

Lesson 7 Database Searching by Keyword and Motif

Assignment

1. (Before class) Read GCG tutorial chapters 6 and 10.
2. (After class) Read the section in the GCG Program manual, "Defining patterns", which occurs in the Motifs and Findpatterns chapters.
3. <http://www.ncbi.nlm.nih.gov/entrez/query/static/help/helpdoc.html>
http://www.ncbi.nlm.nih.gov/BLAST/blast_FAQs.html#Short

Summary of Commands:

Note: In this document different fonts have different meanings:

Times is used to explain commands.

Courier is used to indicate commands and command options.

Courier italics are used to indicate command parameters, for example, filenames.

Courier bold is used to indicate commands that are not displayed.

Courier bold italics are used to indicate computer-generated output.

Helvetica is used to indicate menu items.

stringsearch

Searches for text strings in annotation of sequence libraries.

,
stringsearch -match=or

Stringsearch AND operator. Does Stringsearch search with ",", standing for OR.

lookup

Also searches for text strings in annotation of sequence libraries, but is faster and somewhat more flexible, albeit somewhat less sensitive, than Stringsearch. Specifies GCG list file.

@filename

Boolean operators used with Lookup:

&

AND

Exempli gratia:

Chromosome 13.

Chromosome & 13

OR

|

F'rinstance:

13 | 13q

"13" or "13q".

To perform an Entrez search:
Click on <http://www.ncbi.nlm.nih.gov/>
Click on Entrez

Boolean operators used with Entrez:

AND

Chromosome AND 13

OR

13 OR 13q

<http://srs6.ebi.ac.uk/>

motifs

motifs -freq

findpatterns

findpatterns -nomon

findpatterns -mism=1

Pattern codes (GCG convention):

X

N

R

Y

(symbol(s)){#}

TAAC{5}GT

TAA(AT){3}GT

(symbol(s)){#1,#2}

TAA(AT){2,6}GT

(symbol(s)){#, }

AND

Exempli gratia:

Chromosome 13.

OR

F'instance:

"13" or "13q".

Performs an SRS search.

Compares sequence to
Prosit library of protein
structure motifs. Excludes
frequently occurring motifs.
Compares sequence to
Prosit library of protein
structure motifs. Includes
frequently occurring motifs.
Searches library for exact
agreement with one or more
sequence patterns.
Searches library for exact
agreement with one or more
sequence patterns but
supresses screen output of
sequences searched.
Searches library for
agreement with one or more
sequence patterns, allowing
one mismatch.

Any amino acid.

Any nucleotide.

Purines (A,G).

Pyrimidines (T,C).

Repeats symbols (#) times.

Exempli gratia:

TAACCCCGT

TAAATATATGT

Repeats symbol(s) between
#1 and #2 times.

Exempli gratia:

"TAA" then "AT" repeated
between 2 and 6 times, then
"GT".

Repeats symbol(s) for at
least # times.

Exempli gratia:

CATG{2, }A

(symbol(s)){, # }

ACT{,10 }A

(symbol1, symbol2...)

F(Q,A)S

~(symbol(s))

GC~ACG

GC~(A,T)CG

<symbols

<ATCG

symbols>

ATCG>

PHIBLAST

http://www.ncbi.nlm.nih.gov/BLAST/PHI-BLAST_search

Pattern codes in the NCBI convention are the same as above except as illustrated by the following examples

F[QA]S

“CAT” then “G” repeated at least two times, then “A”.

Repeats symbol(s) for between 0 and # times.

F’rinstance:

“AC” then “T” repeated between 0 and 10 times, then “A”.

Either symbol1, symbol2, etc..

Exempli gratia:

“F”, then “Q” or “A”, followed by “S”.

Any symbol but the one following.

Exempli gratia:

“GC”, any symbol other than “A”, then “CG”.

“GC”, any symbol other than “A” or “T”, CG.

Symbols occurring at the beginning of a sequence only.

Exempli gratia:

ATCG occurring at the beginning of a sequence only.

Symbols occurring at the end of a sequence only.

Exempli gratia:

ATCG occurring at the end of a sequence only.

The user specifies a protein query sequence and a protein motif. PHIBLAST (Pattern Hit Initiated BLAST) searches only those members of a protein database that contains the specified motif. This technique reduces the number of sequences being searched and hence raises the sensitivity of the search. To run PHIBLAST go to:

and click on:

“F”, then “Q” or “A”, followed by “S”.

Also in NCBI format, repeating symbols must be written out explicitly:, i.e. :

TAACCCCGT

Web sites:

<http://expasy.hcuge.ch/sprot/prosite.html>
Prosite

<http://www.ncbi.nlm.nih.gov/Entrez/>
Entrez

<http://srs6.ebi.ac.uk/>
SRS

http://www.ncbi.nlm.nih.gov/cgi-bin/BLAST/nph-psi_blast?Jform=1
Phiblast

Lab or homework:

1. Work through the GCG tutorial chapter 6.1, using Stringsearch, Lookup, and Entrez. Suggestion: since Stringsearch takes a long time to run, Stringsearch in one session, and simultaneously, lookup in another.

2. (optional) Obtain the sequence file of aardvark myoglobin from Swissprot. Do your searching with:

- A. Lookup.
- B. Stringsearch.
- C. Entrez.

3. Do the exercises in the GCG tutorial chapter 10.

4. Find the rare motifs in Chicken SRC protein. Find the common motifs in the same protein. Convince yourself that the uncommon motif found is indeed matched by the protein. If you have trouble doing this on the screen, print the output file out (either from cuccfa or from your own Mac or PC) and examine the match on paper. As there may be delays in getting a printed copy to work with in the lab, the pattern-match is reproduced below:

Protein_Kinase_Atp

```
(L,I,V)G~(P)G~(P)(F,Y,W,M,G,S,T,N,H)(S,G,A)~(P,W)(L,I,V,C,A,T)~(P,D)x(G,S,T,A,C,L,I,V,M,F,
Y)x{5,18}(L,I,V,M,F,Y,W,C,S,T,A,R)(A,I,V,P)(L,I,V,M,F,A,G,C,K,R)K
(L)G~PG~P(F)(G)~(P,W)(V)~(P,D)x(G)x{7}(V)(A)(I)K
272: RLEVK LGQGCFGEVWMGTWNGTTRVAIK TLKPG
```

5. (optional) Do the exercises in the GCG tutorial chapter 6.2.

6. Translate the following into motif notation: “A peptide begins with an alanine, a glycine, and an arginine, followed by between 20 and 30 amino acids of any kind, followed by a glutamine and arginine pair repeated between 15 and 92 times and terminating in a tryptophan”.

7. Obtain the file ced4.tfa from /usr2/seq/seqclass/lab7 . Perform a Blastp and a Phiblast search with this sequence and the pattern (G,A)x4GK(S,T). What sequence has a marginal hit with Phiblast but not with Blastp?

Bibliography

1. GCG manual on Stringsearch, Lookup, Motifs, and Findpatterns.