

# Laborator 2 – Bioinformatică

## I. Simularea procesului de transcriptie din ADN in ARN si pregatirea procesului de translatie

a. Să se implementeze o funcție  $ADN(n)$  care să genereze o secvență aleatoare de ADN ca o succesiune de baze azotate *Adenina*, *Guanina*, *Timina*, *Ciozina* sub forma unei liste de tupluri de câte 3 baze azotate. Lista va avea dimensiunea  $n$ , iar succesiunea de baze este obținută utilizând modulul *random*.

Apelul  $ADN(5)$  întoarce, de exemplu:

$[(A, T, C), (G, A, A), (C, A, A), (T, C, A), (A, G, C)]$

### Indicii:

- să se utilizeze funcțiile  $seed(x)$  și  $randint(a,b)$  din modulul *random*;
- $seed(x)$  trebuie apelată o singură dată pentru a iniția algoritmul de generare de numere aleatoare (pseudoaleatoare);
- pentru ca la fiecare apel al funcției  $ADN()$  să se genereze o altă secvență de numere aleatoare, este recomandat să se seteze ca *seed* timpul sistemului;
- timpul sistemului în secunde este întors de funcția  $time()$  care se găsește în modulul *time*;
- funcția  $seed(time())$  trebuie apelată o singură dată în funcția  $ADN()$ ;
- numerele aleatoare se generează cu funcția  $randint(a,b)$  care întoarce un numar aleator între  $a$  și  $b$ .

b. Să se scrie o funcție  $Afisare(secv)$  care primește ca parametru o secvență precum cea de la punctul a. și o transformă într-un șir de caractere eliminând parantezele și virgulele.

Pentru exemplul de mai sus,  $Afisare(secv)$  va întoarce: ATCGAACAATCAAGC

c. Să se implementeze o funcție  $ARN(secADN)$  care generează o secvență de ARN mesager, care copiaza secvența de ADN primită ca parametru. Secventa de ADN este primita în direcția 5'3'. Tratatii ambele cazuri, considerand ca secventa care se doreste a fi copiata se poate afla fie pe catena primita ca parametru, fie pe cea complementara. Secvența de ARN este memorată în direcția 5'3' tot într-o listă de tupluri de baze azotate precum cea de ADN. Un tuplu de 3 baze azotate reprezintă un codon care va fi tradus într-un aminoacid în procesul de translație care are loc la nivelul ribozomilor.

Observație: Bazele complementare care formează legături sunt A-T și G-C. ARN conține Adenina, Guanina și Citozina, iar Timina este înlocuită cu Uracil.

## II. Numărarea perechilor de baze A-T și C-G dintr-o secvență

1. Să se scrie o funcție pbADN(n) care primește ca parametru o secvență de ADN sub forma unui șir de caractere. Funcția trebuie să numere perechile de baze azotate A=T, respectiv C≡G, conținute în regiunea de ADN reprezentată prin secvența primită ca parametru. Funcția va întoarce un tuplu (at, cg) în care:

- at reprezintă numărul de perechi de baze azotate A=T,
- gt reprezintă numărul de perechi de baze azotate C≡G.

2. Sa se studieze si sa se testeze urmatorul cod. Este o implementare optima? Sa se caute in documentatie (<http://docs.python.org/library>) la Built-in Functions functia *filter*. O modalitate de a defini functii intr-o forma foarte comprimata este utilizarea expresiei lambda (<http://docs.python.org/reference/expressions.html#lambda>).

```
secv = 'AAATTCTCTGGTAGA'
def number2(secv):
    n = len(filter(lambda x: x == 'A' or x == 'T', secv))
    return (n, len(secv)-n)
```

## III. Simularea procesului de translație – traducerea unei secvențe de ADN într-o secvență proteică

1. Nucleotidele sunt codificate standard prin literele: A, T, C, G (în cazul ADN) și A, U, C, G (în cazul ARN). O nucleotidă conține o bază azotată, de aceea denumirea nucleotidelor pleacă de la baza azotată pe care o conține. Acizii nucleici reprezintă lanțuri de nucleotide.

Știind că un codon (trei nucleotide) codifică un aminoacid, să se explice de ce codul genetic standard poate fi exprimat atât sub formă de codoni din ARN, cât și sub formă de codoni din ADN.

