# Bioinformatika és genomanalízis az orvostudományban

## Keresés adatbázisokban

Cserző Miklós

2018

# A mai előadás

- Szekvenciakeresés:
  - FASTA
  - BLAST család
  - HMMER
- Párhuzamos változatok a keresésre
- A keresési eredmények értékelése
- Szöveges keresés: eTBLAST
- Adatbányászat

# Az adatbázis keresés jelentősége

➢ Össze tudunk hasonlítani két szekvenciát

➢ El tudjuk dönteni, hogy közös őstől származnak-e

➢ Tudunk illeszteni több, hasonló szekvenciát

➢ De hogyan találjuk meg a hasonlókat több millió szekvencia közt???

# Mit keresünk?

- ➢ Keresés az annotációban
    - ➢ Kódok
    - ➢ Kulcsszavak
    - ➢ Azonosított funkciók
- ➢ Ez alapkövetelmény minden adatbázissal szemben
- ➢ Keresés a szekvenciában
- ➢ Ez a biológiai adatbázisok jellegzetessége

# Hogy lehet gyorsan keresni?

➢ A szekvenciát feltördeljük rövid, összefüggő karakter-sorozatokra – *szavakra* (*wordsize, k-mer, k-tuple*)

➢ Az adatbázist indexeljük: elkészítjük a *k-merek* listáját és hozzárendeljük a helyüket

➢ Ezt csak egyszer kell megcsinálni

➢ A kereső szekvenciát is felbontjuk *k-merekre*

➢ És így keresünk az adatbázisban

GTCTGACAGCAGCCGCTGCAGCAGCTGCTGCTGCTACCAATGCAG
CTATTGCTGAAGCAA

GTC:1

GTCTGACAGCAGCCGCTGCAGCAGCTGCTGCTGCTACCAATGCAG
CTATTGCTGAAGCAA

GTC:1    TCT:2

GTCTGACAGCAGCCGCTGCAGCAGCTGCTGCTGCTACCAATGCAG
CTATTGCTGAAGCAA

GTC:1    TCT:2    CTG:3

GTCTGACAGCAGCCGCTGCAGCAGCTGCTGCTGCTACCAATGCAG
CTATTGCTGAAGCAA


 Táblázatos forma (lookup table):


AAG:55, AAT:39, ACA:6, ACC:36, AGC:8:11:20:23:44:56, ATG:40,
ATT:48, CAA:38:58, CAG:7:10:19:22:43, CCA:37, CCG:13, CGC:14,
CTA:34:46, CTG:3:16:25:28:31:52, GAA:54, GAC:5, GCA:9:18:21:42:57,
GCC:12, GCT:15:24:27:30:33:45:51, GTC:1, TAC:35, TAT:47, TCT:2,
TGA:4:53, TGC:17:26:29:32:41:50, TTG:49

# A FASTA család

➤ Honlap: http://fasta.bioch.virginia.edu/

➤ Szolgáltatás itt és az EBI oldalán is

➤ Letölthető Linuxra és Windowsra is

➤ Dokumentáció és tutorial is elérhető

➤ Az egyik legnépszerűbb kereső

➤ Az algoritmus régi, de azóta is folyamatosan fejlesztik a programot

# Az algoritmus négy lépése

a) Gyors keresés táblázat alapján (lookup table)

b) Pontozás mátrix alapján

c) Kiválasztás

d) S-W illesztés, végeredmény

**FASTA Algorithm**

(a) Sequence B / Sequence A
Find runs of identities

(b) Sequence B / Sequence A
Re-score using PAM matrix
Keep top scoring segments.

(c) Sequence B / Sequence A
Apply "joining threshold" to eliminate segments that are unlikely to be part of the alignment that includes highest scoring segment.

(d) Sequence B / Sequence A
Use dynamic programming to optimise the alignment in a narrow band that encompasses the top scoring segments.

