

Multiple sequence alignment

Monday, December 6, 2010

Bioinformatics
J. Pevsner
pevsner@kennedykrieger.org

Multiple sequence alignment: today's goals

- to define what a multiple sequence alignment is and how it is generated; to describe profile HMMs
- to introduce databases of multiple sequence alignments
- to introduce ways you can make your own multiple sequence alignments
- to show how a multiple sequence alignment provides the basis for phylogenetic trees

Page 179

Multiple sequence alignment: outline

[1] Introduction to MSA

Exact methods
Progressive (ClustalW)
Iterative (MUSCLE)
Consistency (ProbCons)
Structure-based (Expresso)
Conclusions: benchmarking studies

[2] Hidden Markov models (HMMs), Pfam and CDD

[3] MEGA to make a multiple sequence alignment

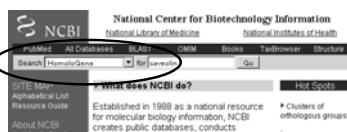
[4] Multiple alignment of genomic DNA

Multiple sequence alignment: definition

- a collection of three or more protein (or nucleic acid) sequences that are partially or completely aligned
- homologous residues are aligned in columns across the length of the sequences
- residues are homologous in an evolutionary sense
- residues are homologous in a structural sense

Page 180

Example: someone is interested in caveolin



Step 1: at NCBI change the pull-down menu to HomoloGene and enter caveolin in the search box

Step 2: inspect the results. We'll take the first set of caveolins. Change the Display to Multiple alignment.

Gene	Accession	Species
1: HomoloGene 1330: Gene conserved in Eutherians	CAV1	caveolin 1, caveolin protein, 22kDa
	CAV1	caveolin 1, caveolin protein, 22kDa
	CAV1	caveolin 1, caveolin protein, 22kDa
	CAV1	caveolin 1, caveolin protein, 22kDa
2: HomoloGene 7295: Gene conserved in Eutherians	CAV3	caveolin 3
	CAV3	caveolin 3
	CAV3	caveolin 3
	CAV3	caveolin 3
	CAV3	caveolin 3
	CAV3	caveolin 3
	CAV3	caveolin 3
	CAV3	caveolin 3

Multiple sequence alignment: outline

[1] Introduction to MSA

Exact methods
 Progressive (ClustalW)
 Iterative (MUSCLE)
 Consistency (ProbCons)
 Structure-based (Expresso)
 Conclusions: benchmarking studies

[2] Hidden Markov models (HMMs), Pfam and CDD

[3] MEGA to make a multiple sequence alignment

[4] Multiple alignment of genomic DNA

[5] Introduction to molecular evolution and phylogeny

Multiple sequence alignment: methods

Progressive methods: use a guide tree (related to a phylogenetic tree) to determine how to combine pairwise alignments one by one to create a multiple alignment.

Examples: CLUSTALW, MUSCLE

Multiple sequence alignment: methods

Example of MSA using ClustalW: two data sets

Five distantly related globins (human to plant)

Five closely related beta globins

Obtain your sequences in the FASTA format!
 You can save them in a Word document or text editor.
 Visit www.bioinfbook.org for web documents 6-3 and 6-4

Feng-Doolittle MSA occurs in 3 stages

[1] Do a set of global pairwise alignments
 (Needleman and Wunsch's dynamic programming algorithm)

[2] Create a guide tree

[3] Progressively align the sequences

Progressive MSA stage 1 of 3: generate global pairwise alignments

SeqA Name	Len(aa)	SeqB Name	Len(aa)	Score
1 beta_globin	147	2 myoglobin	154	25
1 beta_globin	147	3 neuroglobin	151	15
1 beta_globin	147	4 soybean	144	13
1 beta_globin	147	5 rice	166	21
2 myoglobin	154	3 neuroglobin	151	16
2 myoglobin	154	4 soybean	144	8
2 myoglobin	154	5 rice	166	12
3 neuroglobin	151	4 soybean	144	17
3 neuroglobin	151	5 rice	166	18
4 soybean	144	5 rice	166	43

best score

Number of pairwise alignments needed

For n sequences, $(n-1)(n) / 2$

For 5 sequences, $(4)(5) / 2 = 10$

For 200 sequences, $(199)(200) / 2 = 19,900$

Page 185

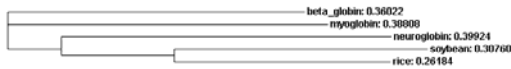
Feng-Doolittle stage 2: guide tree

- Convert similarity scores to distance scores
- A tree shows the distance between objects
- Use UPGMA (defined in the phylogeny lecture)
- ClustalW provides a syntax to describe the tree

Page 187

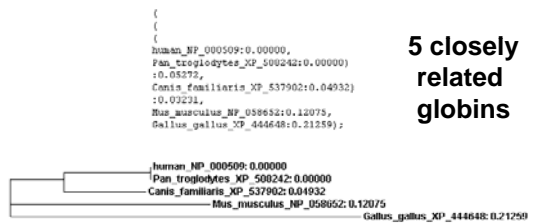
Progressive MSA stage 2 of 3: generate a guide tree calculated from the distance matrix (5 distantly related globins)

```
{
  beta_globin:0.36022,
  myoglobin:0.38808,
  {
    neuroglobin:0.39924,
    {
      soybean:0.30760,
      rice:0.26184
    }
  }
  :0.13652)
:0.06560);
```



