



Bioinformatika és genomanalízis az orvostudományban

Keresés adatbázisokban

Cserző Miklós

2020

<https://semmelweis.zoom.us/j/96102872458?pwd=Rk1PL2tqS21sdIUwc3B4eDFCZkNKQT09>



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-mer*ek listáját és hozzárendeljük a helyüket
- Ezt csak egyszer kell megcsinálni
- A kereső szekvenciát is felbontjuk *k-mer*ekre
- És így keresünk az adatbázisban



GCTGACAGCAGCCGCTGCAGCAGCTGCTGCTGCTACCAATGCAG
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

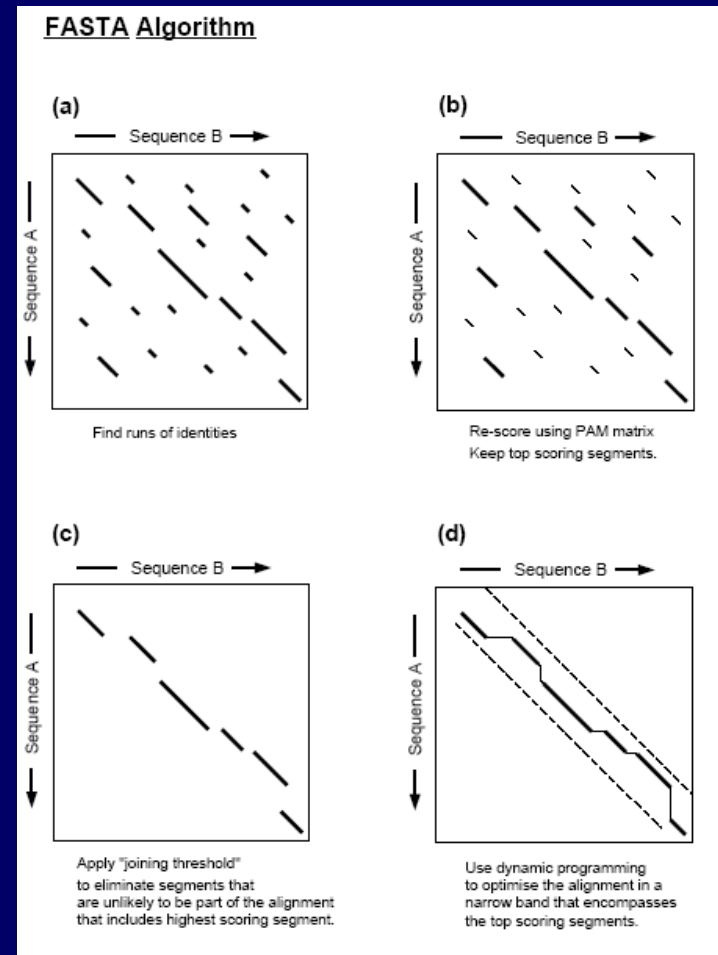


A FASTA család

- 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

- Gyors keresés táblázat alapján (lookup table)
- Pontozás mátrix alapján
- Kiválasztás
- S-W illesztés, végeredmény





A programcsalád tagjai

- A legfrissebb verzió a 3.6-os
 - „fasta”: egy szekvenciát összehasonlít egy adatbázis 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

STEP 1 - Select your databases

PROTEIN DATABASES

1 Databank Selected X Clear Selection

- UniProt Knowledgebase
- UniProtKB/Swiss-Prot
- UniProtKB/Swiss-Prot isoforms
- UniProtKB/TrEMBL
- ▶ UniProtKB Taxonomic Subsets
- ▶ UniProt Clusters
- ▶ Patents
- ▶ Structure
- ▶ Other Protein Databases

STEP 2 - Enter your input sequence

Enter or paste a PROTEIN sequence in any supported format:

or Upload a file: Browse... No file selected.

Paraméterek

A program neve

Táblázat

Statisztika

STEP 3 - Set your parameters

PROGRAM
FASTA

MATRIX: BLOSUM50 GAP.OPEN: -10 GAP.EXTEND: -2 KTUP: 2 EXPECTATION UPPER VALUE: 10 EXPECTATION LOWER VALUE: 0 (default)

DNA STRAND: N/A HISTOGRAM: no FILTER: none STATISTICAL ESTIMATES: Regress

SCORES: 50 ALIGNMENTS: 50 SEQUENCE RANGE: START-END DATABASE RANGE: START-END MULTI HSPs: no

SCORE FORMAT: Default ANNOTATION FEATURES: no

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

Mátrix paraméterek

Eredmény

Maszkolás

FASTA Sequence Comparison at the U. of Virginia

UVa FASTA Server

About

- Getting started
- fasta_guide.pdf

Other FASTA Servers

- EMBL-EBI
- KEGG (Japan)

References

- FASTA
- FASTX/FASTY
- Statistics
- 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

The FASTA programs find regions of local or global (new) similarity between Protein or DNA sequences, either by searching Protein or DNA databases, or by identifying local duplications within a sequence. Other programs provide information on the statistical significance of an alignment. Like BLAST, FASTA can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

<p>Protein</p> <ul style="list-style-type: none"> Protein-protein FASTA Protein-protein Smith-Waterman (ssearch) (New) Global Protein-protein (Needleman-Wunsch) (ggsearch) (New) Global/Local protein-protein (glsearch) Protein-protein with unordered peptides (fasts) Protein-protein with mixed peptide sequences (fastf) 	<p>Nucleotide</p> <ul style="list-style-type: none"> Nucleotide-Nucleotide (DNA/RNA fasta) Ordered Nucleotides vs Nucleotide (fastm) Un-ordered Nucleotides vs Nucleotide (fasts)
<p>Translated</p> <ul style="list-style-type: none"> Translated DNA (with frameshifts, e.g. ESTs) vs Proteins (fastx/fasty) Protein vs Translated DNA (with frameshifts) (tfastx/tfasty) Peptides vs Translated DNA (tfasts) 	<p>Statistical Significance</p> <ul style="list-style-type: none"> Protein vs Protein shuffle (prss) DNA vs DNA shuffle (prsd) Translated DNA vs Protein shuffle (prtx)
<p>Local Duplications</p> <ul style="list-style-type: none"> Local Protein alignments with plots (lalign/plalign) Local DNA alignments with plots (lalign/plalign) 	





UVA FASTA Downloads

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

metabolic pathways poster

Most Visited sajtó 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 (gi121694)** to the PIR1 Annotated protein sequence database.

- What is the highest scoring non-homolog? (How would you confirm that your candidate non-homolog was truly unrelated?)
- 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*)
- Compare the expectation (E0) 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?
- 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 E0-value for the highest scoring unrelated sequence with the different matrices?

Look at the distribution of scores and the E0-value of the highest scoring unrelated sequence when the gap-open/gap-ext penalties are small (-7/-1).
- Try the search with [ssearch](#) (Smith-Waterman). Again, look at the E0-values for distant homologs and the highest scoring unrelated sequence.
- (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 Sequence Comparison

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

metabolic pathways poster

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

Choose: (A) Program, (B) Query (sequence/accession), (C) Database and (D) Start Search:

Annotate Query Sequence
 Annotate Database Sequences

(A) Program: FASTA: protein:protein

Compare your own sequences:

(B) Query sequence: FASTA format Use Subset range

```
>FASTA exercise
MVDFFYYLPGSSPCRSVIMTAKAVGVELNKKLLNLQAGEHLKPEFLKINPOHTIPTLVLDNG
FALWESRAIQ
VYLVKEYGKTDLSLYPKCPKRAVINQRLYFDMGTLYQSFANYYPQVFAKAPADPEAFKK
IEAAFEFLNT
FLEGQDYAAGDSLTVADIALVATVSTFEVAKFEISKYANVNRWYENAKKVTGWEENWAG
CLEFKKYFE
```

[Entrez protein sequence browser](#)
[Entrez DNA sequence browser](#)

Or upload query from file:

Protein DNA (both-strands) DNA (forward only) DNA (rev-comp only)

(C) Database: **(D) Start Search**

Protein: PIR1 Annotated (rel. 66) DNA: GB170.0 Primate

Exclude low complexity (seg)

Comments (optional):

Other search options: **Output limits:** Show Histogram

Scoring matrix: open: ext: Ktup: Statistical estimates
 Blossum50 (20%) -10 -2 ktup = 2 Default

E(): Best E():

Alignment Options: Highlight similarities differences compact differences.