# A programcsalád tagjai

- ➢ A legfrissebb verzió a 3.6-os
    - ➢ „fasta": egy szekvenciát összehasonlít egy databázis szekvenciái ellenében (fehérjét fehérjével vagy DNS-t DNS-sel), gyors algoritmust használ – táblázatos keresés
    - ➢ „ssearch": S-W algoritmust használ a keresésre (fehérjét fehérjével vagy DNS-t DNS-sel), lassabb, de pontosabb

# Folytatás …

- ➢ „ggsearch": global:global keresés (fehérje és DNS)
- ➢ „glsearch": global:local keresés (fehérje és DNS)
- ➢ „fastx": lefordított DNS szekvenciát keres fehérje adatbázis ellen, három frame-et fordít, toldás és frame-eltolás is megengedett
- ➢ „fasty": mint a „fastx", de kódonon belüli eltolás is engedett

# Folytatás ...

➢ „tfastx", „tfasty": fehérje szekvenciát keres DNS adatbázis ellen

➢ „fastf", „tfastf": peptid keverék listáját keresi fehérje, illetve DNS adatbázis ellen, mintha részlegesen emésztett minta lenne

➢ „fasts", „tfasts": peptid fragmens listát keres fehérje, illetve DNS adatbázis ellen, mintha tömegspektrométerből származó minta lenne

➢ „lalign": többszörös illesztő program

# A FASTA felület – EBI

# Paraméterek

File  Edit  View  History  Bookmarks  Tools  Help

UVA FASTA Downloads

fasta.bioch.virginia.edu/fasta_www2/fasta_class.shtml

metabolic pathways poster

Most Visited  sajto  SOTE  Logins  Tool  DAS  IT  Biosites  Library  Athénba mentem  Post-Card-iff  post-card-iff

- FASTS/FASTF

**Software**

- FASTA v36 ChangeLog
- Downloads
- Sequence Libraries
- Developer Mailing list

**Other resources**

- CHAPS - Convert HMMs and Profiles
- Near optimal alignments
- FASTA Exercises
- NCBI BLAST server
- EMBL-EBI Server

Most of the searches in this exercise should be done against a small protein database, e.g. the **PIR1** database available at the FASTA WWW site. Searching a small database makes it practical to consider each of the high scoring similarities, and to evaluate further whether they are likely to be biologically meaningful.

---

**Identifying homologs and non-homologs; effects of scoring matrices and algorithms**

**1.** Use the FASTA search page to compare Drosophila glutathione transferase GSTT1_DROME (gi|121694) to the PIR1 Annotated protein sequence database.

a. What is the highest scoring non-homolog? (How would you confirm that your candidate non-homolog was truly unrelated?)

b. Note that this drosophila glutathione transferase shares significant similarity with both sequences from bacteria (SSPA_SHIFL, stringent starvation protein) and mammals. How might you test whether the stringent starvation protein is homologous to glutathione transferases? (*Hint - search* **SwissProt** *for a more comprehensive view of the family*)

c. Compare the expectation (E()) value for the distant relationship between GSTT1_DROME and GSTM2_RAT (class-mu). How would you demonstrate that GSTT1_DROME is homologous to GSTM2_RAT?

d. Examine how the expectation value changes with different scoring matrices (BLOSUM62, BlastP62, PAM250) and different gap penalties. (The default scoring matrix for the FASTA programs is BLOSUM50, with gap penalties of -10 to open a gap and -2 for each residue in the gap - e.g. -12 for a one residue gap).

What happens to the E()-value for the highest scoring unrelated sequence with the different matrices?

Look at the distribution of scores and the E()-value of the highest scoring unrelated sequence when the gap-open/gap-ext penalties are small (-7/-1).

e. Try the search with **ssearch** (Smith-Waterman). Again, look at the E()-values for distant homologs and the highest scoring unrelated sequence.

f. (*optional*) Try the search with *ktup=1* (What is ktup?). **FASTA** uses the *ktup* parameter to adjust the sensitivity and speed of the search. With *ktup=2*, **FASTA** looks for "pairs" of matched identical residues to find regions of similarity. *ktup=1* looks for singly-aligned residues, and thus takes longer.