Să se descrie pe scurt procesele prin care se face trecerea de la informația genetică din ADN la proteine (transcripția și translația): ADN → ARN (ARNm, ARNt, ARNr) → secvență proteică (catena polipeptidică – lanț de aminoacizi). Codul genetic este ilustrat în figurile următoare:

A doua bază							
T		C		A			
G							
Prima bază	T	TTT (Phe/F)	TCT (Ser/S)	TAT (Tyr/Y)	TGT (Cys/C)	T	
		Fenilalanină	Serină	Tirozină	Cisteină		
		TTC (Phe/F)	TCC (Ser/S)	TAC (Tyr/Y)	TGC (Cys/C)	C	
		Fenilalanină	Serină	Tirozină	Cisteină		
		TTA (Leu/L)	TCA (Ser/S)	TAA (Stop)	TGA (Stop)	A	
	C	Leucină	Serină				
		TTG (Leu/L)	TCG (Ser/S)	TAG (Stop)	TGG (Trp/W)	G	
		Leucină	Serină		Triptofan		
		CTT (Leu/L)	CCT (Pro/P)	CAT (His/H)	CGT (Arg/R)	T	
		Leucină	Prolină	Histidină	Arginină		
		CTC (Leu/L)	CCC (Pro/P)	CAC (His/H)	CGC (Arg/R)	C	
		Leucină	Prolină	Histidină	Arginină		
		CTA (Leu/L)	CCA (Pro/P)	CAA (Gln/Q)	CGA (Arg/R)	A	
		Leucină	Prolină	Glutamină	Arginină		
		CTG (Leu/L)	CCG (Pro/P)	CAG (Gln/Q)	CGG (Arg/R)	G	
		Leucină	Prolină	Glutamină	Arginină		
	A	ATT (Ile/I)	ACT (Thr/T)	AAT (Asn/N)	AGT (Ser/S)	T	
		Izoleucină	Treonină	Asparagină	Serină		
		ATC (Ile/I)	ACC (Thr/T)	AAC (Asn/N)	AGC (Ser/S)	C	
		Izoleucină	Treonină	Asparagină	Serină		
		ATA (Ile/I)	ACA (Thr/T)	AAA (Lys/K)	AGA (Arg/R)	A	
	G	Izoleucină	Treonină	Lizină	Arginină		
		ATG (Met/M)	ACG (Thr/T)	AAG (Lys/K)	AGG (Arg/R)	G	
		Metionină (Start)	Treonină	Lizină	Arginină		
		G	GTT (Val/V)	GCT (Ala/A)	GAT (Asp/D)	GGT (Gly/G)	T
			Valină	Alanină	Acid aspartic	Glicină	
	GTC (Val/V)		GCC (Ala/A)	GAC (Asp/D Acid aspartic	GGC (Gly/G)	C	
	Valină		Alanină		Glicină		
	GTA (Val/V)		GCA (Ala/A)	GAA (Glu/E)	GGA (Gly/G)	A	
	Valină		Alanină	Acid glutamic	Glicină		
		GTG (Val/V)	GCG (Ala/A)	GAG (Glu/E)	GGG (Gly/G)	G	
		Valină	Alanină	Acid glutamic	Glicină		

*Codonii ADN care codifică cei 20 de aminoacizi*

A doua bază								
U		C		A		G		
Prima bază	U	UUU (Phe/F)	UCU (Ser/S)	UAU (Tyr/Y)	UGU (Cys/C)	U	A treia bază	
		Fenilalanină	Serină	Tirozină	Cisteină			
		UUC (Phe/F)	UCC (Ser/S)	UAC (Tyr/Y)	UGC (Cys/C)	C		
		Fenilalanină	Serină	Tirozină	Cisteină			
		UUA (Leu/L)	UCA (Ser/S)	UAA (Stop)	UGA (Stop)	A		
	Leucină	Serină						
	UUG (Leu/L)	UCG (Ser/S)	UAG (Stop)	UGG (Trp/W)	G			
	Leucină	Serină		Triptofan				
	C	CUU (Leu/L)	CCU (Pro/P)	CAU (His/H)	CGU (Arg/R)	U		
		Leucină	Prolină	Histidină	Arginină			
		CUC (Leu/L)	CCC (Pro/P)	CAC (His/H)	CGC (Arg/R)	C		
		Leucină	Prolină	Histidină	Arginină			
		CUA (Leu/L)	CCA (Pro/P)	CAA (Gln/Q)	CGA (Arg/R)	A		
	Leucină	Prolină	Glutamină	Arginină				
	CUG (Leu/L)	CCG (Pro/P)	CAG (Gln/Q)	CGG (Arg/R)	G			
	Leucină	Prolină	Glutamină	Arginină				
	A	AUU (Ile/I)	ACU (Thr/T)	AAU (Asn/N)	AGU (Ser/S)	U		
		Izoleucină	Treonină	Asparagină	Serină			
		AUC (Ile/I)	ACC (Thr/T)	AAC (Asn/N)	AGC (Ser/S)	C		
		Izoleucină	Treonină	Asparagină	Serină			
		AUA (Ile/I)	ACA (Thr/T)	AAA (Lys/K)	AGA (Arg/R)	A		
	Izoleucină	Treonină	Lizină	Arginină				
	AUG (Met/M)	ACG (Thr/T)	AAG (Lys/K)	AGG (Arg/R)	G			
	Metionină (Start)	Treonină	Lizină	Arginină				
	G	GUU (Val/V)	GCU (Ala/A)	GAU (Asp/D)	GGU (Gly/G)	U		
		Valină	Alanină	Acid aspartic	Glicină			
		GUC (Val/V)	GCC (Ala/A)	GAC (Asp/D Acid aspartic	GGC (Gly/G)	C		
		Valină	Alanină		Glicină			
		GUA (Val/V)	GCA (Ala/A)	GAA (Glu/E)	GGA (Gly/G)	A		
	Valină	Alanină	Acid glutamic	Glicină				
	GUG (Val/V)	GCG (Ala/A)	GAG (Glu/E)	GGG (Gly/G)	G			
	Valină	Alanină	Acid glutamic	Glicină				

