

58093 String Processing Algorithms

Lectures, Autumn 2013, period II
Juha Kärkkäinen

Contents

0. Introduction
1. Sets of strings
 - Search trees, string sorting, binary search
2. Exact string matching
 - Finding a pattern (string) in a text (string)
3. Approximate string matching
 - Finding in the text something that is similar to the pattern
4. Suffix tree and suffix array
 - Preprocess a long text for fast string matching and all kinds of other tasks

1

2

0. Introduction

Strings and sequences are one of the simplest, most natural and most used forms of storing information.

- natural language, biosequences, programming source code, XML, music, any data stored in a file
- Many algorithmic techniques are first developed for strings and later generalized for more complex types of data such as graphs.

The area of algorithm research focusing on strings is sometimes known as *stringology*. Characteristic features include

- Huge data sets (document databases, biosequence databases, web crawls, etc.) require efficiency. Linear time and space complexity is the norm.
- Strings come with no explicit structure, but many algorithms discover implicit structures that they can utilize.

3

4

Strings

An **alphabet** is the set of **symbols** or **characters** that may occur in a string. We will usually denote an alphabet with the symbol Σ and its size with σ .

We consider three types of alphabets:

- **Ordered alphabet:** $\Sigma = \{c_1, c_2, \dots, c_\sigma\}$, where $c_1 < c_2 < \dots < c_\sigma$.
- **Integer alphabet:** $\Sigma = \{0, 1, 2, \dots, \sigma - 1\}$.
- **Constant alphabet:** An ordered alphabet for a (small) constant σ .

5

6

A string is a **sequence** of symbols. The set of all strings over an alphabet Σ is

$$\Sigma^* = \bigcup_{k=0}^{\infty} \Sigma^k = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots$$

where

$$\begin{aligned} \Sigma^k &= \overbrace{\Sigma \times \Sigma \times \dots \times \Sigma}^k \\ &= \{a_1 a_2 \dots a_k \mid a_i \in \Sigma \text{ for } 1 \leq i \leq k\} \\ &= \{(a_1, a_2, \dots, a_k) \mid a_i \in \Sigma \text{ for } 1 \leq i \leq k\} \end{aligned}$$

is the set of strings of length k . In particular, $\Sigma^0 = \{\varepsilon\}$, where ε is the **empty string**.

We will usually write a string using the notation $a_1 a_2 \dots a_k$, but sometimes using (a_1, a_2, \dots, a_k) may avoid confusion.

7

8

About this course

On this course we will cover a few cornerstone problems in stringology. We will describe several algorithms for the same problems:

- the best algorithms in theory and/or in practice
- algorithms using a variety of different techniques

The goal is to learn a toolbox of basic algorithms and techniques.

On the lectures, we will focus on the clean, basic problem. Exercises may include some variations and extensions. We will mostly ignore any application specific issues.

The alphabet types are really used for classifying and analysing **algorithms** rather than alphabets:

- Algorithms for ordered alphabet use only character comparisons.
- Algorithms for integer alphabet can use more powerful operations such as using a symbol as an address to a table or arithmetic operations to compute a hash function.
- Algorithms for constant alphabet can perform almost any operation on characters and even sets of characters in constant time.

The assumption of a constant alphabet in the analysis of an algorithm often indicates one of two things:

- The effect of the alphabet on the complexity is complicated and the constant alphabet assumption is used to simplify the analysis.
- The time or space complexity of the algorithm is heavily (e.g., linearly) dependent on the alphabet size and the algorithm is effectively unusable for large alphabets.

An algorithm is called **alphabet-independent** if its complexity has no dependence on the alphabet size.

There are many notations for strings.

When describing algorithms, we will typically use the array notation to emphasize that the string is stored in an array:

$$\begin{aligned} S &= S[1..n] = S[1]S[2] \dots S[n] \\ T &= T[0..n] = T[0]T[1] \dots T[n-1] \end{aligned}$$

Note the half-open range notation $[0..n)$ which is often convenient.

In an abstract context, we often use other notations, for example:

- $\alpha, \beta \in \Sigma^*$
- $x = a_1 a_2 \dots a_k$ where $a_i \in \Sigma$ for all i
- $w = uv$, $u, v \in \Sigma^*$ (w is the **concatenation** of u and v)

We will use $|w|$ to denote the **length** of a string w .

Individual characters or their positions usually do not matter. The significant entities are the substrings or factors.

Definition 0.1: Let $w = xyz$ for any $x, y, z \in \Sigma^*$. Then x is a **prefix**, y is a **factor (substring)**, and z is a **suffix** of w . If x is both a prefix and a suffix of w , then x is a **border** of w .