Page 186

Seq#	Name	Len(aa)	Seq#	Name	Len(aa)	Score
1	human_NP_000509	147	2	Pan_troglodytes_XP_508242	147	100
1	human_NP_000509	147	3	Canis_familiaris_XP_537902	147	89
1	human_NP_000509	147	4	Mus_musculus_NP_058652	147	80
1	human_NP_000509	147	5	Gallus_gallus_XP_444648	147	69
2	Pan_troglodytes_XP_508242	147	3	Canis_familiaris_XP_537902	147	89
2	Pan_troglodytes_XP_508242	147	4	Mus_musculus_NP_058652	147	80
2	Pan_troglodytes_XP_508242	147	5	Gallus_gallus_XP_444648	147	69
3	Canis_familiaris_XP_537902	147	4	Mus_musculus_NP_058652	147	78
3	Canis_familiaris_XP_537902	147	5	Gallus_gallus_XP_444648	147	71
4	Mus_musculus_NP_058652	147	5	Gallus_gallus_XP_444648	147	66



5 closely related globins

Page 188

Feng-Doolittle stage 3: progressive alignment

- Make a MSA based on the order in the guide tree
- Start with the two most closely related sequences
- Then add the next closest sequence
- Continue until all sequences are added to the MSA
- Rule: "once a gap, always a gap."

Page 188

Clustal W alignment of 5 distantly related globins

CLUSTAL W (1.83) multiple sequence alignment

```
beta_globin      -----MVLITPEEKSAVTAIAG--KVVNDEVGGEALGRLLVYVFTORFF 43
cytoglobin      MEKVFQEMIEERFSEELSEAEKAVQAMARLYANCDVGOAIVRFFVVFSAKQVF 60
myoglobin      -----HGLSDGEWQLVNWVKYVADIPHGQEVLRFGHGFETLEKF 44
neuroglobin     -----HEPFEFLRQSRVAVSSPFLHGTVLFARLFALEPFLPLF 42
leghemoglobin   -----MGFTFEQALVNSNEIFKQNFV--YSLVPTTILKQAPAKQHF 43
                : : : * . . . : : * *
beta_globin     E5-FGLSTPDAVHGMFFVFAHGRKVLGAFSDGLA---HLDMRGTFATLSELHCDRIHV 99
cytoglobin      S3-FKMEPELEKESFQIGKACVNHGALHIVVHLEDFKIVSSTLALVGMALKEIV 119
myoglobin      DK-FEHLKSEDEKASEDLKSHGATVLTALGGLK---KGGHAEKIKPLAQSHATKIKI 100
neuroglobin     QVNCQFQSPFECCLSPFLSHIDVNHVLDAAVTVPELSELEEVLAELGDSHRAVQ V 101
leghemoglobin   S---FLKDSAEVYDFPLQAHAEKVFQVHVDISAIQLRASGEVYVLDGATLGAITHIQGVV 99
                : . . * * * : . . :
beta_globin     DPENFRLLGNVLYCVLAIHFQEFPPVQAAYQKVVAGVAMALAHKTH----- 147
cytoglobin      EPVFFKILSGVILEVVAEFAFDFFPETQANAKLGLIYSHVTAAYKVGWVQVPMAT 179
myoglobin      FFKYLEFSEICIIIVLQSKRFGDFADAGANRALELFRGASNTYKELQFQG----- 154
neuroglobin     KLSSTVYQESLLVMLKELGHPAFAFRANGLVYGVYQANSGRDE----- 151
leghemoglobin   DP-HVYVYKALELTFEASGEKNSSEELTANVAYEGLASIAEKAMH----- 146
                : : : : * . . . :
beta_globin     -----
cytoglobin      TTPATLFSGGP 190
myoglobin      -----
neuroglobin     -----
leghemoglobin   -----
```

