included. For this purpose, the following criteria for selection and algorithm for set construction are proposed:

2.2.1 *Criteria.*

- (1) Can not be implemented in three or fewer instructions using functions already in set
- (2) Can be implemented in the least amount of instructions among all remaining functions

Author's address:

2.2.2 Algorithm.

- (1) Functions `slice`, `regex`, and `concatenate` are added to the set
- (2) The function that falls under all criteria in 2.2.1 is added to the set
- (3) The second step is repeated until no fitting functions remain
- (4) Then, the set is hedged by deleting functions that do not fit requirements outlined in step two

2.3 Analysis of the final selected functions

After executing the algorithm, following functions were selected:

- (1) `concatenate`—concatenates two strings
- (2) `slice`—returns substring specified by arguments
- (3) `compare`—compares two strings
- (4) `capitalize`—converts the first character to the uppercase, and all other to the lowercase
- (5) `format`—returns formatted string
- (6) `hash`—returns the hash of the string
- (7) `is-alphabetic`—checks if the string contains only letters
- (8) `is-lowercase`—checks if the string contains only lowercase letters
- (9) `is-uppercase`—checks if the string contains only uppercase letters
- (10) `swap-case`—converts all lowercase characters to uppercase, and vice versa
- (11) `to-float`—converts string to float
- (12) `to-int`—converts string to integer

2.4 Functions in other languages

Libraries in different languages have different specifics. The most complete are the Java and Python language libraries. Libraries contain almost all the functions collected in the common table of 4 libraries. They are missing only the most specific of all the assembled functions. Because of this, they may have redundant functions that have one result. C++ library is the most minimalistic. Perhaps this is due to its age, because the language appeared almost 10 years earlier than the others in the selection. It lacks some common functions, but has a few unusual ones, mostly memory-related. For example, `capacity` - this function for strings is not available in Java, Python or Ruby. Ruby has the most balanced library. It has many features but almost no repetitive features like Java. This is because both Python and C++ have influenced the language. The main information support is statistics collected from various libraries. It can be used to calculate the importance of a function from the frequency of its use by other programmers. The choice is also based on the algorithm from part 2.2. Every language has functions not added to the EO language library:

Methods Necessary for String Classes in Object-Oriented Programming Languages

Function name	Libraries	Reasons of excluding
capacity	C++	Presented only in one library, has limited functionality, used only in 0,02% of general string functions usage
fill-left	Ruby	Presented only in one library, can be implemented in three or fewer instructions, used only in 0,002% of general string functions usage
fill-right	Ruby	Presented only in one library, can be implemented in three or fewer instructions, used only in 0,002% of general string functions usage
is-title	Python	Presented only in one library, can be implemented in three or fewer instructions, used only in 0,0001% of general string functions usage
repeat	Java, Python	Presented only in 2 libraries, can be implemented in three or fewer instructions, used only in 0,002% of general string functions usage
partition	Python, Ruby	Presented only in 2 libraries, used only in 0,003% of general string functions usage
shrink-to-fit	C++	Presented only in one library, has limited functionality, used only in 0,003% of general string functions usage
strip	Java, Python, Ruby	Can be implemented in three or fewer instructions, used only in 0,015% of general string functions usage
starts-with	Java, Python, Ruby	Can be implemented in three or fewer instructions, used only in 0,01% of general string functions usage

2.5 StringBuilder implementation

After discussing the details of the implementation of the functions, we pointed out that the focus was not on the efficiency of the functions, but purely on their availability and usability. Therefore it was decided at this stage of development to remove the implementation of the StringBuilder class from the general task list and to implement functions without efficient development. In the future, it may be possible to refine these functions with the addition of effective use by implementing the StringBuilder.

2.6 GitHub analysis

The team analyzed the libraries of several languages: Ruby, C++, Java, and Python. The data was collected from GitHub with an advanced search among all the code. The result was a list of the

most popular functions for each language involved. Some of the most popular functions have already been implemented, as well as functions related to regular expressions. In C++, the selected functions cover 76% of the uses of functions with strings in general. In Python this number is 44%, Ruby - 33%, Java - 38%. The small percentage in the last three languages can be explained by a large number of redundant functions in that language.

3 CONCLUSION

Therefore, after analyzing the selection of functions, we chose several: concatenate, substring, compare, capitalize, format, hash, is-alphabetic, is-lowercase, is-uppercase, to-float, swap-case, and to-int. In selecting the functions, we were guided by features such as sufficient quantity for the Turing machine, the algorithm in the "Populating the Set" section, and statistics on GitHub.