

Motif-based Hidden Markov Models for Multiple Sequence Alignment

William N. Grundy • Charles P. Elkan
Dept. of Computer Science & Engineering
University of California, San Diego

Abstract

- Protein families are well characterized by a collection of motifs (Sonnhammer & Kahn 1994), sometimes referred to as the “common core” (Chothia & Lesk 1986).
- These motifs can have structural and functional significance, and they may frequently be operated upon as units by diverse evolutionary mechanisms.
- The quality of a multiple alignment depends upon how accurately it identifies an ordered series of motifs (McClure et al. 1994, 1995).
- Hidden Markov models (HMMs) provide a theoretically sound modeling paradigm for collections of motifs for which efficient algorithms exist.
- Meta-MEME (Grundy et al. 1996, 1997) is a software toolkit that builds left-to-right, motif-based HMMs that focus upon the common core.
- Meta-MEME has been shown to detect remote homologies using smaller training sets than are required by standard HMMs.
- In an analysis of four protein families, Meta-MEME alignments are shown to be of higher quality than those produced by standard HMMs.
- In a previous analysis of nine other multiple alignment methods (McClure et al. 1994), only one method yields significantly higher quality alignments than Meta-MEME.

Meta-MEME

- **Step 1:** MEME builds motif models with maximal posterior probability, given a set of related sequences.
- **Step 2:** MAST uses only the motifs that appear in more than half the training sequences to search a database for the typical order and spacing of those motifs.
- **Step 3:** Meta-MEME constructs a hidden Markov model that combines the motif models in a linear fashion according to the typical schema.

Methods

- Four families were analyzed: globins, eukaryotic kinases, aspartic acid proteases, and the RH domain of the RNA-directed DNA polymerase (reverse transcriptase).
- Each family was represented by a test set of twelve sequences (McClure et al. 1994).
- Sequence identity among the sequences was 10-70% for the globins and 8-30% for all other families.
- Standard HMMs (Krogh et al. 1994) were created using the HMMER v1.8 tool `hmmt` (Eddy 1995).
- Meta-MEME and HMMER models were aligned using `hmma` with default parameter settings.

Multiple Alignments

- Pre-determined test motifs are indicated by color.
- The x's above the alignments indicate Meta-MEME motif positions or HMMER match state positions.
- Amino acids are unaligned outside of motifs in Meta-MEME alignments or within insertions in HMMER alignments.

Meta-MEME Alignment of Globin Sequences

```
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
HAHU  vlspa.....dktnVKAAGWKVgahag.....eygaEALERMFSLFPTTKTYFphfd...
HAOR  mltada.....ekkeVTALWGKAaaghge.....eygaEALERLRFQAFPTTKTYFshfd...
HADK  vlsaa.....dktnVKGVFSKIGghae.....eygaETLERMFIAYPQTKTYFphfd...
HBHU  vhltp.....eeksaVTALWGKVNvde.....vggeALGRLLVVPWTQRFFesfgdls
HBOR  vhlsg.....geksaVTNLWGKVNine.....lggeALGRLLVVPWTQRFFeafgdls
HBDK  vhwta.....eekqlITGLWGKVNvad.....cgaEALARLLIVYPWTQRFFasfgnls
MYHU  glsdg.....ewqlVLNVWGKVeadip.....ghgqEVLIRLRFKGGHPETLEKFDkfkhlk
MYOR  glsdg.....ewqlVLKVWGKVEgdip.....ghgqEVLIRLRFKTHPETLEKFDkfkglk
IGLOB mkffavlalcivgaiaspladeaslVQSSWKAVshn.....evEILAAVFAAYPDIQNKfsqf...
GPUGNI altek.....qealLKQSWEVLkqnip.....ahslRRLFALIEAAPESKYVFsflkds.
GPYL  gvl.....tdvQVALVKSsfefnanipknthRFFTLVLEIAPGAKDLFsflkgs.
GGZLB mldqqtin.....iikatvPVLKEHGVTit.....tTFYKNLFAKHPEVRPLFdmg...
```

```
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
HAHU  ...lshGSAQVKGHGKVKVadalt.....navahVDDMPNALsa.....1SDLHAHAKLRVDPVNFk....
HAOR  ...lshGSAQIKAHGKVKVadals.....taagHFDDMDSALsa.....1SDLHAHAKLRVDPVNFk....
HADK  ...lshGSAQIKAHGKVKVaaalv.....eavnHVDDIAGALsk.....1SDLHAQKLRVDPVNFk....
HBHU  tpdavmGNPKVKAHGKVKVlgafs.....dglahLDNLKGTfFat.....1SELHCDKLHVDPENFr....
HBOR  sagavmGNPKVKAHGAKVltsfg.....dalkNLDDLKGTfFak.....1SELHCDKLHVDPENFn....
HBDK  sptailGNPMVRAHGKVKVltsfg.....davkNLDNIKNTfFaq.....1SELHCDKLHVDPENFr....
MYHU  sedemkASEDLKKHGATVltalggi..lkkkghHEAEIKPLAq.....SHATKHKIPVKYLEF.....
MYOR  tedemkASADLKKHGTVltalgni..lkkkqgHEAELKPLAq.....SHATKHKISIKFLEY.....
IGLOB ...agkDLASIKDTGAFat.....HATRIVSFLse.....vIALSGNTSNAAAVNSlvsklg
GPUGNI .neipeNNPKLKAHAAVifkticesatelrqkgHAVWDNNTLkr.....LGSIHLLKNKITDPHFevmk...
GPYL  .sevpqNNPDLQAHAAGKvfklt.....yaaaIQLEVNGAVasdatlksLGSVHVSKGVVDAHFpVvk...
GGZLB ...rqeSLEQPKALAMTVlaa.....aqNIENLPAILpav...kKiAVKHCQAGVAAAHYpi.....
```

```
.....XXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.....
HAHU  .....lshCLLVTLAAHLPAEFTPAVHASLDKfl.....asVSTVLTskYR
HAOR  .....lLAHCILVVLARHCPGEFTPSAHAAMDkfl.....skVATVLTskYR
HADK  .....fLGHCFLVVVAIHHPAALTPEVHASLDKfm.....caVGAVLTAKYR
HBHU  .....lLGNVLVCVLAHFFGKEFTPPVQAAYQkv.....agVANALAHKYH
HBOR  .....rLGNVLIVVLAHFHFKDFSPVQAAWQklv.....sgVAHALGHKYH
HBDK  .....lLGDILIIIVLAHFHTKDFTPVQAAWQklv.....rvVAHALARKYH
MYHU  .....ISECIIQVLSKHPGDFGADAQGAMNKalelf..rkdmasNYKELGFQ
MYOR  .....ISEAIIHVLQSKHSADFGADAQAAMGKalelf..rndmaAKYKEFGFQ
IGLOB ddhkargvsaaqfGEFRTALVAYLQANVSWGDNVAAAWNKal.....dNTFAIVPRL
GPUGNI .....gaLLGTIKEAIKENWSDEMGQAWTEAYNql.....VATIKAEMKE
GPYL  .....eaILKTIKEVVGDKWSEELNTAWTIAYDEla.....iIIKKEMKDA
GGZLB .....VGQELLAGAIKEVLGDAATDDILDWGWKAYgviadvfiqVEADLYAQAVE
```