Fig. 6.3
Page 187

Clustal W alignment of 5 closely related globins

```

CLUSTAL W (1.83) multiple sequence alignment

Human_Hp_000509      NVHITPEEKSAVTALMGKVWVDEVGGELGRLLVVFTQRFDFSGDLS 50
Pan_trogodytes_XP_508242 NVHITPEEKSAVTALMGKVWVDEVGGELGRLLVVFTQRFDFSGDLS 50
Canis_familiaris_XP_537902 NVHITAEKSLVSLGKQVWVDEVGGELGRLLVVFTQRFDFSGDLS 50
Mus_musculus_NF_058652 NVHITAEKSAVSLGKQVWVDEVGGELGRLLVVFTQRFDFSGDLS 50
Gallus_gallus_XP_444648 NVHITAEKSLITLQKQVWVDEVGGELGRLLVVFTQRFDFSGDLS 50
*** * * * * : : * : : * : * : : * : : : * : : * : :

Human_Hp_000509      TFDVNGNPKYKMHGKQVLAFTDGLAHLELNGKGFATLSELMCKLHVD 100
Pan_trogodytes_XP_508242 TFDVNGNPKYKMHGKQVLAFTDGLAHLELNGKGFATLSELMCKLHVD 100
Canis_familiaris_XP_537902 TFDVNSMAYKMHGKQVLAFTDGLAHLELNGKGFATLSELMCKLHVD 100
Mus_musculus_NF_058652 SASALRQPKYKMHGKQVLAFTDGLAHLELNGKGFATLSELMCKLHVD 100
Gallus_gallus_XP_444648 SFTALGNPKYKMHGKQVLAFTDGLAHLELNGKGFATLSELMCKLHVD 100
I : * : * : * : * : * : : * : : : * : : : * : : * : :

Human_Hp_000509      FENFLLGNLVLCVLAHSPGKFTFPAQAYQKVVAVAHALAKHTH 147
Pan_trogodytes_XP_508242 FENFLLGNLVLCVLAHSPGKFTFPAQAYQKVVAVAHALAKHTH 147
Canis_familiaris_XP_537902 FENFLLGNLVLCVLAHSPGKFTFPAQAYQKVVAVAHALAKHTH 147
Mus_musculus_NF_058652 FENFLLGNLVLCVLAHSPGKFTFPAQAYQKVVAVAHALAKHTH 147
Gallus_gallus_XP_444648 FENFLLGNLVLCVLAHSPGKFTFPAQAYQKVVAVAHALAKHTH 147
**** * : : * : * : * : * : * : * : * : * : * : * :
    
```

* asterisks indicate identity in a column

Fig. 6.5
Page 189

Why "once a gap, always a gap"?

- There are many possible ways to make a MSA
- Where gaps are added is a critical question
- Gaps are often added to the first two (closest) sequences
- To change the initial gap choices later on would be to give more weight to distantly related sequences
- To maintain the initial gap choices is to trust that those gaps are most believable

Additional features of ClustalW improve its ability to generate accurate MSAs

- Individual weights are assigned to sequences; very closely related sequences are given less weight, while distantly related sequences are given more weight
- Scoring matrices are varied dependent on the presence of conserved or divergent sequences, e.g.:

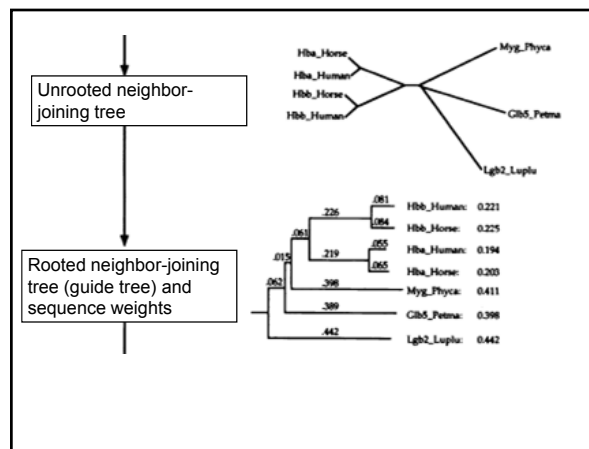
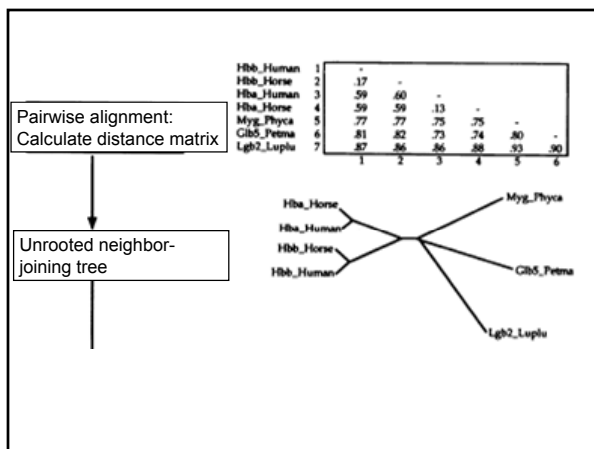
PAM20	80-100% id
PAM60	60-80% id
PAM120	40-60% id
PAM350	0-40% id

- Residue-specific gap penalties are applied



Figure 1. The basic progressive alignment procedure, illustrated using a set of 7 globin of horse variety structure. The sequence names are from those that ClustalW uses: Hba_Human, Hbb_Human, Hbb_Horse, Hbb_Horse, Hba_Horse, Myg_Physa, Cgb5_Petma, Lgb2_Luplu, Hba_Horse, Hbb_Human, Hbb_Horse, Hba_Human, Myg_Physa, Cgb5_Petma, Lgb2_Luplu. In the distance matrix, the mean number of differences per residue.

