

Progetto ASD 2018/19

DNA \rightarrow "sequenziazione"
sequencing

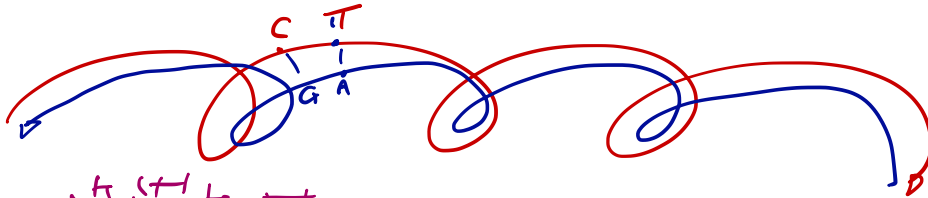
DNA \rightarrow sequenza sull'alfabeto

$\Sigma = \{A, C, G, T\}$
basi

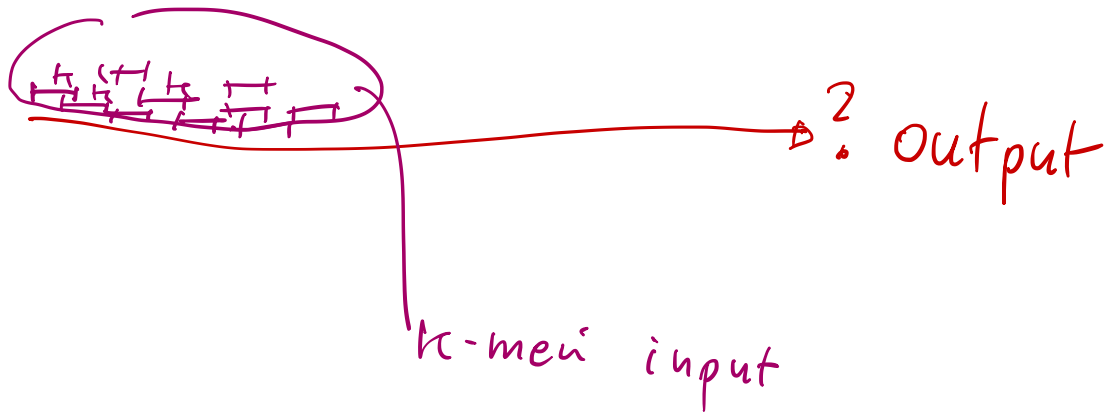
HTS = High-throughput sequencing

k-mero = sequenza di k simboli di Σ
k-mer

(N = simbolo
incerto)



$k \sim 10^2 - 10^3$ ordine # basi \rightarrow milioni di basi



Nostre terminologia:

INPUT = file contenente un numero elevato di stringhe
sull'alfabeto Σ , ciascuna di lunghezza k

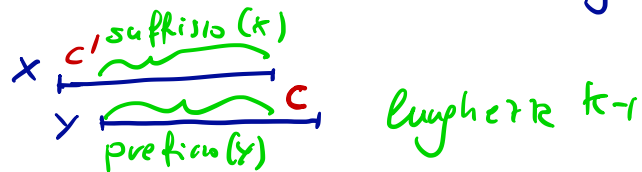
Scopo del progetto:

▷ grafo di de Bruijn: Σ, k

$G = (V, E)$ orientato

$$V = \Sigma^k$$

$$E = \{ (x, y) : x, y \in \Sigma^k \text{ e } \left. \begin{array}{l} \text{suffisso}(x) = \\ \text{prefisso}(y) \end{array} \right\} \text{ di lunghezza } k-1$$

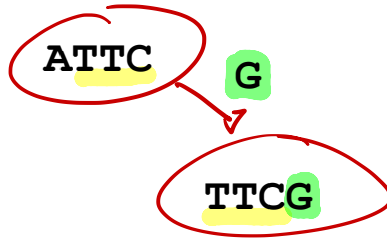


$$L: E \rightarrow \Sigma$$

$$L(x, y) = c \quad \text{t.c.} \quad \begin{array}{l} x = c'd \\ y = dc \end{array}, \quad |d| = k-1$$

Esempio:

$k=4$



Grafo regolare in cui ogni nodo ha $|\Sigma|$ archi uscenti e $|\Sigma|$ archi entranti

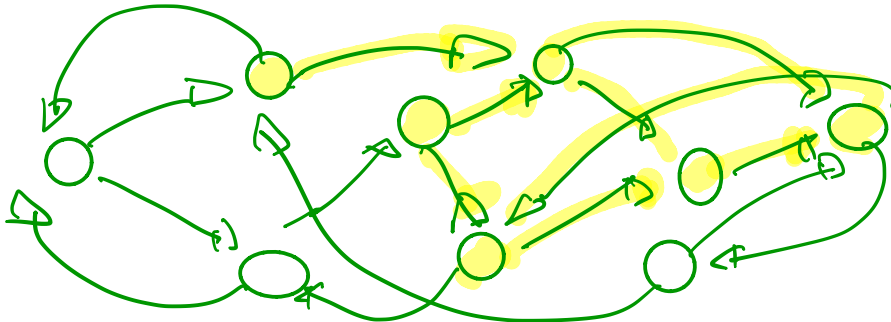
Dato $G=(V,E)$ e un sottoinsieme $V' \subseteq V$

sottografo indotto $G[V'] = (V', E[V'])$

$$\text{t.c. } E[V'] = \{(x,y) \in E : x,y \in V'\}$$

(note: alcuni nodi potrebbero essere isolati e quindi possono essere scartati a priori)

$|\Sigma|=2$



NOTA:

V' è dato dal file di INPUT

Una tecnica comune di sequenziamento, percorre un opportuno cammino (con nodi ripetuti) nel grafo, sovrapponendo i rispettivi k -meri.



$$dC, c_2$$

$$|α| = k$$

$$|α'| = k-1$$

$$|α''| = k-2$$

In realtà ci sono diverse complicazioni.

Task Dato il file di input:

- ① Costruire il prefisso di de Bruijn dell'input
- ② Usare opportune strutture di dati per rispondere alle seguenti richieste (query):

2.a. Data una stringa P su Σ , di lunghezza arbitraria m , stabilire se P appare lungo un qualche cammino (il problema è interessante per $m > k$)

2.b. Data P su Σ , trovare il più lungo prefisso di P che soddisfa 2.a (chiaro che se P appare, abbiamo che è lui la risposta)

2.c. Risolvere la 2.a dove P può avere un errore: uno dei simboli di P non corrisponde, ma gli altri sì.
 ATCC corrisponde ad A Gcc