

## I.5

### BAZA DE DATE INTEGRATĂ NCBI. PROGRAMUL BLAST

În această material ne vom familiariza cu programul BLAST. Prin intermediul lui vom face căutări în bazele de date după secvențe de nucleotide și aminoacizi.

#### 1. Introducere

Programul BLAST (**B**asic **L**ocal **A**lignment **S**earch **T**ool) este un program popular de comparare a secvențelor de ADN, ce face parte dintr-un pachet de programe destinat căutării de secvențe proteice, accesibil în diverse forme la diferiți furnizori, sau prin intermediul NCBI, care mai oferă și Entrez, un instrument de meta-căutare care acoperă mare parte a bazelor de date de la NCBI, inclusiv cele care găzduiesc structuri tridimensionale a proteinelor, genoamele complete ale organismelor și trimiteri la jurnale științifice care însoțesc intrările din bazele de date.

Asocierea dezvoltărilor tehnologiei de calcul și moleculare deschide noi oportunități cercetărilor genetice. Folosirea combinată a informației oferită de secvențe, a instrumentelor de calcul, a bazelor de date și a biologiei tradiționale crește speranța înțelegerii funcției și reglajelor tuturor genelor și proteinelor, precum și a descifrării funcțiilor celulei.

#### 2. Prezentarea programului BLAST

BLAST reprezintă instrumentul de căutare a alinamentului local de bază, fiind un set de programe de căutare a similarităților, creat pentru identificarea clasificării și a omologilor potențiali pentru o secvență dată.

Pentru a înțelege mai bine programele BLAST, trebuie cunoscute aspectele de bază ale aliniamentelor secvențelor. Acestea sunt folosite în special pentru găsirea potențialilor omologi ce vor fi folosiți ulterior pentru prezicerea posibilelor funcții ale secvenței necunoscute sau pentru modelarea structurii sale tridimensionale.

Programele BLAST folosesc un algoritm heuristic care identifică aliniamentele locale, găsind omologii cu secvențele cele mai apropiate, într-un timp eficient.

Serverul BLAST suportă o varietate de programe analitice care sunt fie accesate prin rețeaua Internet, fie instalate în rețele locale pentru a mări viteza de analiză. Programul BLAST bazal nu permite introducerea gap-urilor în aliniamentele sale ceea ce va reduce sensibilitatea căutării. Cu toate acestea, datele de ieșire din program oferă aliniamente regionale multiple, care pot fi folosite pentru a anticipa gap-urile din secvența de interes și cea din baza de date. În continuare sunt enumerate programele BLAST și utilizarea lor.

- a) **BLASTp**: acest program permite utilizatorului să caute similaritățile dintre secvența unei proteine necunoscute și secvențele proteinelor dintr-o bază de date.
- b) **BLASTx**: permite compararea secvențelor traduse în aminoacizi ale nucleotidelor cu secvențele proteinelor din bazele de date.  
Secvența nucleotidică de interes este tradusă inițial în toate cele 6 catene de citire **ORF** (**O**pen **R**eadin**G** **F**rame) posibile. Acest program este folosit în special pentru găsirea erorilor de secvențializare a nucleotidelor, prin compararea secvenței de nucleotide tradusă în aminoacizii săi proteici potențiali dintr-o bază de date cu secvențe proteice.
- c) **BLASTn**: cu ajutorul acestui program se compară o secvență nucleotidică de interes cu secvențele din bazele de date nucleotidice.
- d) **tBLASTn**: permite căutarea similarităților dintre o secvență proteică și secvențele traduse (translate) ale nucleotidelor dintr-o bază de date.

Secvențele nucleotidice dintr-o bază de date sunt traduse inițial în fiecare din cele 6 catene de citire posibile și sunt apoi comparate cu secvența proteinei de interes. Acest program este util pentru găsirea erorilor de secvențializare în proteine prin compararea secvenței proteinei respective cu omologii săi potențiali obținuți prin traducerea secvențelor nucleotidice dintr-o bază de date.