[FASTA program information](#) | [Download FASTA](#) | [About the Author](#)

Copyright © 1988, 2006 by William R. Pearson and the University of Virginia. All rights reserved. The FASTA program and documentation may not be sold or incorporated into a commercial product, in whole or in part, without written

```

File Edit View History Bookmarks Tools Help
FASTA UVA FASTA Downloads RecName: Full=Glutathione S-transfe... FASTA results
fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi
metabolic pathways poster
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:=====*
110     1      1:=====*
112     1      1:=====*
114     0      1:=====*
116     1      0:=====*
118     0      0:=====*

      inset = represents 1 library sequences
  
```



```

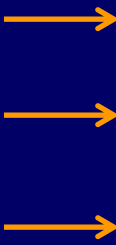
File Edit View History Bookmarks Tools Help
FASTA UVA FASTA Downloads RecName: Full=Glutathione S-transfe... FASTA results
fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi
metabolic pathways poster

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

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:           *
120 0 0:           *
122 0 0:           *
124 0 0:           *
126 0 0:           *
128 0 0:           *
130 0 0:           *
132 1 0:=          *=
134 0 0:           *
136 1 0:=          *=
138 0 0:           *
>140 5 0:=          =====
S190221 residues in 13351 sequences
Statistics: Expectation n fit: rho(ln(x))= 7.1276+/-0.00254; mu= 7.1152+/- 0.130
mean_var=50.7864+/-10.249, 0's: 7 Z-trim(89.1): 32 B-trim: 0 in 0/51
Lambda= 0.179970
statistics sampled from 1889 (1896) to 1889 sequences
Kolmogorov-Smirnov statistic: 0.0423 (N=19) at 52
Algorithm: FASTA (3.8 Nov 2011) [optimized]
Parameters: BL50 matrix (15:-5)XS, open/ext: -10/-2
ktup: 2, E-join: 1 (0.464), E-opt: 0.2 (0.142), width: 16
Scan time: 0.270

The best scores are:
opt bits E(13351) % id % sim alen
sp|P20432|GSTT1_DROME Glutathione S-transferase 1-1 (GS (209) 1399 370.6 1.6e-103 1.000 1.000 209 align
sp|P04907|GSTF3_MAIZE Glutathione S-transferase III (GS (222) 173 52.2 1.2e-07 0.264 0.557 212 align
sp|P12653|GSTF1_MAIZE Glutathione S-transferase I (GST- (214) 151 46.5 5.9e-06 0.276 0.525 181 align
sp|P0ACAS|SSPA_EC057 Stringent starvation protein A gi| (212) 140 43.7 4.2e-05 0.263 0.593 118 align
sp|P00502|GSTA1_RAT Glutathione S-transferase alpha-1 ( (222) 139 43.4 5.4e-05 0.286 0.566 182 align
sp|P14942|GSTA4_RAT Glutathione S-transferase alpha-4 ( (222) 97 32.5 0.1 0.282 0.563 174 align
sp|P08010|GSTM2_RAT Glutathione S-transferase Mu 2 (GST (218) 93 31.5 0.21 0.221 0.517 145 align
sp|P09211|GSTP1_HUMAN Glutathione S-transferase P (GST (210) 82 28.6 1.4 0.193 0.532 171 align
sp|P09457|ATPO_YEAST ATP synthase subunit 5, mitochondr (212) 79 27.8 2.5 0.286 0.556 63 align
sp|P00925|ENO2_YEAST Enolase 2 (2-phosphoglycerate dehy (437) 83 28.6 3 0.264 0.536 125 align
sp|P0A4L1|THIO1_ANASP Thioredoxin 1 (TRX-1) (Thioredoxi (107) 72 26.3 3.6 0.246 0.596 57 align
sp|P21163|ENGF_ELMIR 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 initl: 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)
    
```



File Edit View History Bookmarks Tools Help

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

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

Most Visited sajto SOTE Logins Tool DAS IT Biosites Library Athéna 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 initl: 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 MVDFFYLPGSSPCRSVMITAKAVGVELNKKLLNLQAGEHLKPEFLKINPQHTIPTLVNDNGFALWESRAIQVYLVEKYGKI
  : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
  sp|P20 MVDFFYLPGSSPCRSVMITAKAVGVELNKKLLNLQAGEHLKPEFLKINPQHTIPTLVNDNGFALWESRAIQVYLVEKYGKI
  10 20 30 40 50 60 70 80
  90 100 110 120 130 140 150 160
  FASTA DSYLYPKCKKRAVINQRLYFDMGTLQYSFANYYPQVFAKAPADPEAFKKIEAAFEFLNIFLEGQDYAAGDSLTVADIAL
  : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
  sp|P20 DSYLYPKCKKRAVINQRLYFDMGTLQYSFANYYPQVFAKAPADPEAFKKIEAAFEFLNIFLEGQDYAAGDSLTVADIAL
  90 100 110 120 130 140 150 160
  170 180 190 200
  FASTA VATVSTFEVAKFEISKYANVNRWYENAKKVTIPGWEENWAGCLEFKKYFE
  : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
  sp|P20 VATVSTFEVAKFEISKYANVNRWYENAKKVTIPGWEENWAGCLEFKKYFE
  170 180 190 200
```

```
>>sp|P04907|GSTF3_MAIZE Glutathione S-transferase III (GST-III) ( (222 aa)
  initn: 182 initl: 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 MVDFFYLPGSSPCRSVMITAKAVGVELNKKLLNLQAGEHLKPEFLKINPQHTIPTLVNDNGFALWESRAIQVYLVEKYG
  : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
  sp|P04 MAPLKLYGMPLSPNVVRRVATVINEKGLDFEIVPVDLTTGAHKQPDFLALNPPFGQIPALVDGDVLFESRAINRYIAKYA
  10 20 30 40 50 60 70 80
  80 90 100 110 120 130 140
  FASTA K--TDSLYPKCKKRAVINQRLYFDMGTLQYSFANYYPQVF-----AKAPADPEAFKKIEAAFEFLNIF---LEGQ
  : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
  sp|P04 SEGTDLL----PATASAAKLEWLVES--HHFHPNASPLVFLQLLVRPLGGAPDAARVVEKHAQQLAKVLDVVEAHLARN
  90 100 110 120 130 140 150
  150 160 170 180 190 200
  FASTA DYAAGDSLTVADI--ALVATVSTFEVAKFE--ISKYANVNRWYENAKKVTIPGWEENWAGCLEFKKYFE
  : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
  sp|P04 KYLAGDEFLLADANHALLPALTSARPPRPGCVARRPHVKAWE--AIAARPAFKQTVAAIpppppsaA
  160 170 180 190 200 210 220
```



```

File Edit View History Bookmarks Tools Help
FASTA UVA FASTA Downloads x RecName: Full=Glutathione S-transfe... x FASTA results
fasta.bioch.virginia.edu/fasta_www2/fasta_www.cgi
metabolic pathways poster