Example 0.2: Let $w = \text{bonobo}$. Then

- $\epsilon, \text{b}, \text{bo}, \text{bon}, \text{bono}, \text{bonob}, \text{bonobo}$ are the prefixes of w
- $\epsilon, \text{o}, \text{bo}, \text{obo}, \text{nobo}, \text{onobo}, \text{bonobo}$ are the suffixes of w
- $\epsilon, \text{bo}, \text{bonobo}$ are the borders of w
- $\epsilon, \text{b}, \text{o}, \text{n}, \text{bo}, \text{on}, \text{no}, \text{ob}, \text{bon}, \text{ono}, \text{nob}, \text{obo}, \text{bono}, \text{onob}, \text{nobo}, \text{bonob}, \text{onobo}, \text{bonobo}$ are the factors of w .

Note that ϵ and w are always suffixes, prefixes, and borders of w . A suffix/prefix/border of w is **proper** if it is not w , and **nontrivial** if it is not ϵ or w .

9

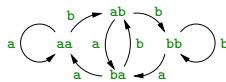
A **De Bruijn sequence** B_k of order k for an alphabet Σ of size σ is a **cyclic string** of length σ^k that contains every string of length k over the alphabet Σ as a factor exactly once. The cycle can be opened into a string of length $\sigma^k + k - 1$ with the same property.

Example 0.4: De Bruijn sequences for the alphabet $\{a, b\}$:

$$\begin{aligned} B_2 &= \text{aabb(a)} \\ B_3 &= \text{aaababbb(aa)} \\ B_4 &= \text{aaaabaabbababbbb(aaa)} \end{aligned}$$

De Bruijn sequences are not unique. They can be constructed by finding Eulerian or Hamiltonian cycles in a **De Bruijn graph**.

Example 0.5: De Bruijn graph for the alphabet $\{a, b\}$ that can be used for constructing B_2 (Hamiltonian cycle) or B_3 (Eulerian cycle).



11

An efficient execution of the operations requires that the set is stored as a suitable data structure.

- A **binary search tree** supports the basic operations in $\mathcal{O}(\log n)$ time and range searching in $\mathcal{O}(\log n + r)$ time, where n is the size of the set and r is the size of the result.
- An **ordered array** supports lookup and range searching in the same time as binary search trees. It is simpler, faster and more space efficient in practice, but does not support insertions and deletions.
- A **hash table** supports the basic operations in constant time (usually randomized) but does not support range queries.

A data structure is called **dynamic** if it supports insertions and deletions (tree, hash table) and **static** if not (array). Static data structures are constructed once for the whole set of objects. In the case of an ordered array, this involves another important operation, **sorting**. Sorting can be done in $\mathcal{O}(n \log n)$ time using comparisons and even faster for integers.

13

Trie

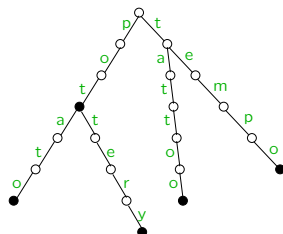
A simple but powerful data structure for a set of strings is the **trie**. It is a **rooted tree** with the following properties:

- Edges are labelled with symbols from an alphabet Σ .
- For every node v , the edges from v to its children have different labels.

Each node represents the string obtained by concatenating the symbols on the path from the root to that node.

The trie for a string set \mathcal{R} , denoted by $\text{trie}(\mathcal{R})$, is the smallest trie that has nodes representing all the strings in \mathcal{R} . The nodes representing strings in \mathcal{R} may be marked.

Example 1.1: $\text{trie}(\mathcal{R})$ for $\mathcal{R} = \{\text{pot}, \text{potato}, \text{pottery}, \text{tattoo}, \text{tempo}\}$.



15

Some Interesting String

The **Fibonacci strings** are defined by the recurrence:

$$\begin{aligned} F_0 &= \epsilon \\ F_1 &= \text{b} \\ F_2 &= \text{a} \\ F_i &= F_{i-1}F_{i-2} \text{ for } i > 2 \end{aligned}$$

The infinite Fibonacci string is the limit F_∞ . For all $i > 1$, F_i is a prefix of F_∞ .

Fibonacci strings have many interesting properties:

- $|F_i| = f_i$, where f_i is the i th Fibonacci number.
- F_∞ has exactly $k + 1$ distinct factors of length k .
- For all $i > 1$, we can obtain F_i from F_{i-1} by applying the substitutions $\text{a} \mapsto \text{ab}$ and $\text{b} \mapsto \text{a}$ to every character.