See Thompson et al. (1994) for an explanation of the three stages of progressive alignment implemented in ClustalW



ProbCons uses an HMM to make alignments

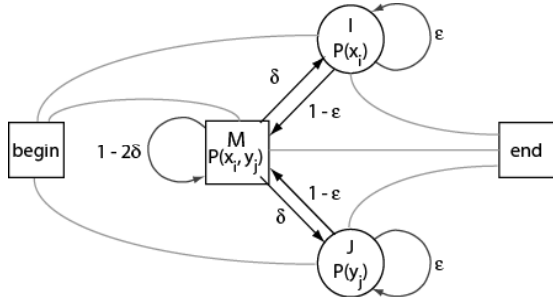


Fig. 5.12
Page 158

ProbCons—consistency-based approach

Sequence x x_i
 Sequence y y_j
 Sequence z z_k

If x_i aligns with z_k
 and z_k aligns with y_j
 then x_i should align with y_j

ProbCons incorporates evidence from multiple sequences to guide the creation of a pairwise alignment.

Page 193

ProbCons output for the same alignment:
 consistency iteration helps

```
(c)
PROBCONS
beta globin M-----VHLTPEEKSAVTLAKGKVNVD--EVGGEALGRLLVYPWQRFES-FG
myoglobin M-----GLSDGQWQLVLYWKGVEADI FGGQEVLI RLFKGFPELTKFK-FK
neuroglobin M-----ERFPELIGQRWAVRSPLEHGTVLPARLFALEPDLPLFOYMKR
soybean M-----VAPFEEGQALVSSFEAPKANI PCVIVVYVFTLEKAPAKGKFFP-LA
rice MALVEDNNAVAVFSSEKQALVYKSHALKKDSANIALRFLKIFEVAPSASQMFSLR
* : : : : : : : : : : : *
beta globin DLSTFDAMNPKFVAAGKRVLGAFSDGAAHLD---NRK---GTFATISELECDKLAADP
myoglobin HIKKEDDMASIEDRGGATVLTALGDI---LREKSHH---AEIPLAQSRAPYRHPV
neuroglobin QFSPDCLSSPEFLDHRVNLVIDAATNVDLSLE---EVLALGRHRVA-GKL
soybean MVDP---TNFKLGHAEKLFALVRSAGQLKAGTVV---ADAAIGSVRAOK-ATD
rice NSDVP--LEKNFKLTHAMSVFVNTCEAAQLRKAGKVTVRDTTLERLGATHLEY-GVGD
. : . . . : : : : : : : : : :
beta globin ENFRLLGNLVCLAHNF-GKEFPFPVQAAYQKVIAGVANALAK-----YH
myoglobin KYLEFISECIIGVQSRH-PDGFADAGGAMNKALELFRKIMANFKELGFGQ
neuroglobin SFFSTVQESLLYMLKCL-GPAFTPATRAAMSQLYGAVVAQMSRG--W-DGE
soybean POFVVVEALKETKAIV-GDWSDLSRAKEVAYDELAALAK-----KA
rice AHFVVVFALDDTTRKEVPADMSGPAKMSAWSEAYDHLVAALKE---MKPAE
: : : : : : : : : : : : : : :
```

Multiple sequence alignment: outline

- [1] Introduction to MSA
 - Exact methods
 - Progressive (ClustalW)
 - Iterative (MUSCLE)
 - Consistency (ProbCons)
 - Structure-based (Expresso)
 - Conclusions: benchmarking studies
- [2] Hidden Markov models (HMMs), Pfam and CDD
- [3] MEGA to make a multiple sequence alignment
- [4] Multiple alignment of genomic DNA

Access to Tcoffee: <http://tcoffee.org>

- Make a MSA
- MSA w. structural data
- Compare MSA methods
- Make an RNA MSA
- Combine MSA methods
- Consistency-based
- Structure-based
- Back translate protein MSA

Page 194

APDB ClustalW output:
 Tcoffee can incorporate structural information into a MSA

```
T-COFFEE, Version 4.71(Thu Nov 16 15:00:43 2006)
Cubic Retriever
CPU TIME:0 sec.
# APDB Evaluation: Color Range Blue [0 % -- 100 %]-Red
# Sequence Score: APDB
# Local Score: APDB

SCORE=47
*
EAD AV: 6000
2hbb : 224
1W5A : 213
2h81 : 210
10JGA : 194
1FSL : 157
2hbb -----HVVLEKSRKFTALRS--LHDEVSGEISQELVIVV
1W5A HEWVQNEIEEERREELITLGGKSRKSRKSRKSRKSRKSRKSRKSRKSRK
2h81 -----HQLSDQWQLVLYWKGVEADIFGGQEVLI RLFKGFPELTKFK
10JGA -----HEPELIGQRWAVRSPLEHGTVLPARLFALEPDLPLFOYMKR
1FSL -----MALVEDNNAVAVFSSEKQALVYKSHALKKDSANIALRFLKIFEVAPSASQMFSLR
```

Protein Data Bank accession numbers

Multiple sequence alignment: outline

- [1] Introduction to MSA
 - Exact methods
 - Progressive (ClustalW)
 - Iterative (MUSCLE)
 - Consistency (ProbCons)
 - Structure-based (Expresso)
- [2] Hidden Markov models (HMMs), Pfam and CDD
- [3] MEGA to make a multiple sequence alignment
- [4] Multiple alignment of genomic DNA

Multiple sequence alignment: methods

How do we know which program to use?

