

Laborator 4 - Bioinformatică

Formatul GenBank

1. conține informații despre genomul secvențiat;
2. conține informații despre autori;
3. secvența de ADN este ultima intrare în fișier și începe după cuvântul ORIGIN;
4. secvența de ADN se termină la apariția șirului "//";
5. secvența de ADN (o catenă) este scrisă pe linii; fiecare linie a secvenței de ADN conține 60 de nucleotide grupate câte 10, iar grupurile de nucleotide sunt separate printr-un spațiu (ultima linie din secvență poate avea mai puțin de 60 de nucleotide);
6. genele codificate în secvența de ADN sau în complementara acesteia sunt semnalate de cuvântul *gene* (când este precedat de spații) și/sau CDS;
7. imediat după semnalarea genei, apare locația acesteia; locația unei gene în secvența de ADN se poate găsi în următoarele situații:
 - *6981..9121* → începe la poziția 6981 din secvența de ADN din fișier și se termină la poziția 9121
 - *complement(2630..3110)* → începe la poziția 2630 și se termină la poziția 3110 din secvența de ADN complementară celei din fișier (catena complementară din dublul helix)
 - *join(1223..1567, 2110..2304)* → gena este alcătuită din două segmente din secvența de ADN găsită în fișier – segmentul care începe la poziția 1223 și se termină la poziția 1567 și segmentul care începe la poziția 2110 și se termină la poziția 2304 (gena poate fi alcătuită și din mai mult de două segmente, acestea fiind separate prin ",")
 - *complement (join(1223..1567, 2110..2304))* → gena este alcătuită din mai multe secvențe de pe catena de ADN complementară celei găsite în fișier
 - *7193..>8061* → semnul ">" arată că nu se știe exact unde se termină gena; pentru a extrage locația unei regiuni codificatoare trebuie tratate toate aceste situații
8. Cuvintele cheie încep de la caracterul al 6-lea (coloana 6 în fișier) și au o lungime maximă de 15 caractere. Locația începe de la caracterul al 22-lea (coloana 22 în fișier). (<ftp://ftp.ncbi.nih.gov/genbank/gbrel.txt>)

Aplicații. Formatul GenBank

1. Să se acceseze baza de date *Genome* de la NCBI (National Center for Biotechnology Information): <http://www.ncbi.nlm.nih.gov/sites/genome> și să se salveze secvența cu numărul de acces NC_014378 în format GenBank. Pentru afișarea întregii secvențe în fișierul GenBank, înainte de salvarea fișierului, să se selecteze opțiunea „Show sequence” de la „Customize view”.
2. Parsarea GenBank

Să se implementeze funcția *extractDNA(fin, fout)*, care citește secvența de ADN dintr-un fișier în format GenBank ('*sequence.gb*') și o scrie apoi într-un fișier în format FASTA ('*sequence.adn*').
3. Pe baza scheletului de cod din fișierul *schelet_genbank.py*, să se extragă din secvența de ADN toate genele indicate în fișierul GenBank și să se scrie câte una pe linie în fișierul '*sequence.genes*'. Trebuie tratate toate cazurile posibile ale poziției unei gene: complement, join, etc. Să se utilizeze fișierul FASTA creat la 2.

Trebuie completate următoarele funcții:

- *join(str, fADN)* – care tratează cazul în care o genă este alcătuită din mai multe subsecvențe din ADN și întoarce gena în urma concatenării secvențelor care o compun.
- *extractGeneFromFile(a,b,fADN)* – care extrage din fișierul cu numele *fADN* (format FASTA) secvența de ADN cuprinsă între poziția *a* și poziția *b* și întoarce gena ca un string.

Observații:

Să se utilizeze funcțiile *split()* și *strip()*:

- *sir.split(',')* – care întoarce o listă cu subșirurile din șirul *sir* separate prin separatorul ','
- *sir.strip()* – care elimină toate spațiile, taburile, caracterele '\n' de la capetele șirului *sir*
- *sir.lstrip()* – care elimină toate spațiile de la stânga șirului
- *sir.rstrip()* – care elimină toate spațiile de la dreapta șirului
- *sir.strip(chars)* – elimină de la capetele șirului *sir* toate caracterele conținute în șirul *chars*
- *sir.lstrip(chars)* – elimină din stânga șirului *sir* toate caracterele conținute în șirul *chars*
- *sir.rstrip(chars)* – elimină din dreapta șirului *sir* toate caracterele conținute în șirul *chars*
- *sir.replace('c1','c2')* – înlocuiește caracterul *c1* cu *c2* din șirul *sir*
- *sir.upper()* – transformă literele din șirul *sir* în litere mari
- *sir.find('subs')* – întoarce indicele de început al subșirului *subs* în șirul *sir*