La adresa <https://blast.ncbi.nlm.nih.gov/Blast.cgi> se găsesc opțiunile programului BLAST:

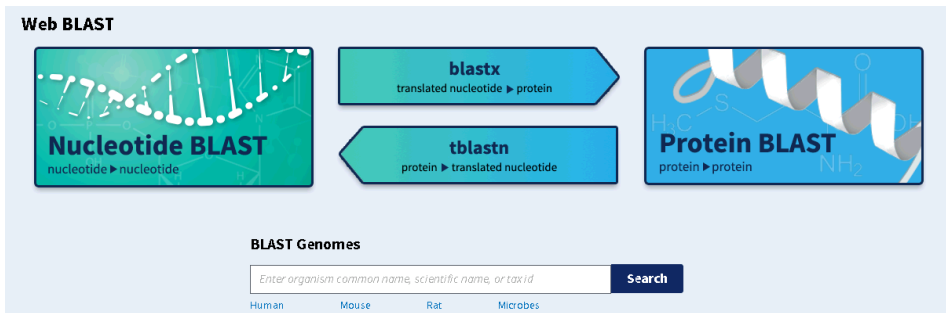


Figura 1. Interfața BLAST

### 3. Exemplu de căutare

Dacă dorim să căutăm secvența de gene care codează receptorul pentru endotelină la om (figura 2) introducem în căsuța search următoarele: **Endothelin receptor**.

Din mulțimea de răspunsuri vom selecta varianta receptorului pentru specia umană (homo sapiens). Va fi afișat cromozomul pe care se află secvența de nucleotide (cromozomul 40) – figura 3, localizarea genei pe cromozom, secvența de nucleotide și secvența codantă.

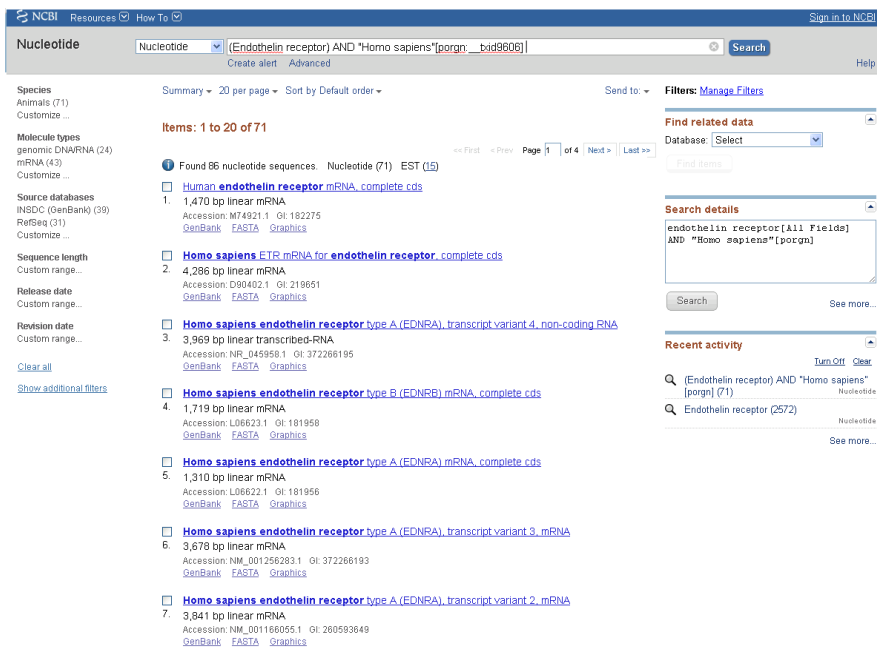


Figura 2. Exemplu de căutare a unei secvențe de nucleotide

```

FEATURES             Location/Qualifiers
     source            1..3841
                       /organism="Homo sapiens"
                       /mol_type="mRNA"
                       /db_xref="taxon:9606"
                       /chromosome="4"
                       /map="4q31.22-q31.23"
     gene              1..3841
                       /gene="EDNRA"
                       /gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
                       /note="endothelin receptor type A"
                       /db_xref="GeneID:1909"
                       /db_xref="HGNC:HGNC:3179"
                       /db_xref="MIM:131243"
     exon              1..460
                       /gene="EDNRA"
                       /gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
                       /inference="alignment:Splign:1.39.8"
     misc feature      30
                       /gene="EDNRA"
                       /gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
                       /note="alternative transcription initiation site"
     STS               46..236
                       /gene="EDNRA"
                       /gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
                       /standard_name="SHGC-67921"
                       /db_xref="UniSTS:38684"
     exon              461..950
                       /gene="EDNRA"
                       /gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
                       /inference="alignment:Splign:1.39.8"
     misc feature      501..503
                       /gene="EDNRA"
                       /gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
                       /note="upstream in-frame stop codon"
     CDS               531..1487
                       /gene="EDNRA"
                       /gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
                       /note="isoform b precursor is encoded by transcript
variant 2; endothelin receptor subtype A;