There are benchmarking multiple alignment datasets that have been aligned painstakingly by hand, by structural similarity, or by extremely time- and memory-intensive automated exact algorithms.

Some programs have interfaces that are more user-friendly than others. And most programs are excellent so it depends on your preference.

If your proteins have 3D structures, **use these** to help you judge your alignments. For example, try Expresso at <http://www.tcoffee.org>.

Page 196

Strategy for assessment of alternative multiple sequence alignment algorithms

- [1] Create or obtain a database of protein sequences for which the 3D structure is known. Thus we can define "true" homologs using structural criteria.
- [2] Try making multiple sequence alignments with many different sets of proteins (very related, very distant, few gaps, many gaps, insertions, outliers).
- [3] Compare the answers.

Page 196

BaliBase: comparison of multiple sequence alignment algorithms

```
Name          hiv-1 protease
Number of sequences      4
Alignment Length       106
Longest Sequence        104
Shortest Sequence        90
Average Percent Identity 49
Maximum Percent Identity 86
Minimum Percent Identity 35

Sequence Name  SWISSPROT Accession
1nb3          P32542
7upj3         P03266
pol_sivce    P17083
POL_SIVK     P05897

Family 1nb3 7upj3 pol_sivce POL_SIVK

1nb3 1  VTRLEKRPFTIVLINDTLNVLSTGADTSLTThvYrIkYvqK.YQ
7upj3 1  pQSLMRKPPVTAIIEGQFVILLSTGADDIIVAG...1e1.gm.YE
pol_sivce 1  pQTLRQKPLPQVYQKCEALLSTGADDTVIE...1e1.qm.HE
POL_SIVK 1  pQSLMRKPPVTAIIEGQFVILLSTGADDIIVT...1e1.gm.YE

1nb3 50  GTGIGQVQVQVTFP...TFVTKQKGRHRTSLVADIPVTLGGDILQDL
7upj3 44  PFTVQIQGF INLELVNVEIVLNRKVPATITGDTPIINIFGRHILTAI
pol_sivce 44  PERIQIQGF INKQVWHPVRIIEGRFPVTVLQVPTPNIIGRHILTAI
POL_SIVK 44  PFTVQIQGF INKQVWHPVRIIEGRFPVTVLQVPTPNIIGRHILTAI

1nb3 99  GATLV
7upj3 94  GRLSL
pol_sivce 94  QCLV
POL_SIVK 94  GRLSL

Key
Alpha helix  RED
Beta strand  GREEN
Coiled coil  UNDERSCORE
```

Page 196

Multiple sequence alignment: methods

Benchmarking tests suggest that ProbCons, a consistency-based/progressive algorithm, performs the best on the BALiBASE set, although MUSCLE, a progressive alignment package, is an extremely fast and accurate program.

ClustalW is the most popular program. It has a nice interface (especially with ClustalX) and is easy to use. But several programs perform better. There is no one single best program to use, and your answers will certainly differ (especially if you align divergent protein or DNA sequences)

Page 196

Multiple sequence alignment: review questions

- [1] Explain how ClustalW works.
- [2] Name two alternative methods. How do they work? What is the evidence that they are better than ClustalW?
- [3] What does it mean to do a benchmarking study, and why is it important?
- [4] Why is it important that various programs that make MSAs often give different answers?

Multiple sequence alignment: outline

- [1] Introduction to MSA
 - Exact methods
 - Progressive (ClustalW)
 - Iterative (MUSCLE)
 - Consistency (ProbCons)
 - Structure-based (Expresso)
 - Conclusions: benchmarking studies
- [2] Hidden Markov models (HMMs), Pfam and CDD
- [3] MEGA to make a multiple sequence alignment
- [4] Multiple alignment of genomic DNA

Multiple sequence alignment to profile HMMs

- ▶ Hidden Markov models (HMMs) are "states" that describe the probability of having a particular amino acid residue at arranged in a column of a multiple sequence alignment
- ▶ HMMs are probabilistic models
- ▶ HMMs may give more sensitive alignments than traditional techniques such as progressive alignment

Page 197

Structure of a hidden Markov model (HMM)

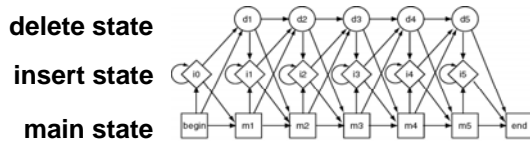


Fig. 5.12
Page 158

Pfam (protein family) database is a leading resource for the analysis of protein families

<http://pfam.sanger.ac.uk/>

The screenshot shows the Pfam website interface. At the top, there are navigation links: HOME | SEARCH | BROWSE | FTP | HELP. The main heading is 'Pfam 24.0 (October 2009, 11912 families)'. Below this, there is a description: 'The Pfam database is a large collection of protein families, each represented by multiple sequence alignments and hidden Markov models (HMMs). More...'. There are several 'QUICK LINKS' and 'VIEW A...' options, along with a 'JUMP TO' field for searching by accession ID.

