

III. Formate de fișiere pentru secvențe biologice. Baze de date biologice

1. Noțiuni introductive

Bioinformatica se ocupă cu prelucrarea cantității enorme de date biologice. În urma secvențierii unui genom sau a unei proteine, secvența respectivă, alcatuită din cele patru nucleotide, în cazul unui genom, sau din aminoacizi, în cazul unei proteine, este stocată într-o bază de date publică pentru a putea fi studiată de numeroși cercetători.

Există diferite programe de analiză și prelucrare a secvențelor biologice. Astfel, pentru ca datele referitoare la o secvență să poată circula între diferite aplicații informatice, au fost definite mai multe formate de fișiere utilizate în bioinformatică, precum:

FASTA
ASN. 1
GenBank
PDB
XML

Pentru a putea publica descoperirea unei noi gene sau a unui fragment genomic, trebuie trimisă secvența către bazele de date biologice publice pentru a putea obține un număr de acces („accession number”). Numărul de acces este unic pentru fiecare secvență trimisă, pentru a putea fi identificată printre numeroasele secvențe descoperite zilnic de cercetători din toată lumea.

2. Baza de date biologice

NCBI (National Center for Biotechnology Information) pune la dispoziție baze de date biologice, împreună cu un SGBD (care include motorul de căutare *Entrez* și diferite programe și utilitare pentru gestiunea și prelucrarea înregistrărilor): <http://www.ncbi.nlm.nih.gov>

Secvențele de nucleotide (ADN) sau cele de aminoacizi (proteine) au un număr de acces unic într-o bază de date. Motorul de căutare *Entrez* caută informații disponibile în mai multe baze de date, astfel se poate obține aceeași secvență de mai multe ori cu numere de acces diferite. NCBI pune la dispoziție o bază de date *RefSeq* în care sunt incluse informații non-redundante.

Secvențele din *RefSeq* au numere de acces de forma: NT_123456 – regiuni de cromozomi; NC_123456 – cromozom; NP_123456 – proteina; NG_123456 – gena;

Genomurile secvențiate până la ora actuală pot fi găsite pe NCBI: <http://www.ncbi.nlm.nih.gov/sites/genome>

- a. Să se acceseze următoarea adresă de la NCBI (National Center for Biotechnology Information):
<http://www.ncbi.nlm.nih.gov/guide/genomes-maps/>
- b. Să se acceseze baza de date *Genome*.
- c. Să se acceseze un grup de organisme și, apoi, un număr de acces (*accession number*)
- d. pentru un genom cromozomial al unui organism, de exemplu: NC_005562 .
- e. Să se acceseze baza de date *RefSeq*: NC_005562.
- f. Să se salveze informația găsită, în formatele: FASTA, ASN. 1, GenBank.
- g. Să se salveze doua secvente genomice diferite (unul procariot si unul eucariot) în fiecare din cele trei formate.

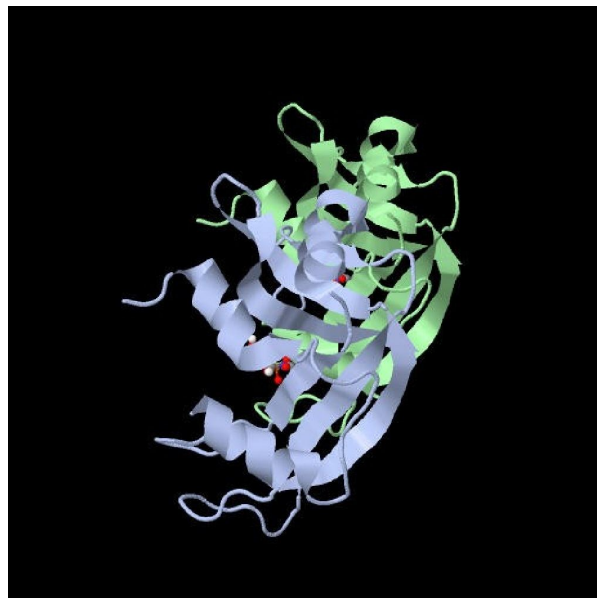
3. Formatul FASTA (Fast Alignment)

FASTA reprezintă un format de fișier foarte utilizat pentru date de tip secvență, precum lanțuri de nucleotide sau de aminoacizi. Formatul este foarte simplu:

- pe prima linie conține un header, caracterul ">", urmat de informații despre secvență (numărul de identificare, specia);
- următoarele linii reprezintă secvența de aminoacizi sau nucleotide.
- liniile au maxim 120 de caractere, dar în general nu depășesc 80. Fiecare linie se termină cu caracterul de sfârșit de linie "\n".