---

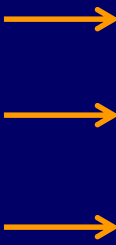**2.** Do the same search (121694) using the Course BLAST WWW page.

File Edit View History Bookmarks Tools Help

UVA FASTA Downloads | RecName: Full=Glutathione S-transfe... | FASTA results

fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi

metabolic pathways poster

Most Visited | sajto | SOTE | Logins | Tool | DAS | IT | Biosites | Library | Athénba mentem | Post-Card-iff | post-card-iff

Search Databases with FASTA | Find Duplications | Search Status

```
# fasta36 -p -q -w 80 -m 9i -m 6 -H -f -10 -S -g -2 TMP.q A 2

FASTA searches a protein or DNA sequence data bank
 version 36.3.6 Sep, 2012(preload9)
Please cite:
 W.R. Pearson & D.J. Lipman PNAS (1988) 85:2444-2448

Query: TMP.q
  1>>>FASTA exercise - 209 aa
Library: PIR1 Annotated (rel. 66)
  5190221 residues in 13351 sequences

         opt      E()
< 40      8     0:===
  42      2     0:=            one = represents 3 library sequences
  44     10     1:*===
  46     13     5:=*===
  48     40    15:====*=========
  50     54    37:============*=====
  52     89    70:======================*======
  54     69   107:======================           *
  56    142   142:=======================================*
  58    170   166:=================================================*=
  60    179   177:==================================================*=
  62    155   176:=================================            *
  64    129   165:=====================================         *
  66    141   148:=====================================   *
  68    131   129:==========================================*=
  70     77   109:=========================           *
  72     88    90:=============================*
  74     69    74:======================= *
  76     79    59:=====================*=======
  78     47    47:================*
  80     61    38:=============*========
  82     27    30:=========*
  84     22    23:=======*
  86     15    18:=====*
  88     16    14:====*=
  90     14    11:===*=
  92      8     9:==*
  94     14     7:==*==
  96      6     5:=*
  98      2     4:=*
 100      4     3:*=
 102      0     2:*
 104      3     2:*
 106      0     1:*
 108      2     1:*            inset = represents 1 library sequences
 110      1     1:*
 112      1     1:*        :*
 114      0     1:*        :*
 116      1     0:=        *=
 118      0     0:          *
```

File Edit View History Bookmarks Tools Help

UVA FASTA Downloads | RecName: Full=Glutathione S-transfe... | FASTA results | +

fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi | metabolic pathways poster