160 170 180 190 200 210

>>sp|P0ACA5|SSPA_ECO57 Stringent starvation protein A gi|8117473 (212 aa)
  initn: 138 initl: 115 opt: 140 Z-score: 198.1 bits: 43.7 E(13351): 4.2e-05
  Smith-Waterman score: 140; 26.3% identity (59.3% similar) in 118 aa overlap (43-160:49-161)
  Entrez Lookup Re-search database General re-search
FASTA 10 20 30 40 50 60 70 80
FASTA DFYYLPGSSPCRSVIMTAKAVGVELNKKLLNLQAGEHLKPEFLKINPQHTIPTLVNDGFALWESRAIQVYLVEKYGKTD
sp|P0A SVMTLFSGPTDIYSHQVRIVLAEKGVSFIEHVEKDNPPQDLIDLNPQSVPTLVDRLETLWESRIIMEYLDERFPHPP-
  10 20 30 40 50 60 70 80
FASTA 90 100 110 120 130 140 150 160
FASTA LYPKCPKRAVINQRLYFDMGTLYQSFNYYYPQVFAKAPADPEAFKKIEAAFEFLNTFLEGQDYAAGDSLTVADIALVA
  .. : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
sp|P0A LMPVYVPVARG--ESRLY--MHRIEKDWTLMNTIINGSASEADAARKQLREELLAIAIPVFGQKPYFLSDEFSLVDCYLP
  90 100 110 120 130 140 150 160
FASTA 170 180 190 200
FASTA TVSTFEVAKFEISKYANVNRWYENAKKVTIPGWEENWAGCLEFHKYFE
sp|P0A LLWRLPQLGIEFSGPGAKELKGYMTRVFERDSFLASLTEAEREMRLGRS
  170 180 190 200 210

>>sp|P00502|GSTA1_RAT Glutathione S-transferase alpha-1 (Glutath
  initn: 109 initl: 62 opt: 139 Z-score: 196.2 bits: 43.4 E(13351): 5.4e-05
  Smith-Waterman score: 139; 29.6% identity (56.6% similar) in 182 aa overlap (1-175:6-173)
  Entrez Lookup Re-search database General re-search
FASTA 10 20 30 40 50 60 70
FASTA MVDFYYLPGSSPCRSVIMTAKAVGVELNKKLLNLQAGEHLKPEFLKINPQ---HTIPTLVNDGFALWESRAIQVY
  .. : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
sp|P00 MSGKPVHYFNARGRMCEIRWLLA--AAGVFDEKFI--QSPEDL--EKLKDGNDLMPDQVPMVEIDGMKLAQTRAILNY
  10 20 30 40 50 60 70
FASTA 80 90 100 110 120 130 140
FASTA LVEKYGKTDLSLYPKCPKRAVINQRL--YFDMGTLYQSFNYYYPQVFAK-APADPEAFKKIEAAFE-FLNTFLEGQDYA
  .. : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
sp|P00 IATKY----DLYGKMKERALIDMYTEGILDLEMIMQLVICPPDQKEAKTALAKDRTHNRYLPAFEKVLKS--HGQDYL
  80 90 100 110 120 130 140
FASTA 150 160 170 180 190 200
FASTA AGDSLTVADIALVAIVSTFEVAKFEISKYANVNRWYENAKKVTIPGWEENWAGCLEFHKYFE
  .. : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : : :
sp|P00 VGNRLTRVDIHLELL--LVVEFDASLLTSFPPLLKAFKSRISLSPNVKFLQPGSQRKLPVDAKQIEEARKIFKF
  150 160 170 180 190 200 210 220