*Codonii ARN care codifică cei 20 de aminoacizi*

2. Să se genereze o secvență de ADN aleatoare sub forma unui șir de caractere sau a unei liste de tupluri (ca la 1.), având  $n=3*k$  nucleotide. Știind că un codon codifică un aminoacid, să se construiască secvența proteică corespunzătoare secvenței de ADN generate, sub forma unui șir de caractere în care aminoacizii sunt separati printr-un spatiu. Să se afișeze înșiruirea de aminoacizi în două moduri:

a. utilizând denumirile lor prescurtate din codul genetic standard:

```
stdcode = {
'TTT':'Phe', 'TTC':'Phe', 'TTA':'Leu', 'TTG':'Leu',
'TCT':'Ser', 'TCC':'Ser', 'TCA':'Ser', 'TCG':'Ser',
'TAT':'Tyr', 'TAC':'Tyr', 'TAA':'Stop', 'TAG':'Stop',
'TGT':'Cys', 'TGC':'Cys', 'TGA':'Stop', 'TGG':'Trp',

'CTT':'Leu', 'CTC':'Leu', 'CTA':'Leu', 'CTG':'Leu',
'CCT':'Pro', 'CCC':'Pro', 'CCA':'Pro', 'CCG':'Pro',
'CAT':'His', 'CAC':'His', 'CAA':'Gln', 'CAG':'Gln',
'CGT':'Arg', 'CGC':'Arg', 'CGA':'Arg', 'CGG':'Arg',

'ATT':'Ile', 'ATC':'Ile', 'ATA':'Ile', 'ATG':'Met',
'ACT':'Thr', 'ACC':'Thr', 'ACA':'Thr', 'ACG':'Thr',
'AAT':'Asn', 'AAC':'Asn', 'AAA':'Lys', 'AAG':'Lys',
'AGT':'Ser', 'AGC':'Ser', 'AGA':'Arg', 'AGG':'Arg',

'GTT':'Val', 'GTC':'Val', 'GTA':'Val', 'GTG':'Val',
'GCT':'Ala', 'GCC':'Ala', 'GCA':'Ala', 'GCG':'Ala',
'GAT':'Asp', 'GAC':'Asp', 'GAA':'Glu', 'GAG':'Glu',
'GGT':'Gly', 'GGC':'Gly', 'GGA':'Gly', 'GGG':'Gly'
}
```

b. utilizând codificarea standard de o literă pentru aminoacizi – în acest caz, dacă în secvența de ADN se generează codoni de stop care nu codifică aminoacizi, se va afișa secvența proteică până la codonul de Stop.

```
aminocode = {
'Gly' : 'G', 'Ala' : 'A', 'Pro' : 'P', 'Val' : 'V',
'Ile' : 'I', 'Leu' : 'L', 'Phe' : 'F', 'Met' : 'M',
'Ser' : 'S', 'Cys' : 'C', 'Thr' : 'T', 'Asn' : 'N',
'Gln' : 'Q', 'His' : 'H', 'Tyr' : 'Y', 'Trp' : 'W',
'Asp' : 'D', 'Glu' : 'E', 'Lys' : 'K', 'Arg' : 'R'
}
```

**Observație.** Se poate utiliza funcția `sir.split(' ')` care întoarce o listă cu subșirurile din șir separate de un separator dat ca parametru (în acest caz, separatorul este un spațiu).

3. Se da o secvență de ARN sub forma unei liste de tupluri generata aleator. Să se construiască secvența proteică corespunzătoare. Un tuplu va contine 3 simboluri, fiind corespunzator unui codon ARN alcatuit din 3 nucleotide.

$[(A, U, C), (G, A, A), (C, A, A), (U, C, A), (A, G, C)]$

## IV. Formate de fișiere pentru secvențe și 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, alcătuită 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*) pentru un genom cromozomial al unui organism, de exemplu: NC\_005562 .

d. Să se acceseze baza de date *RefSeq*: NC\_005562.

e. Să se salveze informația găsită, în formatele: FASTA, ASN. 1, GenBank.

f. Să se salveze doua secvențe genomice diferite (unul procariot și unul eucariot) în fiecare din cele trei formate.