Page 198

PFAM HMM for lipocalins: resembles a position-specific scoring matrix

PFAM HMM for lipocalins: resembles a position-specific scoring matrix

20 amino acids

position

```

HMMER2.0 (1.14)
NAME: lipocalin
ACC: PF00041
DESC: Lipocalin / cytosolic Entropulin binding protein family
LENG: 157
ALPH: amino
DF: 40
IC: 40
EAP: 788
COP: hmmdslid -e -F HMM_msa REE3::msa
COE: hmmdslidslidc ---seed 0 REE::msa
INDEL: 28
DATE: Thu May 29 15:09:42 2003
COPIES: 2018
SA: 9.0 9.0 9
TC: 9.0 9.0 9
IC: 9.0 9.0 9
BT:
---
-4 -1000 -1000 -8455 -4 -8455 -4
SILET: -4 -8455
MUSK: 895 -1558 85 338 -294 453 -1158 197 249 902 -1085 -142 -21 -333 45 531 205 394 -1998 -644
INDEL: 28
EFG: -9.782953 0.648810
MSE:
a-c b-d e-f g-h i-jk l-mn o-pq r-st u-v w-x y-z
-1000 -1000 233 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
-4 -10723 -11765 -894 -1115 -701 -1170 -1000 -805
-2 -1118 -2845 -4645 -3095 -6645 -3155 801 -8454 681 -2138 -633 -6714 -1152 -6214 -1588 -1107 1171 -3396 2245
-1 -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805
- 3 1990 -4020 -10318 1 -4548 -3888 -28 -4307 -2055 799 -3156 801 -3901 44 -1145 866 707 335 -4454 -485
- -110 -100 233 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -11765 -894 -1115 -701 -1170 -1000 -805
-702 -614 -278 394 -281 -747 -640 7120 -6440 -4374 -6440 -6440 -6210 -1424 -681 -1107 -7337 -6445
233 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
-11765 -894 -1115 -701 -1170 -1000 -805

```

PFAM HMM for lipocalins: GXW motif

PFAM HMM for lipocalins: GXW motif

20 amino acids

G W

```

22 528 -4297 968 1874 -4618 -3799 -2458 -1465 1714 -649 -3387 -1158 -1992 2 745 249 -1764 -1449 -4461 -3798
- -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805
23 -107 2119 -8999 -216 -1038 144 -607 391 550 929 1274 1490 -4823 -3745 -6023 -247 -105 37 -1824 -444
- -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805
24 277 -4634 802 -293 -6453 2399 674 -2039 -4313 -3345 -995 -1139 3507 -143 -808 -1123 -685 -6479 -3797
- -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805
25 -170 240 -2007 -1144 -4501 1893 -230 -486 339 451 -3389 -180 1882 699 331 148 -1444 1338 975 -3284
- -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805
26 661 -626 735 -792 349 -645 -1600 -241 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805
27 -680 -4338 449 1812 56 9779 98 -132 244 -9716 -1365 -147 -6895 941 1782 -43 -1100 -105 -288 -796
- -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805
28 1482 -3278 -204 -400 393 -4554 -3409 -1990 -1311 -1689 2388 -4651 -3749 -4802 -470 181 1354 -1427 1874
- -149 -500 333 43 -281 399 104 -614 230 -644 -700 278 394 45 96 359 117 -189 -294 -249
- -10723 -11765 -894 -1115 -701 -1170 -1000 -805