De exemplu, un fisier FASTA care descrie secventa unei proteine, arata astfel:

```
>sp|P61823|RNASE1_BOVIN Ribonuclease pancreatic OS=Bos taurus
GN=RNASE1
MALKSLVLLSLLVLVLLLVVRVQPSLGKETAAAKFERQHMDSSSTAASS
SNYCNQMMKSRNL
TKDRCKPVNTFVHESLADVQAVCSQKNVACKNGQTNCYQSYSTMSIT
DCRETGSSKYPNC AYKTTQANKHIIVACEGNPYVPVHFDASV
```



Ribonucleaza pancreatică bovină

Literele reprezintă aminoacizii proteinei secvențiate. În acest caz este reprezentată enzima ribonucleaza pancreatică bovină. Ribonucleazele sunt proteine enzimatice care degradează moleculele de ARN care nu mai sunt utile în cadrul celulei.

Să se descarce o secvență de nucleotide și o secvență proteică de pe www.ncbi.nlm.nih.gov în format FASTA. **Să se implementeze în Python următoarele cerințe:**

- a. Să se citească secvența de nucleotide, linie cu linie, și să se afișeze într-un fișier text .fasta, concomitent cu citirea.
- b. Să se concateneze două fișiere FASTA într-un singur fișier FASTA.
- c. Să se citească secvența proteinei și să se memoreze într-o variabilă, sub forma unui șir de caractere fără spații. Să se afișeze variabila obținută.

4. Formatul ASN. 1

Seq-entry ::= set { informații și secvențe
}

1. conține informații despre genomul secvențiat;
2. conține informații despre autorii care au participat la secvențiere;
3. secvența de ADN apare imediat după apariția șirului **ncbi2na** ' ;
4. secvența de ADN se termină la întâlnirea caracterului ' ;
5. pot exista mai multe regiuni de ADN care împreună alcătuiesc cromozomul respectiv, fiecare regiune de ADN fiind semnalată de prezența câmpului ncbi2na;
6. secvența de ADN este afișată într-un mod mai economic (fiind redusă la jumătate ca dimensiune), astfel:
 - cele 4 nucleotide din ADN sunt codificate de 2 biți, astfel:
A – 00
C – 01
G – 10
T – 11
 - se afișează câte 2 nucleotide (4 biți) ca un număr Hexazecimal:
0000 – 0
0001 – 1
0010 – 2
0011 – 3
0100 – 4
0101 – 5
0110 – 6
0111 – 7
1000 – 8
1001 – 9

1010 – A
1011 – B
1100 – C
1101 – D
1110 – E
1111 – F

De exemplu, următoarea secvență de nucleotide devine următorul șir în format ASN. 1: ATTCGAAC → 3C81

- a. **Să se citească dintr-un fișier format ASN. 1** prima secvență de ADN care apare și să se scrie într-un alt fișier în format FASTA.

Observații:

- Pentru citire se utilizează doar funcția *readline()*. Citirea se face linie cu linie, concomitent cu scrierea. Nu se citește toată secvența de ADN într-o variabilă, deoarece este prea mare.
- Nu se utilizează cicluri *for/while* pentru efectuarea conversiei, ci se utilizează funcțiile:
 - *map(f, data)* – aplică funcția *f* pe elementele din *data* (*data* poate fi lista, sir de caractere, etc.) și întoarce rezultatele într-o lista
 - *sir.join(lista)* – concatenează stringurile din lista *lista* și pune între ele șirul *sir* pentru care este apelată – va întoarce șirul obținut în urma concatenării.
 - Exemplu:
a= ['1','2','3','4','5','6']
print '&'.join(a)
print ".join(a)
- Pentru citirea primei secvențe de ADN (prima regiune de ADN) **nu** se utilizează cicluri imbricate. Pentru rezolvarea cerinței a. nu se utilizeaza cicluri *for/while* imbricate.

- b. Sa se gaseasca un fisier ASN.1 care contine mai multe regiuni.**
Să se citească toate regiunile de ADN din fișierul în format ASN. 1 și să se scrie un fișier multisecvență FASTA, de forma:

```
>Regiunea1  
ATTTGGG.....  
.....
```

```
>Regiunea2  
CCCTTAA.....  
.....
```

Pentru rezolvarea cerinței **b.** se utilizează maxim 2 cicluri *for/while* imbricate.

Observații:

- Funcția *sir.find('xyz')* întoarce indicele de început al subșirului 'xyz' în șirul *sir*.
- Să se utilizeze dicționarul *code* și funcția *conv* pentru conversie și să se implementeze funcția: *DNAfromASN1(fin, fout)* care primește ca intrare fișierul ASN. 1 din care va citi și un fișier de ieșire în care va scrie.

```
code = {  
'0' : 'AA', '1' : 'AC', '2' : 'AG', '3' : 'AT', '4' : 'CA', '5' : 'CC', '6' : 'CG', '7' : 'CT', '8'  
      : 'GA', '9' : 'GC', 'A' : 'GG', 'B' : 'GT', 'C' : 'TA', 'D' : 'TC', 'E' : 'TG', 'F' : 'TT'  
}
```

```
def conv(x):  
    return code[x]
```

```
def DNAfromASN1(fin, fout):  
    pass
```

```
if __name__ == '__main__':  
    DNAfromASN1('seq.asn1', 'seq.fasta')
```