Most Visited | sajto | SOTE | Logins | Tool | DAS | IT | Biosites | Library | Athénba mentem | Post-Card-iff | post-card-iff

```
sp|P21163|PNGF_ELIMR Peptide-N(4)-(N-acetyl-beta-D-gluc ( 354)   80 27.9      4 0.307 0.557   88 align
sp|P23400|TRXM_CHLRE Thioredoxin M-type, chloroplast pr ( 140)   71 25.9    6.1 0.274 0.524   84 align
sp|P01577|IFNB3_BOVIN Interferon beta-3 precursor        ( 186)   72 26.1    7.4 0.345 0.509   55 align
sp|P17472|VGLB_EHV4 Glycoprotein B precursor             ( 919)   83 28.3    7.9 0.269 0.551   78 align


>>>FASTA, 209 aa vs A library

>>sp|P20432|GSTT1_DROME Glutathione S-transferase 1-1 (GST class-t        (209 aa)
 initn: 1399 init1: 1399 opt: 1399  Z-score: 1964.9  bits: 370.6 E(13351): 1.6e-103
Smith-Waterman score: 1399; 100.0% identity (100.0% similar) in 209 aa overlap (1-209:1-209)
Entrez Lookup   Re-search database   General re-search
              10        20        30        40        50        60        70        80
FASTA    MVDFYYLPGSSPCRSVIMTAKAVGVELNKKLLNLQAGEHLKPEFLKINPQHTIPTLVDNGFALWESRAIQVYLVEKYGKT
         ::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
sp|P20  MVDFYYLPGSSPCRSVIMTAKAVGVELNKKLLNLQAGEHLKPEFLKINPQHTIPTLVDNGFALWESRAIQVYLVEKYGKT
              10        20        30        40        50        60        70        80


              90       100       110       120       130       140       150       160
FASTA    DSLYPKCPKKRAVINQRLYFDMGTLYQSFANYYYPQVFAKAPADPEAFKKIEAAFEFLNTFLEGQDYAAGDSLTVADIAL
         ::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::::
sp|P20  DSLYPKCPKKRAVINQRLYFDMGTLYQSFANYYYPQVFAKAPADPEAFKKIEAAFEFLNTFLEGQDYAAGDSLTVADIAL
              90       100       110       120       130       140       150       160


             170       180       190       200
FASTA    VATVSTFEVAKFEISKYANVNRWYENAKKVTPGWEENWAGCLEFKKYFE
         :::::::::::::::::::::::::::::::::::::::::::::::::::
sp|P20  VATVSTFEVAKFEISKYANVNRWYENAKKVTPGWEENWAGCLEFKKYFE
             170       180       190       200



>>sp|P04907|GSTF3_MAIZE Glutathione S-transferase III (GST-III) (        (222 aa)
 initn: 182 init1: 142 opt: 183  Z-score: 258.0  bits: 54.8 E(13351): 1.9e-08
Smith-Waterman score: 183; 26.4% identity (55.7% similar) in 212 aa overlap (4-199:6-210)
Entrez Lookup   Re-search database   General re-search
              10        20        30        40        50        60        70
FASTA      MVDFYYLPGSSPCRSVIMTAKAVGVELNKKLLNLQAGEHLKPEFLKINPQHTIPTLVDNGFALWESRAIQVYLVEKYG
           .: .:  :  . :....  .: .: .:.:: .::   ..:::::: :... :.. ::.
sp|P04  MAPLKLYGMPLSPNVVRVATVLNEKGLDFEIVPVDLTTGAHKQPDFLALNPFGQIPALVDGDEVLFESRAINRYIASKYA
              10        20        30        40        50        60        70        80


        80        90       100       110                 120       130       140
FASTA   K--TDSLYPKCPKKRAVINQRLYFDMGTLYQSFANYYYPQVF--------AKAPADPEAFKKIEAAFEFLNTF---LEGQ
        .  :: :      :   .. ..... . :       : ::      . ::    .: .: . : . :...    :  .
sp|P04  SEGTDLL----PATASAAKLEVWLEVES--HHFHPNASPLVFQLLVRPLLGGAPDAAVVEKHAEQLAKVLDVYEAHLARN
                 90       100       110       120       130       140       150


       150       160       170       180       190       200
FASTA   DYAAGDSLTVADI--ALVATVSTFEVAKFE-ISKYANVNRWYENAKKVTPGWEENWAGCLEFKKYFE
        : ::: .:.::   ::.  ....  .  ..   .:. :.: :   . :.....  :
sp|P04  KYLAGDEFTLADANHALLPALTSARPPRPGCVAARPHVKAWWE-AIAARPAFQKTVAAIplppppsssA
           160       170       180       190       200       210       220


>>sp|P12653|GSTF1_MAIZE Glutathione S-transferase I (GST-I) (GST-2        (214 aa)
```

# A BLAST család

- Honlap: http://blast.ncbi.nlm.nih.gov/Blast.cgi
- Letülthető Linuxra és Windowsra is
- Web alapú szolgáltatás elérhető
- Alapos dokumentáció a honlapon
- Igen népszerű, megbízható
- Régóta van jelen az irodalomban, folyamatosan fejlesztik