```

PFAM GCG MSF format

↓

```

1                               50
A1AG_HUMAN_38-183 QITQWF.YI ASAFPNEEYN .KSVQEQAT FFFTFPKRTE DTIFLR.EYQ
A1AG_RABIT_38-183 QLSKHF.FT ASAFPNPKYK .QLVQNTAA FFFTAIKKE DTLLR.EYQ
A1AH_MOUSE_39-184 WLSDKF.FI GAAVLNPDYK .QEQIQTVV FFLTLNLI DTIELR.EYH
A1AG_RAT_39-183 WLSDKF.YM GAADFPVK .QAVQTIQE VFYLPNLI DTIELR.EFQ
APHR_CRICR_21-165 ELQGGV.TI VVADNLEKI .KEGGELRFY FRHIDCYKNC SEMEIT.FYV
OBP_RAT_27-170 EVNGDWR.TL VVADNVEKV .AEGGSLRAY FQHGEGDEC QELKI.FNV
PBA5_RAT_27-170 KIEGWR.TV YLAASSVEKI .NEGSPRTY FRHICQK.R CNRINI.FYF
MSP1_MOUSE_32-175 KINGEWH.TI ILASSNREKI .EDGQFRIF LEQIVLE . NSLVLEK.FYH
MUPH_MOUSE_36-179 QISQYF.SI AEASVEEKKI .KENGSMRAF VENITVLE . NSLVFK.FHL
MUP_RAT_33-176 KNGDFV.SI VVASSREKI .KENGSMRAF MOHIVLE . NSLGFK.FRI
OBP_BOVIN_12-156 ELSGPRV.TV YIGSTPEKI .QENGFRTY FRELVDDEK GTVDFY.FSV
COBG_HUMAN_46-188 QFAGTUL.LV AVGSACRFIQ .EQGHRAET TLHVAPOG . TAMAVS.FTR
AMPF_HUMAN_39-188 RIYGVKY.NL AIGSTCPWLK .KIMDRHTVS TLVLEGATE AEISMT.SFR
AMPF_PLEPL_41-189 RFVQTVH.DV ALTSSCPHQ . .RNPADAII QKLVLEKDTG NKLKVT.RTR
L1PO_BUPFA_32-179 KILGWHY.GI GLASSNHWQ . SKKQQLKMC TTYITFTA.D QHLDPV.ATF
PGHD_HUMAN_38-186 KFLORV.SA GLASSNSULR .EKKAALMC KSVVAPAT.D GGLNLT.STF
NGAL_HUMAN_46-195 QFGQKY.VV GLAG.NAILR .EDKDPQKMY AT.IYELK.E DKSYNV.TSV
NGAL_MOUSE_46-197 QFQGVV.VV GLAG.NAVOK . .KTEGFTM YSTIVLEQ.E MNSYNV.TSI
ERBP_RAT_32-176 KFLGFY.EI AFASKNQTPG . .LANKEEM GANVVELK.E NLLALT.TTY
QSP_CHICK_29-173 EVAQKVI.NI ALASNTDFL .REKQKMKV HARISFLG.E DELEVS.YAA
ESP4_LACVV_33-167 KTVQVH.PI GHASLPEVP . .EYEQISP BDHHEVLT.G GQHLT.ANY
CEFA_RAMP1_30-174 KTVGVV.GI AASSNPFQL QMSDHPAP VNTYSLNM . GHRKSS.TSF
LALP_MACUL1_28-171 FSEQTY.VQ VLVW.DREPK .EDEFPDIS PI.VITLHM.H QKREAK.FTY
VGI1_RAT_25-172 DVSOTY.LK AAAM.DKEIP DKFFGSVST PHIKITL.E GNLQVK.FTV

```

Pfam (protein family) database

The screenshot shows the Pfam Logo generation form. The sequence logo displays the relative entropy of residues at each position. Key residues are highlighted: K at position 2, Y at position 3, A at position 4, S at position 11, and N at position 12. The x-axis is labeled 'Position' and the y-axis is 'Relative Entropy'.

PFAM JalView viewer

The screenshot shows the JalView alignment viewer. It displays a multiple sequence alignment of protein sequences. The alignment is color-coded by conservation, with highly conserved regions in red and less conserved regions in blue. The viewer includes a 'Calculate' menu and a 'Help' button.

Average distance tree using PID

The screenshot shows an average distance tree using PID. The tree is rooted and shows the relationships between different protein families. The x-axis represents the distance between sequences, and the y-axis represents the sequence identifiers. The tree is color-coded by conservation.

Sequence analysis

You may use either the [interpro/protein](#) sequence identifier (ID) or accession number (ACC) or the protein sequence itself to request the result reports

Sequence ID or ACC

Sequence

Sequence SMART

Domains detected by SMART

You can search for keywords in the domain annotation

Search keywords

Architecture analysis

You can search for proteins with combinations of specific domains in different species taxonomic range

Domain selection

Taxonomic selection

Alert

SMART: Simple Modular Architecture Research Tool (emphasis on cell signaling)

Page 199

CDD: Conserved domain database (at NCBI): CDD = Pfam + SMART

- Go to NCBI → Domains & Structure (left sidebar)
- Click CDD
- Enter a text query, or a protein sequence

The screenshot shows the NCBI Conserved Domains database search interface. The search bar contains the text 'Conserved Domains' and the search button is highlighted. The interface includes a search bar, a search button, and a list of search results.

CDD entry for "globin"

CDD = PFAM + SMART

CDD = PFAM + SMART

CDD entry for "globin"

CDD entry for "globin"

CDD uses RPS-BLAST: reverse position-specific

Purpose: to find conserved domains in the query sequence

Query = your favorite protein

Database = set of many position-specific scoring matrices (PSSMs), i.e. a set of MSAs

CDD is related to PSI-BLAST, but distinct

CDD searches against profiles generated from pre-selected alignments

Multiple sequence alignment: outline

- [1] Introduction to MSA
 - Exact methods
 - Progressive (ClustalW)
 - Iterative (MUSCLE)
 - Consistency (ProbCons)
 - Structure-based (Expresso)
 - Conclusions: benchmarking studies
- [2] Hidden Markov models (HMMs), Pfam and CDD
- [3] MEGA to make a multiple sequence alignment
- [4] Multiple alignment of genomic DNA

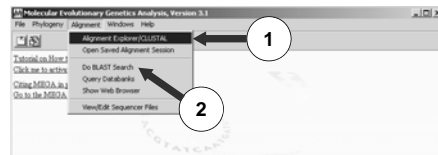
MEGA version 4: Molecular Evolutionary Genetics Analysis