```

**Datele despre secvența de nucleotide  
cromozomul, bratul cromozomului, gena**

Figura 3. Rezultatele căutării

Ca secvență de nucleotide și aminoacizi, găsim următoarea formulare (fig. 4).

```

444 agtctgaggg agtctgaggg agtctgaggg agtctgaggg agtctgaggg
181 agtctctccg ctgctctgac gattctggag agtctgagga gaggctcat ccactccacc
241 cgtctctcgc cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
301 ggtttcttga agtctgagga gattctggag agtctgagga gaggctcat ccactccacc
361 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
421 ggtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
481 taagaagcgc aagaagcgc taagaagcgc taagaagcgc taagaagcgc
541 ttgtctgagc agtctgagga gattctggag agtctgagga gaggctcat ccactccacc
601 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
661 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
721 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
781 ctgtgagcgc agtctgagga gattctggag agtctgagga gaggctcat ccactccacc
841 taagaagcgc aagaagcgc taagaagcgc taagaagcgc taagaagcgc
901 ttgtctgagc agtctgagga gattctggag agtctgagga gaggctcat ccactccacc
961 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1021 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1081 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1141 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1201 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1261 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1321 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1381 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1441 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1501 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1561 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1621 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1681 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1741 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1801 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1861 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1921 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
1981 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
2041 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc
2101 agtctgagga cgtctctgac gattctggag agtctgagga gaggctcat ccactccacc

```

**secvența de nucleotide**

```

531..1487
/gene="EDNRA"
/gene_synonym="ET-A; ETÀ; ETÀ-R; ETAR; ETRA; hET-AR; MFDA"
/note="isoform b precursor is encoded by transcript
variant 2; endothelin receptor subtype A;
endothelin-1-specific receptor; 6 protein-coupled
receptor"
/codon_start=1
/product="endothelin-1 receptor isoform b precursor"
/protein_id="NP_001159527.1"
/db_xref="CCDS:CCDS4818.1"
/db_xref="GeneID:1909"
/db_xref="HGNC:HGNC:3179"
/db_xref="MIM:131243"
/translation="METLCLASPLALVGVISDIPERYSTLNSHVDQFTFRGTE
LSPFLVTHPTNLVPSMSHMYCQQTITSAFYNTIVISCTFIWVWVHALL
RIZVQWCHRMGPHALSLALGDLVYVZLPTNFTFYQVQVWALPQVYCPVLV
CTKTYVITUTCEINRRNSLRLALSEHLKQREKAVKTVCLVFLVAFWPLHLSR
LKKTVVNIQNKIRCLLFLLLDYIDYDILATNSCHPLALVYVFKFKFCQSLC
CCYQSKSLVTPVHSTQIKNQINQINRSHSKDSIN"

```

**secvența de aminoacizi**

Figura 4. Formatul BLAST de vizualizare

## 4. Exerciții practice

NOTA: pentru secvențele-text accesați fișierul *Bioinfo-2018-practic.txt*.



### Exercițiul 1

National Center for Biotechnology Information:

<https://www.ncbi.nlm.nih.gov/>

BLAST - Basic Local Alignment Search Tool

Nucleotide BLAST - Enter Query Sequence:

```
1 ctgaaaccg tatgtatat aattatatac tataaagtaa taatgtatac agtgaatgg
61 atcatgggcc atgtgctttt caaactaatt gtacataaaa caagcatcta ttgaaaatat
121 ctgacaaact catcttttat ttttgatgtg tgtgtgtgtg tgtgtgtgtg ttttttaac
181 aggggatttgg gg
```

(căutarea este comandată cu butonul BLAST)

Observați scorul de potrivire - investigați/căutați informațiile referitoare la gena *Homo sapiens Human cystic fibrosis transmembrane conductance regulator (CFTR)*, Locus AH002646.

Observați adnotările/referințele bibliografice – coloana din dreapta.



### Exercițiul 2

Căutați proteina *insulina*.



### Soluții

National Center for Biotechnology Inform ation: <https://www.ncbi.nlm.nih.gov/>

Protein - Search for "insulin" - alegeți "insulin[*Homo sapiens*]"

Observați adnotările/referințele bibliografice – coloana din dreapta

Observați secvența de aminoacizi ai proteinei – ORIGIN:

```
1 malwmrllpl lallalwgp d paaafvnqhl cgshlvealy lvcgergffy tpkttreaed
61 lqvqqvelgg gpgagslqpl alegslqkrg iveqcctsic slyqlenycn
```

```
Site      order(32,35,37,108)
          /site_type="other"
          /note="putative receptor binding surface"
          /db_xref="CDD:239833"

CDS      1..110
          /gene="INS"
          /coded_by="join(AH002844.2:2424..2610,
          AH002844.2:3397..3542)"
          /note="precursor"
          /db_xref="GDB:600-119-349"

ORIGIN
1 malwmrllpl lallalwgpd paaafvnqhl cgshlvealy lvcgergffy tpktreaed
61 lqvqvelgg gpgagslqpl alegslqkrg iveqcctsic slyqlenycn
//
```

Figura 5. Secvența de aminoacizi ai proteinei insulina umană



### Exercițiul 3

Căutați după secvența de aminoacizi **malwmrllpl**.