# Az algoritmus

1. A szekvencia maszkolása
2. A szekvencia felbontása szavakra
3. A lista szűkítése a nagy pontértékű szavakra
4. Keresés az adatbázisban a lista alapján
5. A találatok kiterjesztése (HSP – high-scoring segment pair)
6. A HSP-k statisztikus kiértékelése
7. HSP-k összefűzése hosszabb illesztéssé

# A programcsomag tagjai

- „blastn": DNS szekvencia keresése DNS adatbázis ellen

- „blastp": fehérje szekvencia fehérje adatbázis ellen

- „psi-blast": fehérjék iteratív keresése fehérje adatbázis ellen

- „blastx": lefordított DNS szekvencia keresése fehérje adatbázis ellen

- „tblastx": lefordított DNS szekvencia keresés lefordított DNS adatbázison

- „tblastn": visszafordított fehérje keresés DNS adatbázis ellen

- „megablast": sok szekvencia keresése egy futás során

# Paraméterek

Statisztika

A program neve

Mátrix paraméterek

Eredmény

Maszkolás

STEP 3 - Set your parameters

PROGRAM
blastp

| MATRIX | GAP OPEN | GAP EXTEND | EXP. THR | FILTER |
|---|---|---|---|---|
| BLOSUM62 | 11 | 1 | 10 (default) | no |

| DROPOFF | SCORES | ALIGNMENTS | SEQUENCE RANGE | GAPALIGN |
|---|---|---|---|---|
| 0 (default) | 50 (default) | 50 (default) | START-END | true |

ALIGNMENT VIEWS
pairwise

STEP 4 - Submit your job

☐ Be notified by email *(Tick this box if you want to be notified by email when the results are available)*

Submit

# Egy további változat

BLAT:

- ➢ A cél a további gyorsítás
- ➢ Az ár: rosszabb érzékenység
  - ➢ Csak a nagyon hasonló szegmenseket talalálja meg: 95% egyezés DNS-re, 80% fehérjére
  - ➢ Rövid egyezéseket nem talál meg
- ➢ Új generációs szekvenálásnál igen hasznos