>>sp|P14942|GSTA4_RAT Glutathione S-transferase alpha-4 (Glutathio
  initn: 80 initl: 53 opt: 106 Z-score: 149.9 bits: 34.8 E(13351): 0.02
  Smith-Waterman score: 132; 28.2% identity (56.3% similar) in 174 aa overlap (4-168:7-168)
  Entrez Lookup Re-search database General re-search
FASTA 10 20 30 40 50 60 70

```




A BLAST család

- Honlap:
<http://blast.ncbi.nlm.nih.gov/Blast.cgi>
- Letölthető 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



COVID-19 is an emerging, rapidly evolving situation.
 Get the latest public health information from CDC: <https://www.coronavirus.gov>.
 Get the latest research from NIH: <https://www.nih.gov/coronavirus>.
 Find NCBI SARS-CoV-2 literature, sequence, and clinical content: <https://www.ncbi.nlm.nih.gov/sars-cov-2/>.

Basic Local Alignment Search Tool

BLAST finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

A new feature was added to Primer-BLAST.

We have added a new function to Primer-BLAST that helps users design primers common for a group of highly similar sequences.

Tue, 29 Sep 2020 12:00:00 EST [More BLAST news...](#)

Web BLAST

Nucleotide BLAST
nucleotide ▶ nucleotide

blastx
translated nucleotide ▶ protein

Protein BLAST
protein ▶ protein

tblastn
protein ▶ translated nucleotide

BLAST Genomes

Enter organism common name, scientific name, or tax id

Human Mouse Rat Microbes

NCBI BLAST+
Protein Nucleotide Vectors Web services Help & Documentation

Tools > Sequence Similarity Searching > NCBI BLAST

Protein Similarity Search

The emphasis of this tool is to find regions of sequence similarity, which will yield functional and evolutionary clues about the structure and function of your novel sequence.

A new, more accurate, search tool combining optimal searching with iterative profile generation and over-extension error prevention is available using [PSI-Search](#).

STEP 1 - Select your databases

PROTEIN DATABASES

1 Databank Selected X Clear Selection

- UniProt Knowledgebase
- UniProtKB/Swiss-Prot
- UniProtKB/Swiss-Prot isoforms
- UniProtKB/TrEMBL
- ▶ UniProtKB Taxonomic Subsets
- ▶ UniProt Clusters
- ▶ Patents
- ▶ Structure
- ▶ Other Protein Databases

STEP 2 - Enter your input sequence

Enter or paste a **PROTEIN** sequence in any supported format:

or upload a file: No file selected.

NCBI BLAST < Sequence Si... x UVA FASTA Downloads x BLAST: Basic Local Alignm... x +

www.ebi.ac.uk/Tools/sss/ncbiblast/

decatlon

- UniProtKB/Swiss-Prot
- UniProtKB/Swiss-Prot isoforms
- UniProtKB/TrEMBL
- UniProtKB Taxonomic Subsets
- UniProt Clusters
- Patents
- Structure
- Other Protein Databases

STEP 2 - Enter your input sequence

Enter or paste a **PROTEIN** sequence in any supported format:

or upload a file: **Browse...** No file selected.

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	COMPOSITION-BASED STATISTICS			
pairwise	F (default)			

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

If you plan to use these services during a course please [contact us](#).

EMBL-EBI Services Research Training Industry About us

Paraméterek

A program neve

Mátrix paraméterek

Eredmény

Statisztika

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	GAP ALIGN
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

Maszkolás



Egy további változat

BLAT:

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



Profilkeresés

- Az aminosavak 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	-----MPREDRATWKS	NYFLKIIQL	LLDDYPKCFIVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76																													
RLA0_HUMAN	-----MPREDRATWKS	NYFLKIIQL	LLDDYPKCFIVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76																													
RLA0_MOUSE	-----MPREDRATWKS	NYFLKIIQL	LLDDYPKCFIVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76																													
RLA0_RAT	-----MPREDRATWKS	NYFLKIIQL	LLDDYPKCFIVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76																													
RLA0_CHICK	-----MPREDRATWKS	NYFMKIIQL	LLDDYPKCFVVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76																													
RLA0_RANSY	-----MPREDRATWKS	NYFLKIIQL	LLDDYPKCFIVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	SALE	76																													
Q7_ZUG3_BRARE	-----MPREDRATWKS	NYFLKIIQL	LLDDYPKCFIVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76																													
RLA0 ICTPU	-----MPREDRATWKS	NYFLKIIQL	LLNDYPKCFIVGADNVGS	KOMQOIRMSLRGK	AVVLMGKNTMMRKAIRGHLENN	PALE	76																													
RLA0_DROME	-----MVRENKAAWKAQY	FIKVVLEFDEF	PKCFIVGADNVGS	KOMONIRTSLRGL	AVVLMGKNTMMRKAIRGHLENN	PQLE	76																													
RLA0_DICDI	-----MSGAG-SKRK	KLFIEKATKLF	TTYDKMIVAEAD	FVGS	SQLOKIRKSIRGI	GAVLMGKNTMIRKVIDRLADSK	PELD	75																												
Q54LP0_DICDI	-----MSGAG-SKRK	NVFIKATKLF	TTYDKMIVAEAD	FVGS	SQLOKIRKSIRGI	GAVLMGKNTMIRKVIDRLADSK	PELD	75																												
RLA0_PLAF8	-----MAKLSKQQKQMY	EKLSSLIQYSK	ILLVHVDNVGS	NOMASVRKSLRGK	ATILMGKNTIRRTALKKNLQAV	PQIE	76																													
RLA0_SULAC	-----MIGLAVTTT	KKIAKWKVDE	VAELTEKLRKHT	IIIANIEGFP	PADKLHEIRKCLR	ADIKVTKNNLNFNIALKNAG	YDTK	79																												
RLA0_SULTO	-----MRIMAVITQER	KIAKWKIEE	VKELEOKLREY	HTIIIANIEGFP	PADKLHDIRKKMRGM	AEIKVTKNTLFGIAAKNAG	LDVS	80																												
RLA0_SULSO	-----MKRLALALKQ	RKVASWKL	EEVKELETEL	IKNSNTILIGN	LEGFPADKLHEIRKCLR	ATIKVTKNTLFGIAAKNAG	IDIE	80																												
RLA0_AERPE	MSVVS	SLVGQMYKRE	KPIPEW	KTMLRELEEL	FSKIRVVLF	ADLTGTP	FVVQRVRK	KLWKK	YPM	MAK	KRIIL	RAMKA	AAGLE	---	LDDN	86																				
RLA0_PYRAE	-----MMLAIG	KRRYVRT	RQYP	PARKVKIV	SEATEL	LQKYP	VYVFL	DLHGLSS	RILHE	YRRL	RRY	GVIK	IKPTL	FKIAFT	KVYGG	---	IPAE	85																		
RLA0_METAC	-----MAEERH	TEHIP	QWKDEI	ENIKEL	IQSHK	VFGM	VGIEG	ILATK	MOKIR	RD	LKDV	AVL	KVSR	NL	TERAL	NQLG	---	ETIP	78																	
RLA0_METMA	-----MAEERH	TEHIP	QWKDEI	ENIKEL	IQSHK	VFGM	VGIEG	ILATK	IQIR	RD	LKDV	AVL	KVSR	NL	TERAL	NQLG	---	ESIP	78																	
RLA0_ARCFU	-----MAAVR	GS---	PPEY	KVRAVEE	IKRM	ISSK	PVVA	IVSFR	NVPAG	OMQ	IRRE	FRGK	AEIK	VV	KNTL	LERAL	DALG	---	GDYL	75																
RLA0_METKA	MAVKAKG	QPPSGYE	PKVAE	WKRRE	VKEL	ELMDE	YENV	GLVD	LEGIP	APQL	QEI	RAKLR	ERTI	IRMS	RNTL	MRIA	LEEK	LDER	---	PELE	88															
RLA0_METTH	-----MAHVAE	WKKKEV	QELHDL	IKGYE	VVGIAN	LADIP	ARQL	QKMR	QTLRDS	ALIRMS	SKKTL	LISL	LALE	KAG	REL	---	ENVD	74																		
RLA0_METTL	-----MITAE	SEHKI	APWKIE	EVNKL	KELLK	NGQ	IVAL	VD	MMEV	PARQL	QEI	IRDK	IR	GTMT	LKMS	RNTL	LIE	RAI	KEVA	ETGN	PEFA	82														
