

# Bioinformatics Algorithms

(Fundamental Algorithms, module 2)

Zsuzsanna Lipták

Masters in Medical Bioinformatics  
academic year 2018/19, II semester

## Strings and Sequences in Computer Science

### Some formalism on strings

- $\Sigma$  a finite set called **alphabet**
- its elements are called **characters** or **letters**
- $|\Sigma|$  is the **size** of the alphabet (number of different characters)
- a **string over  $\Sigma$**  is a finite sequence of characters from  $\Sigma$
- we write strings as  $s = s_1 s_2 \dots s_n$  i.e.  $s_i$  is the  $i$ 'th character of  $s$

N.B.: We number strings from 1, not from 0

2/7

### Some formalism on strings (cont.)

- $|s|$  is the **length** of string  $s$
- $\epsilon$  is the **empty string**, the (unique) string of length 0
- $\Sigma^n$  is the set of strings of length  $n$
- $\Sigma^* = \bigcup_{n=0}^{\infty} \Sigma^n = \Sigma^0 \cup \Sigma^1 \cup \Sigma^2 \cup \dots$  is the set of all strings over  $\Sigma$

3/7

### Some formalism on strings: Examples

#### Examples

- **DNA:**  $\Sigma = \{A, C, G, T\}$ , alphabet size  $|\Sigma| = 4$ ,  
 $s = ACCTG$  is a string of length 5 of  $\Sigma$ , with  
 $s_1 = A, s_2 = s_3 = C, s_4 = T, s_5 = G$ .
- **RNA:**  $\Sigma = \{A, C, G, U\}$ , again alphabet size is 4
- **protein:**  $\Sigma = \{A, C, D, E, F, \dots, W, Y\}$ , alphabet size is 20,  
ANRFYWNL is a string over  $\Sigma$  of length 8
- **English alphabet:**  $\Sigma = \{a, b, c, \dots, x, y, z\}$  of size 26,  
alphabet is a string over  $\Sigma$  of length 8

4/7

### Some formalism on strings

Let  $s = s_1 \dots s_n$  be a string over  $\Sigma$ .

ex.  $s = ACCTG$

- $t$  is a **substring** of  $s$  if  $t = \epsilon$  or  $t = s_i \dots s_j$  for some  $1 \leq i \leq j \leq n$   
(i.e., a "contiguous piece" of  $s$ ) CCT, AC, ...
- $t$  is a **prefix** of  $s$  if  $t = \epsilon$  or  $t = s_1 \dots s_j$  for some  $1 \leq j \leq n$   
(i.e., a "beginning" of  $s$ ) AC, ACCTG, ...
- $t$  is a **suffix** of  $s$  if  $t = \epsilon$  or  $t = s_i \dots s_n$  for some  $1 \leq i \leq n$   
(i.e., an "end" of  $s$ ) CCTG, G, ...
- $t$  is a **subsequence** of  $s$  if  $t$  can be obtained from  $s$  by deleting some  
(possibly 0, possibly all) characters from  $s$  AT, CCT, ...

N.B.

string = sequence, but substring  $\neq$  subsequence!

5/7

### Substrings etc.

N.B.

1. Every substring is a subsequence, but not every subsequence is a substring!  
**Ex.:** Let  $s = ACCTG$ , then ACT is a subsequence but not a substring.
2. Every prefix and every suffix is a substring.
3.  $t$  is substring of  $s \Leftrightarrow t$  is prefix of a suffix of  $s \Leftrightarrow t$  is suffix of a prefix of  $s$

6/7

## Counting substrings, subsequences etc.

### Question

Given  $s = s_1 \dots s_n$ . How many

- prefixes,
- suffixes,
- substrings,
- subsequences

does  $s$  have (exactly, or at most, or at least)?