Example 0.3:

$$\begin{aligned} F_3 &= \text{ab} \\ F_4 &= \text{aba} \\ F_5 &= \text{abaab} \\ F_6 &= \text{abaababa} \\ F_7 &= \text{abaababaabaab} \\ F_8 &= \text{abaababaabaabaabaab} \end{aligned}$$

10

1. Sets of Strings

Basic operations on a set of objects include:

- Insert:** Add an object to the set
- Delete:** Remove an object from the set.
- Lookup:** Find if a given object is in the set, and if it is, possibly return some data associated with the object.

There can also be more complex queries:

- Range query:** Find all objects in a given range of values.

There are many other operations too but we will concentrate on these here.

12

The above time complexities assume that basic operations on the objects including comparisons can be performed in constant time. When the objects are strings, this is no more true:

- The worst case time for a string comparison is the length of the shorter string. Even the average case time for a random set of n strings is $\mathcal{O}(\log_\sigma n)$ in many cases, including for basic operations in a balanced binary search tree. We will show an even stronger result for sorting later. And sets of strings are rarely fully random.
- Computing a hash function is slower too. A good hash function depends on all characters and cannot be computed faster than the length of the string.

For a string set \mathcal{R} , there are also new types of queries:

- Prefix query:** Find all strings in \mathcal{R} that have the query string S as a prefix. This is a special type of range query.
- Lcp (longest common prefix) query:** What is the length of the longest prefix of the query string S that is also a prefix of some string in \mathcal{R} .

Thus we need special set data structures and algorithms for strings.

14

The trie is conceptually simple but it is not simple to implement efficiently. The time and space complexity of a trie depends on the implementation of the **child function**:

For a node v and a symbol $c \in \Sigma$, $\text{child}(v, c)$ is u if u is a child of v and the edge (v, u) is labelled with c , and $\text{child}(v, c) = \perp$ (null) if v has no such child.

As an example, here is the insertion algorithm:

Algorithm 1.2: Insertion into trie
Input: $\text{trie}(\mathcal{R})$ and a string $S[0..m] \notin \mathcal{R}$
Output: $\text{trie}(\mathcal{R} \cup \{S\})$

- (1) $v \leftarrow \text{root}; j \leftarrow 0$
- (2) **while** $\text{child}(v, S[j]) \neq \perp$ **do**
- (3) $v \leftarrow \text{child}(v, S[j]); j \leftarrow j + 1$
- (4) **while** $j < m$ **do**
- (5) Create new node u (initializes $\text{child}(u, c)$ to \perp for all $c \in \Sigma$)
- (6) $\text{child}(v, S[j]) \leftarrow u$
- (7) $v \leftarrow u; j \leftarrow j + 1$
- (8) Mark v as representative of S

16

There are many implementation options for the child function including:

Array: Each node stores an array of size σ . The space complexity is $\mathcal{O}(\sigma N)$, where N is the number of nodes in $\text{trie}(\mathcal{R})$. The time complexity of the child operation is $\mathcal{O}(1)$. Requires an integer alphabet.

Binary tree: Replace the array with a binary tree. The space complexity is $\mathcal{O}(N)$ and the time complexity $\mathcal{O}(\log \sigma)$. Works for an ordered alphabet.

Hash table: One hash table for the whole trie, storing the values $\text{child}(v, c) \neq \perp$. Space complexity $\mathcal{O}(N)$, time complexity $\mathcal{O}(1)$. Requires an integer alphabet.

A common simplification in the analysis of tries is to assume a **constant alphabet**. Then the implementation does not matter: Insertion, deletion, lookup and lcp query for a string S take $\mathcal{O}(|S|)$ time.

Note that a trie is a complete representation of the strings. There is no need to store the strings separately.

Many data structures and algorithms for a string set \mathcal{R} become simpler if \mathcal{R} is prefix free.

Definition 1.3: A string set \mathcal{R} is **prefix free** if no string in \mathcal{R} is a prefix of another string in \mathcal{R} .

If \mathcal{R} is prefix free, the **leaves** of $\text{trie}(\mathcal{R})$ represent exactly \mathcal{R} . This simplifies the implementation of the trie:

- Only internal nodes need the child data structure.
- Only leaves need the representation markers.

There is a simple way to make any string set prefix free:

- Let $\$ \notin \Sigma$ be an extra symbol satisfying $\$ < c$ for all $c \in \Sigma$.
- Append $\$$ to the end of every string in \mathcal{R} .

This has little or no effect on most operations. The length of each string increases by one only, and the additional symbol could be there only virtually.