RLA0_METVA	-----MIDAK	SEHKI	APWKIE	EVNKL	KELLK	NSAN	VIAL	ID	MMEV	PAVQL	QEI	IRDK	IR	DQMT	LKMS	RNTL	LIE	RAI	KEVA	ETGN	PEFA	82														
RLA0_METJA	-----METK	VKAH	VAPWKIE	EVKTL	KG	LKSK	PVVA	IVD	MMDV	PAPQL	QEI	IRDK	IR	DKV	KLRMS	RNTL	LIE	RAI	KEVA	ETGN	PEFA	81														
RLA0_PYRAB	-----MAHVAE	WKKKEV	EELAN	LKSY	PVIAL	VDV	SSMP	PAYPL	SQMR	RLIRE	NGGL	LRSR	NL	LIE	LAIK	KA	QELG	KPE	LE	77																
RLA0_PYRHO	-----MAHVAE	WKKKEV	EELAK	LKSY	PVIAL	VDV	SSMP	PAYPL	SQMR	RLIRE	NGGL	LRSR	NL	LIE	LAIK	KA	QELG	KPE	LE	77																
RLA0_PYRFU	-----MAHVAE	WKKKEV	EELAN	LKSY	PVVAL	VDV	SSMP	PAYPL	SQMR	RLIRE	NGGL	LRSR	NL	LIE	LAIK	KA	QELG	KPE	LE	77																
RLA0_PYRKO	-----MAHVAE	WKKKEV	EELAN	LKSY	PVIAL	VDV	AGV	PAYPL	SKMR	DKLR	GKALL	RSR	NL	LIE	LAIK	RA	QELG	KPE	LE	76																
RLA0_HALMA	-----MSAES	ERKT	ETIPEW	KQEE	VD	AI	EM	IES	SVG	VVNI	IAGIP	SRQL	QDM	RRD	IHGT	AE	LRSR	NL	LIE	RAL	DDVD	---	DGLE	79												
RLA0_HALVO	-----MSESE	VRQTE	EVIP	QWKREE	VDEL	VDF	IES	YES	SVG	VGV	VAGIP	SRQL	QSM	RRLH	GS	AAV	RMSR	NL	LVN	RAL	DEVN	---	DGFE	79												
RLA0_HALSA	-----MSAEE	QRTTE	EVPEW	KRQ	EAEL	VDL	LET	YDS	VG	VVNV	TGIP	SKOL	QDM	RRGL	HGQ	AA	LMSR	NL	LLV	RA	LEE	AG	---	DGLD	79											
RLA0_THEAC	-----MKEV	SQK	KE	LVNE	IT	OR	KAS	RS	V	AI	VD	LAG	IR	ROI	QD	IR	GK	NR	RGK	IN	LKV	IK	KTLL	FK	KALE	NL	GD	---	EKLS	72						
RLA0_THEVO	-----MRKIN	PKK	KE	IV	SELA	QD	IT	KS	K	AV	AI	VD	IK	GV	RT	ROM	QD	IR	AK	NR	DK	VK	IK	VV	KTLL	FK	KAL	S	IND	---	EKLT	72				
RLA0_PICTO	-----MTEPA	QWK	ID	FV	KN	LE	INS	RK	V	AA	IV	SK	L	R	NE	F	Q	K	IR	NS	IR	DK	AR	IK	V	S	R	AR	LL	RL	AI	ENT	GK	---	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

STEP 3 - Set your parameters

PSI-BLAST THRESHOLD
1.0e-3

MATRIX	GAP OPEN	GAP EXTEND	EXPECTED THRESHOLD	FILTER
BLOSUM62	11	1	10.0	no

SCORES	ALIGNMENTS	SEQUENCE RANGE	DROPOFF	FINAL DROPOFF
500	500	START-END	15 (default)	25 (default)

ALIGNMENT VIEW
pairwise

USAGE MODE FOR PHI-BLAST
blastpgp

UPLOAD A CHECKPOINT FILE (ASN.1 Binary Format)
Browse...

UPLOAD A PATTERN FILE FOR PHI-BLAST
Browse...

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

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
- “jackhammer”: interaktív változat



Ész Ventura: Lét... Python Programmi... Top 8 resources for... Python Pandas Tut... Services and Suppo... MetaXpress® and li... Telefonkönyv... Tájékoztató, Kli... FASTA UVA FASTA Ser... HMMER

hmmmer.org

HMMER DOWNLOAD DOCUMENTATION SEARCH PUBLICATIONS BLOG

HMMER: biosequence analysis using profile hidden Markov models

Get the latest version

v3.2.1

[Download source](#)
(archived older versions)

HMMER is used for searching sequence databases for sequence homologs, and for making sequence alignments. It implements methods using probabilistic models called profile hidden Markov models (profile HMMs).

HMMER is often used together with a profile database, such as [Pfam](#) or many of the databases that participate in [Interpro](#). But HMMER can also work with query *sequences*, not just profiles, just like BLAST. For example, you can search a protein query sequence against a database with [phmmer](#), or do an iterative search with [jackhmmmer](#).

HMMER is designed to detect remote homologs as sensitively as possible, relying on the strength of its underlying probability models. In the past, this strength came at significant computational expense, but as of the new HMMER3 project, HMMER is now essentially as fast as BLAST.

HMMER can be downloaded and installed as a command line tool on your own hardware, and now it is also more widely accessible to the scientific community via [new search servers](#) at the European Bioinformatics Institute.

PERFORM A SEARCH

An online interactive [search](#) service is available at the European Bioinformatics Institute. Go there to [search](#) against the latest Uniprot databases.

DOCUMENTATION

The HMMER User's Guide: [\[PDF\]](#).

NEWS

See the blog [Cryptogenomicon](#) for more information and discussion about HMMER3.



The screenshot shows the HMMER web server interface. The browser window has several tabs open, including 'UVA FASTA Server', 'FASTA/SSEARCH/GGSEARCH', 'Help: FASTA (Protein Data)', 'BLAST: Basic Local Alignment', and 'Biosequence analysis using'. The address bar shows 'https://www.ebi.ac.uk/Tools/hmmer/'. The page header includes 'EMBL-EBI', 'Services', 'Research', 'Training', 'About us', and 'EMBL-EBI Hinxton'. The main heading is 'HMMER Biosequence analysis using profile hidden Markov Models'. A navigation menu includes 'Home', 'Search', 'Results', 'Software', 'Help', 'About', and 'Contact'. The 'Quick search' section features a text input field with a placeholder 'Paste in your sequence or use the example', radio buttons for 'Reference Proteomes', 'UniProtKB', 'SwissProt', and 'Pfam', and 'Submit' and 'Reset' buttons. To the right, a text block describes the server as 'fast and sensitive homology searches' and provides links for 'Quickstart tutorial' and 'Online documentation'. The footer contains 'Blog news' (dated August, 2015), 'Download HMMER v3.1b2', and 'Recent papers' (with a link to 'HMMER web server: 2015 update').



EMBL-EBI Services Research Training About us Hinxton

HMMER

Biosequence analysis using profile hidden Markov Models

Home Search Results Software Help About Contact

phmmer hmmscan hmmsearch jackhmmer

protein sequence vs protein sequence database

[Paste a Sequence](#) | [Upload a File](#) | [Accession Search](#)

Paste in your sequence or use the [example](#)

Submit Reset

Sequence Database

Frequently used databases: [Reference Proteomes](#) [UniProtKB](#) [SwissProt](#) [PDB](#) [Ensembl](#)



www.ebi.ac.uk/Tools/hmmer/search/phmmer

Home Search Results Software Help About

Sequence Database

Frequently used databases

Reference Proteomes UniProtKB SwissProt PDB

Representative Sets (UniProt)

rp75 rp55 rp35 rp15

Other databases

QfO Pfamseq

► Restrict by Taxonomy

Cut-Offs

E-value Bit score

Significance E-values: Sequence Hit

Report E-values: Sequence Hit

Customize Results

Select Visible Columns

Row Count Identical Seqs

Secondary Accessions and Ids Number of Hits

Description Number of Significant Hits

Species Bit Score

Kingdom Hit Positions

Known Structure

Rows Per Page

50

100

250

1000

2500

All

Gap Penalties

open extend

Substitution scoring matrix:

Filters

Turn off bias composition filter



Browser address bar: <https://www.ebi.ac.uk/Tools/hmmer/results/3AED41F2-0C87-11E8>