# Profilkeresés

➢ Az aminósavak helyettesítési hajlandósága függ a pozíciótól

➢ Egy jó többszörös illesztés megadja ezt az információt

➢ A többszörös illesztést közvetlenül „profillá" alakítjuk

➢ Ezt használjuk a kereséshez

➢ Megnő az eljárás érzékenysége

```
                                    *        .       :         .           .     *          :  : :        .
Q5E940_BOVIN   ----------------MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE    76
RLA0_HUMAN     ----------------MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE    76
RLA0_MOUSE     ----------------MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE    76
RLA0_RAT       ----------------MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE    76
RLA0_CHICK     ----------------MPREDRATWKSNYFMKIIQLLDDYPKCFVVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--PALE    76
RLA0_RANSY     ----------------MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQQIRMSLRGK-AVVLMGKNTMMRKAIRGHLENN--SALE    76
Q7ZUG3_BRARE   ----------------MPREDRATWKSNYFLKIIQLLDDYPKCFIVGADNVGSKQMQTIRLSLRGK-ZVVLMGKNTMMRKAIRGHLENN--PALE    76
RLA0_ICTPU     ----------------MPREDRATWKSNYFLKIIQLLNDYPKCFIVGADNVGSKQMQTIRLSLRGK-AIVLMGKNTMMRKAIRGHLENN--PALE    76
RLA0_DROME     ----------------MVRENKAAWKAQYFIKVVELFDEFPKCFIVGADNVGSKQMQNIRTSLRGL-AVVLMGKNTMMRKAIRGHLENN--PQLE    76
RLA0_DICDI     ----------------MSGAG-SKRKKLFIEKATKLFTTYDKMIVAEADFVGSSQLQKIRKSIRGI-GAVLMGKKTMIRKVIRDLADSK--PELD    75
Q54LP0_DICDI   ----------------MSGAG-SKRKNVFIEKATKLFTTYDKMIVAEADFVGSSQLQKIRKSIRGI-GAVLMGKKTMIRKVIRDLADSK--PELD    75
RLA0_PLAF8     ----------------MAKLSKQQKKQMYIEKLSSLIQQYSKILIVHVDNVGSNQMASVRKSLRGK-ATILMGKNTRIRTALKKNLQAV--PQIE    76
RLA0_SULAC     -----MIGLAVTTTKKIAKWKVDEVAELTEKLKTHKTIIIANIEGFPADKLHEIRKKLRGK-ADIKVTKNNLFNIALKNAG-----YDTK       79
RLA0_SULTO     ----MRIMAVITQERKIAKWKIEEVKELEQKLREYHTIIIANIEGFPADKLHDIRKKMRGM-AEIKVTKNTLFGIAAKNAG-----LDVS       80
RLA0_SULSO     ----MKRLALALKQRKVASWKLEEVKELTELIKNSNTILIGNLEGFPADKLHEIRKKLRGK-ATIKVTKNTLFKIAAKNAG-----IDIE       80
RLA0_AERPE     MSVVSLVGQMYKREKPIPEWKTLMLRELEELFSKHRVVLFADLTGTPTFVVQRVRKKLWKK-YPMMVAKKRIILRAMKAAGLE---LDDN      86
RLA0_PYRAE     -MMLAIGKRRYVRTRQYPARKVKIVSEATELLQKYPYVFLFDLHGLSSRILHEYRYRLRRY-GVIKIKPTLFKIAFTKVYGG---IPAE      85
RLA0_METAC     ------MAEERHHTEHIPQWKKDEIENIKELIQSHKVFGMVGIEGILATKMQKIRRDLKDV-AVLKVSRNTLTERALNQLG-----ETIP       78
RLA0_METMA     ------MAEERHHTEHIPQWKKDEIENIKELIQSHKVFGMVRIEGILATKIQKIRRDLKDV-AVLKVSRNTLTERALNQLG-----ESIP       78
RLA0_ARCFU     -----MAAVRGS---PPEYKVRAVEEIKRMISSKPVVAIVSFRNVPAGQMQKIRREFRGK-AEIKVVKNTLLERALDALG-----GDYL       75
RLA0_METKA     MAVKAKGQPPSGYEPKVAEWKRREVKELKELMDEYENVGLVDLEGIPAPQLQEIRAKLRERDTIIRMSRNTLMRIALEEKLDER--PELE      88
RLA0_METTH     ---------MAHVAEWKKKEVQELHDLIKGYEVVGIANLADIPARQLQKMRQTLRDS-ALIRMSKKTLISLALEKAGREL--ENVD        74
RLA0_METTL     -------MITAESEHKIAPWKIEEVNKLKELLKNGQIVALVDMMEVPARQLQEIRDKIR-GTMTLKMSRNTLIERAIKEVAEETGNPEFA      82
RLA0_METVA     -------MIDAKSEHKIAPWKIEEVNALKELLKSANVIALIDMMEVPAVQLQEIRDKIR-DQMTLKMSRNTLIKRAVEEVAEETGNPEFA      82
RLA0_METJA     --------METKVKAHVAPWKIEEVKTLKGLIKSKPVVAIVDMMDVPAPQLQEIRDKIR-DKVKLRMSRNTLIIRALKEAAEELNNPKLA      81
RLA0_PYRAB     ------------MAHVAEWKKKEVEELANLIKSYPVIALVDVSSMPAYPLSQMRRLIRENGGLLRVSRNTLIELAIKKAAQELGKPELE       77
RLA0_PYRHO     ------------MAHVAEWKKKEVEELAKLIKSYPVIALVDVSSMPAYPLSQMRRLIRENGGLLRVSRNTLIELAIKKAAKELGKPELE       77
RLA0_PYRFU     ------------MAHVAEWKKKEVEELANLIKSYPVVALVDVSSMPAYPLSQMRRLIRENNGLLRVSRNTLIELAIKKVAQELGKPELE       77
RLA0_PYRKO     ------------MAHVAEWKKKEVEELANIIKSYPVIALVDVAGVPAYPLSKMRDKLR-GKALLRVSRNTLIELAIKRAAQELGQPELE       76
RLA0_HALMA     -----MSAESERKTETIPEWKQEEVDAIVEMIESYESVGVVNIAGIPSRQLQDMRRDLHGT-AELRVSRNTLLERALDDVD-----DGLE     79
RLA0_HALVO     -----MSESEVRQTEVIPQWKREEVDELVDFIESYESVGVVGVAGIPSRQLQSMRRELHGS-AAVRMSRNTLVNRALDEVN-----DGFE     79
RLA0_HALSA     -----MSAEEQRTTEEVPEWKRQEVAELVDLLETYDSVGVVNVTGIPSKQLQDMRRGLHGQ-AALRMSRNTLLVRALEEAG-----DGLD     79
RLA0_THEAC     -------------MKEVSQQKKELVNEITQRIKASRSVAIVDTAGIRTRQIQDIRGKNRGK-INLKVIKKTLLFKALENLGD----EKLS     72
RLA0_THEVO     -------------MRKINPKKEIVSELAQDITKSKAVAIVDIKGVRTRQMQDIRAKNRDK-VKIKVVKKTLLFKALDSIND----EKLT     72
RLA0_PICTO     -------------MTEPAQWKIDFVKNLENEINSRKVAAIVSIKGLRNNEFQKIRNSIRDK-ARIKVSRARLLRLAIENTGK----NNIV     72
ruler          1.......10........20........30........40........50........60........70........80........90
```