Download from www.megasoftware.net

MEGA version 4: Molecular Evolutionary Genetics Analysis

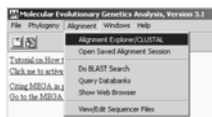


MEGA version 4: Molecular Evolutionary Genetics Analysis

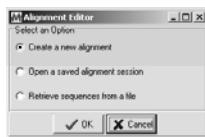


- Two ways to create a multiple sequence alignment**
1. Open the Alignment Explorer, paste in a FASTA MSA
2. Select a DNA query, do a BLAST search

Once your sequences are in MEGA, you can run ClustalW then make trees and do phylogenetic analyses



- [1] Open the Alignment Explorer



- [2] Select "Create a new alignment"



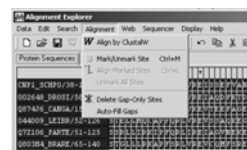
- [3] Click yes (for DNA) or no (for protein)



- [4] Find, select, and copy a multiple sequence alignment (e.g. from Pfam; choose FASTA with dashes for gaps)



- [5] Paste it into MEGA



- [6] If needed, run ClustalW to align the sequences

- [7] Save (Ctrl+S) as .mas then exit and save as .meg

Multiple sequence alignment: outline

- [1] Introduction to MSA
 - Exact methods
 - Progressive (ClustalW)
 - Iterative (MUSCLE)
 - Consistency (ProbCons)
 - Structure-based (Expresso)
 - Conclusions: benchmarking studies
- [3] Hidden Markov models (HMMs), Pfam and CDD
- [4] MEGA to make a multiple sequence alignment
- [5] Multiple alignment of genomic DNA

Multiple sequence alignment of genomic DNA

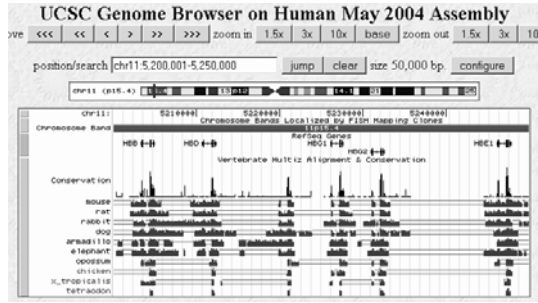
There are typically few sequences (up to several dozen), each having up to millions of base pairs. Adding more species improves accuracy.

Alignment of divergent sequences often reveals islands of conservation (providing "anchors" for alignment).

Chromosomes are subject to inversions, duplications, deletions, and translocations (often involving millions of base pairs). E.g. human chromosome 2 is derived from the fusion of two acrocentric chromosomes.

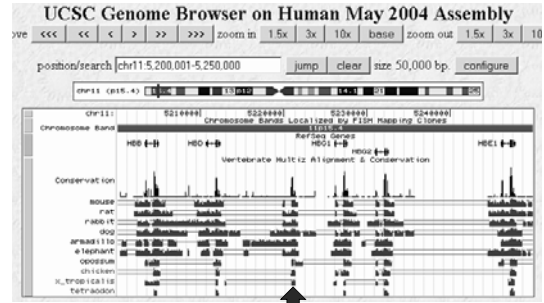
There are no benchmark datasets available.

Multiple alignment of genomic DNA at UCSC
50,000 base pairs (at <http://genome.ucsc.edu>)



Page 205

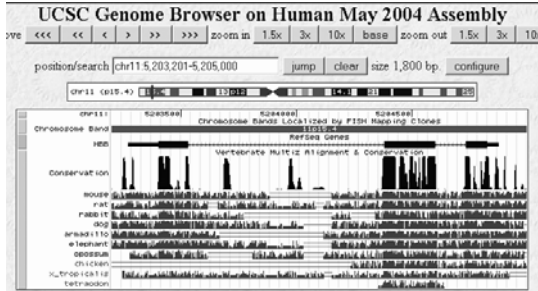
Note conserved regions: exons and regulatory sites
(scale: 50,000 base pairs)



regulatory

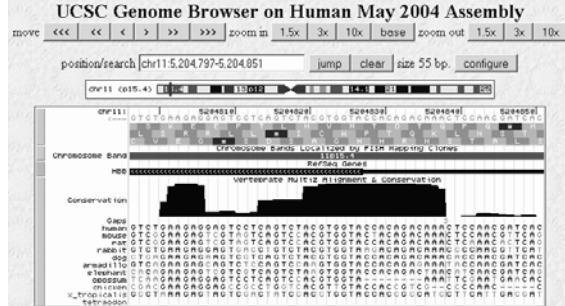
Page 205

Multiple alignment of beta globin gene
scale: 1,800 base pairs



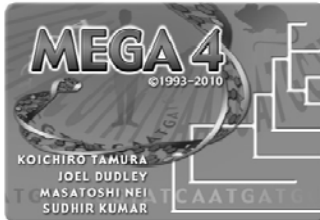
Page 205

Multiple alignment of beta globin gene
scale: 55 base pairs



Page 205

This week: please download MEGA software and paste in a set of protein sequences. We'll use MEGA next week to make phylogenetic trees.



MEGA5

Now includes Maximum Likelihood (ML) methods for tree searching and model testing for DNA and protein alignments, and many more improvements.

Click to DOWNLOAD

Download MEGA

Download from www.megasoftware.net