

# Pomozte Craigu Venterovi hledat geny v DNA

V ústavu Všeobecné Pozemské Genomiky, který nedávno založil světoznámý neúnavný průkopník v biologii Craig Venter, vytvořili nový přístroj, jímž lze bezdotykově během 3 minut přečíst celou DNA libovolné osoby. Přístroj je založen na velmi přesné laserové technologii a využívá kvantové efekty probíhající při kontrakci podkožních svalů. Přístroj, resp. jeho detekční část v podobě malé plastické pistolky stačí zaměřit na např. na vnitřek dlaně a požádat měřenou osobu, aby během tří minut zlehka pohybovala palcem této ruky. Během snímání posílá detekční část bezdrátově vysokou rychlostí data do připojeného počítače, kde pokročilé algoritmy rekonstruují nejvýše do 30 sekund po skončení snímání celou sekvenci DNA.

Při testování tohoto přístroje, které mimochodem probíhalo i v Čechách v Kladně na ČVUT, se výzkumníci zaměřili na čtení určitých genů způsobujících vysokou pravděpodobnost ztráty dětské senné rýmy v adolescenci. Jednotlivé geny v tomto souboru jsou označeny jako S1, S2, ..., SN a jsou identické u všech jednotlivců v populaci. Kromě toho se zjistilo, že se v molekule DNA vyskytují na poměrně libovolném místě, vždy těsně za sebou, ale jejich pořadí může být zcela libovolné.

Výzkum je mírně komplikován nečekanou okolností. Měření kvantových efektů, samo o sobě poměrně citlivé, může být ovlivněno průchodem superrychlých neutronů planetou Zemí a tudíž také měřicí částí zařízení. Kolize tau-neutrina s jádrem nestabilního izotopu astatu 213 užitého v detektoru způsobí chybné čtení jednoho písmene DNA. Protože hustota toku neutronů v daném čase a místě je poměrně dobře tabelována, lze předpovědět, jaký je maximální počet chybně přečtených písmen DNA v předpokládaném souboru genů. Toto číslo je označeno jako MEGSR (Maximum Errors in Gene Set Reading) a je obecně různé pro různá čtení DNA i u téže osoby.

Abychom výzkumníkům pomohli, máme sestavit program, který na základě přečtené sekvence DNA, znalosti genů S1, S2, ..., SN a čísla MEGSR určí, zda daná osoba má tento soubor genů ve své DNA. Z testovacích důvodů se neomezíme jen na standardní čtyři písmena A, C, T, G označující báze DNA, ale budeme předpokládat libovolnou abecedu, z níž je sestavena DNA i jednotlivé geny. Budeme předpokládat, že geny se vždy vyskytují v libovolném pořadí těsně za sebou, mezi sousedními dvěma v dané DNA není žádné další písmeno. Délka každého genu je nejvýše 20 000 znaků, délka DNA (v našem cvičném případě) je nejvýše 1 000 000 znaků, počet genů nepřesáhne 1000. Všechny geny jsou navzájem různé.

## Vstup

První řádek vstupu specifikuje použitou abecedu. Abeceda je zadána jako řetězec bez mezer, v němž jsou všechny znaky navzájem různé. Pořadí znaků v abecedě je nepodstatné. Druhý řádek vstupu obsahuje jediný dlouhý řetězec nad danou abecedou představující přečtenou sekvenci DNA. Na třetím řádku je uvedeno jediné číslo N určující počet genů v souboru. Na dalších N řádcích je uveden seznam genů, na každém řádku jeden gen. Gen, podobně jako DNA, je neprázdným řetězcem znaků nad danou abecedou. Geny jsou v souboru genů indexovány od 1 do N v tom pořadí, v jakém byly načteny. Na posledním řádku vstupu je uvedeno jediné nezáporné celé číslo MEGSR, určující maximální celkový počet chybně přečtených znaků DNA v místě, které je obsazeno geny daného souboru.

## Výstup

Na výstupu je seznam možných pozic daného souboru genů v dané DNA. Každý prvek seznamu je uložen na jednom řádku a má následující strukturu. Necht'  $S_k$  ( $1 \leq k \leq N$ ) je první gen daného souboru bráno od počátku DNA, za kterým následují ostatní geny. Nejprve je uvedena pozice prvního znaku genu  $S_k$  v DNA. Pozice v DNA jsou číslovány odleva počínaje 1. Za číslem pozice následuje N čísel, určujících pořadí genů v DNA. Nejprve je uveden index k a dále indexy jednotlivých dalších genů v tom pořadí, v jakém se geny vyskytují v DNA. Všechny hodnoty na řádku jsou odděleny jednou mezerou.

Prvky seznamu na výstupu jsou uvedeny v rostoucím lexikografickém pořadí, přičemž jeden řádek výstupu pokládáme za vektor s  $N+1$  hodnotami, podle nichž lexikograficky řadíme. Například

```
...
211 2 1 3 4
211 2 4 3 1
211 3 1 2 4
212 1 2 3 4
atd ...
```

Pokud se daný soubor genů v dané DNA i se započtením možných chybných detekcí nemůže vyskytovat, vystoupí jediný řádek obsahující číslo 0.

## Příklad 1

Vstup:

```
AbC
bbbbCAACAAACbbCCAbbAAAb
3
AAA
```

bbCC  
AC  
2

Výstup:

3 2 3 1  
7 3 1 2  
8 1 3 2  
10 1 2 3  
11 3 2 1  
13 2 3 1

## Příklad 2

Vstup:

01  
11011011011  
2  
0110  
010  
1

Výstup:

3 1 2  
3 2 1