|   | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... |
|---|---|---|---|---|---|---|---|---|---|-----|
| A | $a_1$ |   |   |   |   |   |   |   |   |   |
| C |   |   |   |   |   |   |   |   |   |   |
| D |   |   |   |   |   |   |   |   |   |   |
| E |   |   |   |   |   |   |   |   |   |   |
| F |   |   |   |   |   |   |   |   |   |   |
| ... |   |   |   |   |   |   |   |   |   |   |

|     | 1     | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... |
|-----|-------|---|---|---|---|---|---|---|---|-----|
| A   | $a_1$ |   |   |   |   |   |   |   |   |     |
| C   | $c_1$ |   |   |   |   |   |   |   |   |     |
| D   |       |   |   |   |   |   |   |   |   |     |
| E   |       |   |   |   |   |   |   |   |   |     |
| F   |       |   |   |   |   |   |   |   |   |     |
| ... |       |   |   |   |   |   |   |   |   |     |

|     | 1     | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... |
|-----|-------|---|---|---|---|---|---|---|---|-----|
| A   | $a_1$ |   |   |   |   |   |   |   |   |     |
| C   | $c_1$ |   |   |   |   |   |   |   |   |     |
| D   | $d_1$ |   |   |   |   |   |   |   |   |     |
| E   | $e_1$ |   |   |   |   |   |   |   |   |     |
| F   | $f_1$ |   |   |   |   |   |   |   |   |     |
| ... |       |   |   |   |   |   |   |   |   |     |

|     | 1     | 2     | 3 | 4 | 5 | 6 | 7 | 8 | 9 | ... |
|-----|-------|-------|---|---|---|---|---|---|---|-----|
| A   | $a_1$ | $a_2$ |   |   |   |   |   |   |   |     |
| C   | $c_1$ | $c_2$ |   |   |   |   |   |   |   |     |
| D   | $d_1$ | $d_2$ |   |   |   |   |   |   |   |     |
| E   | $e_1$ | $e_2$ |   |   |   |   |   |   |   |     |
| F   | $f_1$ | $f_2$ |   |   |   |   |   |   |   |     |
| ... |       |       |   |   |   |   |   |   |   |     |