### Soluții

Figura 6. Interogarea BD după o secvență de aminoacizi

Observați scorul de potrivire....



### Exercițiul 4

Căutați secvența de aminoacizi:

1 mqnqagasrt stflngnre rplnvfcdme tdgggwlvfq rrmkgqtdfw rdwedyahgf  
 61 gnisgefwlq nealhsltqa gdsyirvdlr agdeavfaqy dsfhvdsaae yyrlhlegyh  
 121 gtagdmsyhgsgsvfsardr dpnslisla vsyrgawwyr nchyanlngl ygstvdhqqv  
 181 swyhwkgfef svpftemklr prnfrspagg g

Și vizualizați structura 3D.



### Soluții

National Center for Biotechnology Information:

<https://www.ncbi.nlm.nih.gov/>

Protein BLAST - Enter Query Sequence....(copy/paste secvență)

Pentru *fibrinogen* [*Homo sapiens*] accesați structura 3D.

The screenshot displays the NCBI Protein BLAST interface. The search query is 'Fibrinogen (Homo sapiens)'. The results section shows a list of protein entries, including 'Fibrinogen (Homo sapiens)' and 'Fibrinogen (Homo sapiens)'. The 3D structure visualization of the protein is shown on the right side of the interface, with a red arrow pointing to it.

Figura 7. Vizualizarea 3D a structurii fibrinogenului uman



### Exercițiul 5

Căutați informații avansate despre fibrinogenul uman.



### Soluții

National Center for Biotechnology Information:

<https://www.ncbi.nlm.nih.gov/>

Protein - Advanced Search for "human fibrinogen".

Protein Advanced Search Builder

(fibrinogen[Protein Name]) AND human[Organism]

Edit Clear

Builder

Protein Name: fibrinogen Show index list

AND Organism: human Show index list

AND All Fields Show index list

Search or Add to history

History

Search	Add to builder	Query	Items found	Time
#9	Add	Search fibrinogen	112309	14:11:24
#8	Add	Search insulin	83725	13:47:56

You are here: NCBI > Proteins > Protein Database

Support Center

Figura 8. Filtru avansat de căutare

Căutați informații și pe Wikipedia:

<https://en.wikipedia.org/wiki/Fibrinogen>



### Exercițiul 6

Căutați secvența:

1 mkwvtfisll flfssaysrg vfrrdahkse vahrfkdige enfkalvlia faqylqqepf  
61 edhvklynev tefaktcvad esaencdksl htlfgdklct vatltretyge madccakqep  
121 ernecflqhk ddnpnlprlv rpevdvmcta fhdneetflk kylieiarrh pyfyapellf  
181 fakrykaaft eccqaadkaa cllpkldeir degkassakq rlkcaslqkf gerafkawav



241 arlsqrpfka efaevsklvt dltkvhtec hgdllcadd radlakyice nqdsissklk  
301 eccekppllek shciaevend empadlpsla adfveskdvc knyaeakdvf lgmflyeyar  
361 rhpdyssvll lrlaktyett lekccaaadp hecyakvfde fkplveepqn likqncelfe  
421 qlgeyqfna llvrytkvp qvstptlvev srnlgkvgs cckhpeakrm pcaedylsvv  
481 lnqlcvlhek tpsdrvtkc cteslvnrrp cfsalevdet yvpkefnaet ftfhadietl  
541 sekerqikkq talvelvkhk pkatkeqlka vmddfaafve kcckaddket cfaeegkklv  
601 aasqaalgl

Observați și analizați rezultatele.

Accesați opțiunea "Analyze your query with SmartBLAST".

## 5. Concluzii

Din acest material am învățat să facem:

- + căutări în bazele de date după secvențe de gene;
- + determinăm cromozomul pe care se află o secvență de nucleotide;
- + căutări în bazele de date după secvențe de nucleotide;
- + căutare în BD după secvența de aminoacizi;
- + găsim gradul de potrivire al secvenței de aminoacizi introdusă.

## Referințe

Gheorghe-Ioan Mihalas, Anca Tudor, Sorin Paralescu – Bioinformatica, Colecția "Științele exacte în cercetarea medicală", Timisoara: Editura Victor Babes, 2011, ISBN 978-606-8054-33-9  
<https://www.ncbi.nlm.nih.gov/>