EMBL-EBI Services Research Training About us

HMMER

Biosequence analysis using profile hidden Markov Models

Home Search **Results** Software Help About Contact

PHMMER Results Search Again

Score Taxonomy Domain Download

Sequence Matches and Features

Pfam GSTU_3 209

hit coverage

hit similarity

disorder coiled-coil tm & signal peptide

[Show hit details](#)

Distribution of Significant Hits

more significant

- Bacteria
- Eukaryota
- Archaea
- Viruses
- Unclassified Sequences
- Other Sequences

« First « Previous **Page 1** of 579 Next » Last »

Significant Query Matches (21376) in *uniprotrefprot* (v.2019_09) Customise



Browser tabs: pets/A, diger, Built-i, Topology, TCB, ks.uui, Neptu, Semm, Re: [AMBE], Bejele, Telefo, UVA F, RecNa, BLAST

Address bar: <https://www.ebi.ac.uk/Tools/hmmer/results/3AED41F2-0C87-11EB>

Navigation: Home Search Results Software Help About Contact

Pfam: GST_C [209]

hit coverage [Heatmap]

hit similarity [Heatmap]

✓ disorder ✓ coiled-coil ✓ tm & signal peptide

Show hit details

Distribution of Significant Hits

Legend: Bacteria (red), Eukaryota (yellow), Archaea (blue), Viruses (orange), Unclassified Sequences (grey), Other Sequences (black)

Page 1 of 579

Significant Query Matches (21376) in *uniprotrefprot* (v.2019_09)

Target	Description	Species	Cross-references	E-value
> A0A182JTU9_9DIPT	Uncharacterized protein	Anopheles christyi		0.0e+00
> A0A0L0BYE0_LUCCU	Uncharacterized protein	Lucilia cuprina		0.0e+00
> B0W6B0_CULQU	Glutathione S-transferase 1-6	Culex quinquefasciatus		5.3e-284
> A0A182QC27_9DIPT	Uncharacterized protein	Anopheles farauti		2.2e-265
> A0A182UJV8_9DIPT	Uncharacterized protein	Anopheles melas		3.5e-261
> A0A1W4UME7_DROFC	uncharacterized protein LOC108087889	Drosophila fusciphila		9.0e-252
> A0A182MSE4_9DIPT	Uncharacterized protein	Anopheles culicifacies		8.2e-245
(show all) alignments	A0A3B8YIY4_BLAGE	Uncharacterized protein (Fragment)	Blattella germanica	8.6e-244

Your search took: 2.51 secs showing rows 1 - 50 of 28945



Browser: <https://www.ebi.ac.uk/Tools/hmmer/results/3AED41F2-0C87-11E8>

Navigation: Home Search Results Software Help About Contact

Filters: Bacteria Eukaryota Archaea Viruses Unclassified Sequences Other Sequences

Page 1 of 579

Significant Query Matches (21376) in *uniprotrefprot* (v.2019_09)

Target	Description	Species	Cross-references	E-value
> A0A182JTU9_9DIPT	Uncharacterized protein	Anopheles christyi		0.0e+00
> A0A0L0BYE0_LUCCU	Uncharacterized protein	Lucilia cuprina		0.0e+00
> B0W6B0_CULQU	Glutathione S-transferase 1-6	Culex quinquefasciatus		5.3e-284
> A0A182QC27_9DIPT	Uncharacterized protein	Anopheles farauti		2.2e-265
> A0A182UJV8_9DIPT	Uncharacterized protein	Anopheles melas		3.5e-261
> A0A1W4UME7_DROFC	uncharacterized protein LOC108087889	Drosophila fusciphila		9.0e-252
> A0A182MSE4_9DIPT	Uncharacterized protein	Anopheles culicifacies		8.2e-245
> A0A2P8XIY1_BLAGE	Uncharacterized protein (Fragment)	Blattella germanica		2.6e-240
> A0A118MLE7_MUSDO	Uncharacterized protein	Musca domestica		3.9e-225
> B4K5W6_DROMO	Uncharacterized protein	Drosophila mojavensis		4.7e-220
> A0A182YGN2_ANOST	Uncharacterized protein	Anopheles stephensi		3.3e-208
> A0A1J1IJQ6_9DIPT	CLUMA_CG013690, isoform A	Clunio marinus		3.4e-200
> A0A0P8ZYE2_DROAN	Uncharacterized protein	Drosophila ananassae		2.7e-199
> A0A1W4UML8_DROFC	uncharacterized protein LOC108087890	Drosophila fusciphila		2.6e-198

(show all) alignments DROFC uncharacterized protein LOC108087890 Your search took: 2.51 secs Drosophila fusciphila showing rows 1 - 50 of 28945



Browser: https://www.ebi.ac.uk/Tools/hmmer/results/3AED41F2-0C87-11EB

Navigation: Home Search Results Software Help About Contact

Target	Description	Species	Cross-references	E-value
▼ B0W6B0_CULQU	Glutathione S-transferase 1-6	Culex quinquefasciatus		5.3e-284

Query	Target Envelope		Target Alignment		Bias	Accuracy	% Identity (count)	% Similarity (count)	Bit Score	E-value		
	start	end	start	end						Ind.	Cond.	
2	202	1	205	1	201	0.08	0.99	65.2 (131)	80.6 (162)	296.2	2.9e-85	1.7e-88

Query: 2 vdfyyllpgsspcrswimtakavgvelnkllnlqagehlpkpeflkinpghctiptlvdnqfalwesraiqvylvekygktd 81
 +dfyyllpgs+pcr+v mta avgveln kl nl ageh+kpeflk+npqh iptlvd +f +wesrai +ylvekygk +

Target: 1 MDFYYLPGSAPCRVQMTAAAVGVELNLKLTNLMAGEHMKPEFLKLNPHQHCIPTLVDSFPVWESRAIMIYLVKEYGKDE 80
 PP 59*****

Query: 82 slypkcckkravingnrlfydmgtylqsfanyvypqvfakapadpeafkkieaafeflntflgqcyagdelvadialv 161
 slypk p+krav+nqrl+fd+gtylq fa+y+ypq+fak pa+ + kk+ ++flntfl g y agd lt+ad+ ++

Target: 81 SLYPKDPKRAVVNQRLLFFDQGTLYQRFADYFYPQIFAKQANADNEKMLDGLDFLNTPLGSKYVAGDQTLTADLTII 160
 PP *****

Query: 162 atvstfevakfeiskyanvnrwyenakvtvgweenwagcl 202
 atvst++vak +++ky nv wy +k pg n agc

Target: 161 ATVSTYDVAKVDLAKYPNVAGWYARLRKEAPGAINEAGCC 201
 PP *****94

This row shows the alignment between your sequence and the matching HMM.

Query consensus of the HMM (query sequence), coloured according to the match: Identical residues █; Similar residues █.

Match: the match between the query sequence (HMM) and target sequence

Target query sequence, coloured according to the posterior probability:
 0% 100%

PP: posterior probability, or the degree of confidence in each individual aligned residue

Query	Target Envelope	Target	Count	Bit Score	E-value	
start	start	start			Ind.	Cond.
39	207	217	223	203.8	5.9e-57	3.5e-60

Query: 39 hlkpeflkinpghctiptlvdnqfalwesraiqvylvekygkdslypkcckkravingnrlfydmgtylqsfanyvypqvf 118
 h + +inpqh iptlvd+gf+lwesraiq+ylveky+k d+slypk p+krav+nqrl+fd+g lvg f y+ypq+f

(show all) alignments Your search took: 2.51 secs showing rows 1 - 50 of 28945



Browser: https://www.ebi.ac.uk/Tools/hmmer/results/3AED41F2-0C87-11E8

Navigation: Home Search Results Software Help About Contact

Target	Description	Species	Cross-references	E-value
> A0A1Y5RYI9_9RHOB	Glutaredoxin 2	Pseudoruegeria aquimaris	xxx [grid] [document]	0.99
> A0A0J5KUX3_PLUGE	Glutathione S-transferase	Pluralibacter gergoviae	xxx [grid] [document]	0.99
> A0A2U2HLL2_9BURK	Maleylacetoacetate isomerase	Massilia glaciei	xxx [grid] [document]	0.99
▼ A0A1C7D796_9SPHN	Maleylpyruvate isomerase	Altererythrobacter namhicola	xxx [grid] [document] [refresh]	0.99

Query		Target Envelope		Target Alignment		Bias	Accuracy	% Identity (count)	% Similarity (count)	Bit Score	E-value	
start	end	start	end	start	end						Ind.	Cond.
6	94	4	180	9	95	0.00	0.84	29.9 (26)	51.7 (45)	19.3	1.8	0.001

Insignificant match:

```

.....*.....*.....*.....*.....*.....*.....*.....*.....*
Query 6 ylpqsspcrsyvtakavqvalnkklilnlcageelkpeelklnpchtitllvdngfallewraalqvylvekvgktdel 84
      + g+s ++ + g+e ++ ++l+ag h f npq +p l +g l +s al +l e y + +l
Target 9 WRSQTSRTRIAL--ELKGLEYSRQEPVDLRAGAHKDEGPTAQNPGGLVPEVLeTPDGAQLSQSMALIEWLEETYPF-PALL 85
PP      5555554444443..45699*****94468999*****9986.5788

.....*
Query 85 akopkkravi 94
      p +ra+
Target 86 PACTTERAIA 95
PP      9888888875
    
```

This row shows the alignment between your sequence and the matching HMM.

- Query: consensus of the HMM (query sequence), coloured according to the match: Identical residues [█]; Similar residues [█];
- Match: the match between the query sequence (HMM) and target sequence
- Target: query sequence, coloured according to the posterior probability: 0% [█] 100%
- PP: posterior probability, or the degree of confidence in each individual aligned residue

Showing rows 28901 - 28945 of 28945



HMMER
Biosequence analysis using profile hidden Markov Models

Home Search **Results** Software Help About Contact

Next release within a week, think about downloading your results

PHMMER Results

Score	Taxonomy	Domain	Download
-------	----------	--------	----------

- **Job:** 76F84AE6-CAFF-11E8-88C2-7DBB53F04F9B.1
- **Started:** 2018-10-08 14:38:40
- **Algorithm:** phmmmer
- **HMMER Options:** -E 1 --domE 1 --incE 0.01 --incdomE 0.03 --mx BLOSUM62 --pextend 0.4 --popen 0.02 --seqdb uniprotkb

▼ **Format**

Text A plain text file containing the hit alignments and scores. 	Tab Delimited A tab delimited text file containing the hit information. No alignments. 	XML An XML file formatted for machine parsing of the data. 	JSON All the results information encoded as a single JSON string.
FASTA Download the significant hits from your search as a gzipped FASTA file. 	Full length FASTA A gzipped file containing the full length sequences for significant search hits. 	Aligned FASTA A gzipped file containing aligned significant search hits in FASTA format. 	STOCKHOLM Download an alignment of significant hits as a gzipped STOCKHOLM file.
ClustalW Download an alignment of significant hits as a gzipped ClustalW file. 	PSI-BLAST Download an alignment of significant hits as a gzipped psiblast file. 	PHYLIP Download an alignment of significant hits as a gzipped phylip file. 	



Paralell keresés

- A szekvenciák sorrendje egy adatbázisban esetleges
- Nem kell sorban haladni a keresés során
- Több processzoros, több magos architektúrán párhuzamosan lehet futtatni a keresést – GPU computing
- A párhuzamos futtatáshoz nemcsak a kódot kell átírni, de az algoritmust is

www.nvidia.com/object/tesla-supercomputing-solutions.html

NCBI BLAST < Sequence Si... x UVA FASTA Downloads x BLAST: Basic Local Alignm... x HMMER x High Performance Compu... x

Search NVIDIA USA - United States

DRIVERS ▾ PRODUCTS ▾ COMMUNITIES ▾ SUPPORT SHOP ABOUT NVIDIA ▾

TESLA WHAT IS GPU COMPUTING? GPU APPLICATIONS SERVERS AND WORKSTATIONS

NVIDIA Home > Products > High Performance Computing [Subscribe](#)

TESLA ACCELERATED COMPUTING

Your Platform for Insight and Discovery

<h3>WHAT IS ACCELERATED COMPUTING?</h3> <p>Learn how GPUs deliver significantly higher application performance ></p>	<h3>HUNDREDS OF APPLICATIONS</h3> <p>See the hundreds of applications accelerated with GPUs ></p>	<h3>GPUS FOR SERVERS & WORKSTATIONS</h3> <p>Accelerate your scientific computation with Tesla GPUs ></p>
<h3>DEVELOPING WITH GPU's</h3> <p>Get thousands of cores working for you ></p>	<h3>ACCELERATE YOUR CODE</h3> <p>Test drive a Tesla K80 GPU Accelerator ></p>	<h3>WHERE TO BUY</h3> <p>Find systems powered by Tesla GPUs ></p>

Ész Vent: Python Prog: Top 8 resour: Python Pand: Services and: MetaXpress®: Telefonk: Tájékozt: UVA FAS: RecNam: downlo: HelioBLAST: live: Itthon: K: GPU X

https://www.nvidia.com/en-us/data-center/gpu-accelerated-applicati

NVIDIA

DATA CENTER PRODUCTS SOLUTIONS APPS FOR DEVELOPERS TECHNOLOGIES

GPU APPLICATIONS CATALOG

HUNDREDS OF APPLICATIONS ACCELERATED

Find out if your application is being accelerated by NVIDIA GPUs. Today, hundreds of applications are already GPU-accelerated and the number is growing. See the list below.

Industry Product Category Keyword search



The screenshot shows the NVIDIA Data Center website. The browser address bar displays <https://www.nvidia.com/en-us/data-center/gpu-accelerated-applications>. The navigation menu includes DATA CENTER, PRODUCTS, SOLUTIONS, APPS, FOR DEVELOPERS, and TECHNOLOGIES. Below the navigation is a search bar with a dropdown menu set to 'All' and a filter set to 'Bioinformatics & Genomics'. The main content area displays a grid of 20 application cards, each with a title and a brief description:

SEQNFIND Accelerated Technology Laborat...	GHOST-Z GPU Akiyama_Laboratory, Tokyo Insti...	GPU-BLAST Carnegie Mellon University	SOAP3 Genomics
G-BLASTN Hong Kong Baptist University	ARIOC Johns Hopkins University	BEAGLE-LIB Open Source	CUDASW++ Open Source
CUSHAW Open Source	MUMMER GPU Open Source	NVBIO Open Source	NVBOWTIE Open Source
PEANUT Open Source	REACTA Open Source	WIDELM Open Source	MCUDA-MEME Open Source
SYNOMICS STUDIO Row Analytics	CAMPAIGN SimTK	SOAP3-DP The University of Hong Kong	UGENE Unipro



The screenshot shows the NVIDIA Data Center website with a grid of bioinformatics tools. A modal window for PEANUT is open, displaying the following information:

- PEANUT** (Open Source)
- * Achieves supreme sensitivity and speed compared to current state of the art read mappers like BWA MEM, Bowtie2 and RazerS3
- * PEANUT reports both only the best hits or all hits
- Read mapper for DNA or RNA sequence reads to a known reference genome.
- [View Quick Start Guide>](#)

The background grid includes tools such as SEQNFIND, G-BLASTN, CUSHAW, PEANUT, SYNOMICS STUDIO, MUMMER GPU, REACTA, CAMPAIGN, NVBIO, WIDELM, SOAP3-DP, NVBOWTIE, MCUDA-MEME, SOAP3, UDASW++, and UGENE.



The screenshot shows the NVIDIA Data Center website. The browser address bar displays <https://www.nvidia.com/en-us/data-center/gpu-accelerated-applications>. The page features a navigation menu with categories: DATA CENTER, PRODUCTS, SOLUTIONS, APPS, FOR DEVELOPERS, and TECHNOLOGIES. Below the navigation is a search bar with filters for 'All' and 'Molecular Dynamics', and a search input field containing 'Keyword search'. The main content area displays a grid of software options, each with a logo and a brief description:

- GROMACS**: FAST. FLEXIBLE. FREE. (Logo: a bird in flight). GROMACS. [View Quick Start Guide>](#)
- NAMD**: Scalable Molecular Dynamics. University of Illinois at Champaign Urbana. [View Quick Start Guide>](#)
- HTMD**: Acellera Ltd.
- ACEMD**: Acellera Ltd.
- GPUGRID.NET**: Acellera Ltd.
- DESMOND**: David E. Shaw Research.
- ESPRESSO**: ESPResSo.
- HALMD**: HALMD.
- CHARMM**: Harvard University.
- MYPRESTO**: N2PC/AIST/JBIC, Japan.
- GENESIS**: RIKEN.
- SOP-GPU**: SOP-GPU.
- LAMMPS**: Sandia National Lab.
- OPENMM**: Stanford University.



The screenshot shows the NVIDIA Data Center website. The browser address bar displays <https://www.nvidia.com/en-us/data-center/gpu-accelerated-applications>. The NVIDIA logo is at the top left. A navigation bar contains: DATA CENTER, PRODUCTS, SOLUTIONS, APPS, FOR DEVELOPERS, and TECHNOLOGIES. Below the navigation bar is a search area with a dropdown menu set to 'All', a search box containing 'Molecular visualization and Doc...', and a search button. The main content area features a grid of application cards:

- MEGADOCK**
- BINDSURF** (Bioinformatics and High Perfor...)
- PIPER PROTEIN DOCK...** (Boston University)
- BUDE** (Bristol University Docking Station)
- FASTROCS** (Open Eye Scientific Software, Inc.)
- MOLEGRO VIRTUAL D...** (QIAGEN)
- PYMOL** (Schrodinger, Inc.)
- AMIRA** (Thermo fisher Scientific)
- VEGA ZZ** (University of California, San Fra...)
- VMD** (University of Illinois)
- INTERACTIVE MOLEC...** (University of Illinois)