|     | 1     | 2     | 3     | 4     | 5     | 6     | 7     | 8     | 9     | ... |
|-----|-------|-------|-------|-------|-------|-------|-------|-------|-------|-----|
| A   | $a_1$ | $a_2$ | $a_3$ | $a_4$ | $a_5$ | $a_6$ | $a_7$ | $a_8$ | $a_9$ |     |
| C   | $c_1$ | $c_2$ | $c_3$ | $c_4$ | $c_5$ | $c_6$ | $c_7$ | $c_8$ | $c_9$ |     |
| D   | $d_1$ | $d_2$ | $d_3$ | $d_4$ | $d_5$ | $d_6$ | $d_7$ | $d_8$ | $d_9$ |     |
| E   | $e_1$ | $e_2$ | $e_3$ | $e_4$ | $e_5$ | $e_6$ | $e_7$ | $e_8$ | $e_9$ |     |
| F   | $f_1$ | $f_2$ | $f_3$ | $f_4$ | $f_5$ | $f_6$ | $f_7$ | $f_8$ | $f_9$ |     |
| ... |       |       |       |       |       |       |       |       |       |     |

# Megvalósítás: PSI-BLAST

➢ Egy BLAST kereséssel megtaláljuk a közeli homológokat

➢ Ezekből többszörös illesztést készítünk

➢ Ebből származtatjuk a kerső profilt a következő BLAST kereséshez

➢ Ezzel bővül a homológok listája a távolabbi rokonokkal

➢ Az eljárást ismételjük

# Paraméterek



Statisztika

Mátrix paraméterek

Eredmény

Maszkolás

# Másik lehetőség HMMER

➢ Honlap: http://hmmer.org/

➢ Letölthető Linux és Windows verzióban

➢ Bőséges dokumentáció a honlapon

➢ Fehérje szekvenciákat kezel

➢ Egyszerre gyors és pontos módszer

➢ Egy hatékony statisztikai modellen alapul – „Markov modell"

# A programcsalád tagjai

- „phmmer": egy vagy több fehérje szekvenciát keres a fehérje adatbázis ellenében

- „hmmscan": fehérje szekvenciákat keres profil adatbázis ellen

- „hmmsearch": profilokat keres fehérje adatbázis ellenében

- "jackhmmer": interaktív változat

# Paralell keresés

➢ A szekvenciák sorrendje egy adatbázisban esetleges

➢ Nem kell sorban haladni a keresés során

➢ Több processzoros, több magos architekturán párhuzamosan lehet futtatni a keresést – GPU computing

➢ A párhuzamos futtatáshoz nemcsak a kódot kell átírni, de az algoritmust is

# Adatbányászat

➢ Keresés szöveges adatbázisokban

➢ MEDLINE: tudományos szövegek kivonatai

➢ Ez is elsődleges adatbázis

➢ Igen nagy és ingyenesen elérhető

➢ Egy cikk mennyire hasonlít egy másikra?

➢ „számítógépes nyelvészet" – adatbányászat

# eTBLAST

- Első fázis: súlyozott kulcsszó keresés
- Ez gyors, de nem túl érzékeny
- Második fázis: „mondat illesztő" lépés
- Ez az érzékenyebb
- Tartalmuk szerint hasonló cikkeket talál
- Javaslatot tesz a megfelelő újságra
- Hasonló érdeklődésű kutatókat azonosít
- http://helioblast.heliotext.com/

# Mit tanultunk ma?

➢ Az adatbázis keresés lényegében nagyléptékű szekvenciaillesztés.

➢ Nagyon kiforrott technika.

➢ Gyakran a bioinformatikai vizsgálódás kiindulópontja.

## Feladat 5.

- Válassz ki egy érdekes cikket és a kivonatával keress hasonlókat az et-blast rendszerben. Mennyire hatékony a szolgáltatás?

- Esetleg fogalmazd meg egy néhány mondatos kivonatban a téged érdeklő problémát és azzal keress.