At the bottom, a green banner reads: **DOWNLOAD CATALOG- LAST UPDATED MARCH 2018**. Below this, it says: 'Download the complete list of GPU-accelerated applications.' with an upward-pointing arrow icon.



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/>



Előadások | Élettani Intézet x Központi Könyvtár - Tudásbázis x UVA FASTA Downloads x NCBI BLAST < Sequence Si... x HelioBLAST by HelioText x +

helioblast.heliotext.com

HelioBLAST

Home HelioBLAST

Ask HelioBLAST

Search in: Medline Would you like to add your own database? [Get in touch with us.](#)

HelioBLAST is a free service provided by HelioText. The HelioBLAST text similarity engine finds text records that are similar to the submitted query. HelioBLAST uses our super-fast search engine, and this service is provided to demonstrate what can be done using text similarity searching. HelioBLAST can be customized to search a particular database or multiple ones; and proprietary databases can be created for individual clients.

[Submit to HelioBLAST](#)

MEDLINE index last updated Friday, October 07, 2016.

NLM data are produced by a U.S. Government agency and include works of the United States Government that are not protected by U.S. copyright law but may be protected by non-US copyright law, as well as abstracts originating from publications that may be protected by U.S. copyright law.

NLM assumes no responsibility or liability associated with use of copyrighted material, including transmitting, reproducing, redistributing, or making commercial use of the data. NLM does not provide legal advice regarding copyright, fair use, or other aspects of intellectual property rights. Persons contemplating any type of transmission or reproduction of copyrighted material such as abstracts are advised to consult legal counsel.

Copyright HelioText 2015. email:kmenier@heliotext.com. Phone: +1-214-924-6334.

Welcome to HelioBLAST

HelioBLAST is a free service provided by HelioText. The HelioBLAST text similarity engine finds text records that are similar to the submitted query. HelioBLAST uses our super-fast search engine, and this service is provided to demonstrate what can be done using text similarity searching. HelioBLAST can be customized to search a particular database or multiple ones; and proprietary databases can be created for individual clients.

Here, your query is searched against the citations (abstract and titles) in Medline/PubMed. **Submissions are limited to 1,000 words**, so we recommend your query should be an abstract or paragraph.



Előadások | Élettani Intézet x Központi Könyvtár - Tudásbázis x UVA FASTA Downloads x NCBI BLAST < Sequence Si... x HelioText | Liberate your insight... x +

helioblast.heliotext.com/search

HelioText

Home HelioBLAST

HelioBLAST results similar to your query

The 50 best matches found by HelioBLAST in 2.829 sec:

Algorithms for recollection of search terms based on the Wikipedia category structure. Score:0.067
 by Vandamme, Stijn; De Turck, Filip
in TheScientificWorldJournal (2014)

Abstract: The common user interface for a search engine consists of a text field where the user can enter queries consisting of one or more keywords. Keyword query based search engines work well when the users have a clear vision what they are looking for and are capable of articulating their query using the same terms as indexed. For our multimedia ... [More>>](#)

Medline PMID: [24616630](#)

Evaluating Open-Source Full-Text Search Engines for Matching ICD-10 Codes. Score:0.066
 by Jurcu, Daniel-Alexandru; Stoicu-Tivadar, Vasile
in Studies in health technology and informatics (2016)

Abstract: This research presents the results of evaluating multiple free, open-source engines on matching ICD-10 diagnostic codes via full-text searches. The study investigates what it takes to get an accurate match when searching for a specific diagnostic code. For each code the evaluation starts by extracting the words that make up its text and continues with building full-text search queries from ... [More>>](#)

Medline PMID: [27350484](#)

Identifying duplicate content using statistically improbable phrases. Score:0.063
 by Errami, Mounir; Sun, Zhaohui; George, Angela C; Long, Tara C; Skinner, Michael A; Wren, Jonathan D; Garner, Harold R
in Bioinformatics (Oxford, England) (2010)

Abstract: Document similarity metrics such as PubMed's 'Find related articles' feature, which have been primarily used to identify studies with similar topics, can now also be used to detect duplicated or potentially plagiarized papers within literature reference databases. However, the CPU-intensive nature of document comparison has limited MEDLINE text similarity studies to the comparison of abstracts, which constitute only a small ... [More>>](#)

Analysis Tools

A few tools to do more:

Find An Expert

Experts are potential reviewers, collaborators or competitors. Experts are identified from their publication history in this search.

1. Garner, Harold R ★★ ★
2. Lyu, Ping-Chiang ★★ ★
3. Errami, Mounir ★★ ★
4. Hanauer, David A ★★ ★
5. Sternberg, Paul W ★★ ★
6. Lo, Wei-Cheng ★★ ★
7. Gibney, Gretchen ★★ ★

Implicit Keywords

Implicit Keywords help identify concepts that were not originally mentioned in the query. Words are extracted from the 50 best matches found by HelioBLAST.

[View Implicit Keywords>>](#)



Előadások | Élettani Intézet x Központi Könyvtár - Tudásbázis x UVA FASTA Downloads x NCBI BLAST < Sequence Si... x HelioText | Liberate your insight... x

helioblast.heliotext.com/search

HelioText

Home HelioBLAST

HelioBLAST results

The 50 best matches found by HelioText

Algorithms for recollection structure.
by Vandamme, Stijn; De Turck, Filip in *TheScientificWorldJournal* (2014)
Abstract: The common user interface for search engines consists of one or more search boxes. The users have a clear vision what they want to search for. The same terms as indexed. For more information see the full text.
Medline PMID: [24616630](#)

Evaluating Open-Source Full-Text Search Engines.
by Jurcu, Daniel-Alexandru; Stoicu-Tivadar, Daniela in *Studies in health technology and informatics* (2010)
Abstract: This research presents a method for evaluating full-text search engines. The method is based on ICD-10 diagnostic codes via full-text search. The method is used to match when searching for a specific term. The method is used to match the words that make up its text.
Medline PMID: [27350484](#)

Identifying duplicate content in a large database.
by Errami, Mounir; Sun, Zhaohui; George, George; Harold R. in *Bioinformatics (Oxford, England)* (2010)
Abstract: Document similarity measures are primarily used to identify studies that are potentially plagiarized papers with similar content. Document comparison has limitations. Document comparison has limitations which constitute only a small ...

Implicit Keyword | Average Frequency

information	1.0
data	0.9
system	0.8
searches	0.7
users	0.7
use	0.6
engines	0.6
web	0.6
used	0.6
user	0.6
queries	0.6
also	0.6
image	0.5
two	0.5
biomedical	0.5
one	0.5
articles	0.5
between	0.4
available	0.4
based	0.4
more	0.4
images	0.4
literature	0.4
all	0.4
medical	0.4
concepts	0.4
pubmed	0.4

Implicit Keywords [close](#)

Analysis Tools

tools to do more:

An Expert
Experts are potential reviewers, collaborators or competitors. Experts are identified from their citation history in this search.

Erner, Harold R ★★★
Lu, Ping-Chiang ★★★
Rami, Mounir ★★
Sauer, David A ★★
Sternberg, Paul W ★★
Wei-Cheng ★★
Woney, Gretchen ★★

Implicit Keywords
Implicit Keywords help identify keywords that were not originally mentioned in the query. Words are selected from the 50 best matches found by HelioBLAST.
[View Implicit Keywords>>](#)



Mit tanultunk ma?

- Az adatbázis keresés lényegében nagyléptékű szekvenciaillesztés.
- Nagyon kiforrott technika.
- Gyakran a bioinformatikai vizsgáldás 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.