HMMER Alignment of Globin Sequences

```
XXXXXXXXXXXXX.XXXXXXXXXXX.XXXXXX.....XXXXXXXXX.XXXXXXXXXX.XXXXXXXXXXXXXX.XXXXXXXXXXX
HAHU V.LSPADKTN.VKAAWGKVG.AHAGE.....YGAEAL.ERMFLSF..PTTKTYFPH.FDLS.HGSA
HAOR M.LTDAEKKE.VTALWGKAA.GHGEE.....YGAEAL.ERLFQAF..PTTKTYFSH.FDLS.HGSA
HADK V.LSAADKTN.VKGVFSKIG.GHAEE.....YGAETL.ERMFIAY..PQTKTYFPH.FDLS.HGSA
HBHU VHLTPEEKSA.VTALWGKVN.VDEVG.....G.EAL.GRLLVVY..PWTQRFES.FGDL.STPD
HBOR VHLSGGEKSA.VTNLWGKVN.INELG.....G.EAL.GRLLVVY..PWTQRFEEA.FGDL.SSAG
HBDK VHWTAEKQL.ITGLWGKVNvAD.CG.....A.EAL.ARLLVIVY..PWTQRFAS.FGNL.SSPT
MYHU G.LSDGEWQL.VLNVWGKVE.ADIPG.....HGQEV.L.IRLFKGH..PETLEKFDK.FKHL.KSED
MYOR G.LSDGEWQL.VLKVWGKVE.GDLPG.....HGQEV.L.IRLFKTH..PETLEKFDK.FKGL.KTED
IGLOB M.KFFAVLALCiVGAIASPLT.ADEASlvqsswkavshNEVEIlaAVFAAY.PDIQNKFSQFaGKDLASIKD
GPUGNI A.LTEKQEAAL.LKQSWEVLK.QNIPA.....HS.LRL.FALIIEA.APESKYVFSF.LKDSNEIPE
GPYL GVLTDVQVAL.VKSSFEEFN.ANIPK.....N.THR.FFTLVLEIaPGAKDLFSF.LKGSSEVPQ
GGZLB M.L.DQQTIN.IIKATVPVLkEHGVT.....ITTF.YKNLFAK.HPEVRPLFDM.GRQ..ESLE
```

```
XXXXX.XXXXXXXXXXXXXXXXXX.XXXXXXXXXXXXXXXXXXXXXX.XXXXXX.XXXXXX.....XXXXXXXXXXXXXXXXXXXXX
HAHU QVKGH.GKKVADA.LTN.....AVA.HVDDMPNA..LSALS.D.LHAHKL...RVDPVNF.KLLSHCLL
HAOR QIKAH.GKKVADA.L.S.....TAAGHFDDMSA..LSALS.D.LHAHKL...RVDPVNF.KLLAHCIL
HADK QIKAH.GKKVAAA.LVE.....AVN.HVDDIAGA..LSKLS.D.LHAQKL...RVDPVNF.KFLGHCF
HBHU AVMGNpKVKAHGK.KVLGA..FSDGLAHLNLDLKG...FATLS.E.LHCDKL...HVDPENF.RL.LGNVL
HBOR AVMGNpKVKAHGA.KVLTS..FGDALKNLDDLKG...FAKLS.E.LHCDKL...HVDPENFNRL..GNVL
HBDK AILGNpMVRAHGK.KVLTS..FGDAVKNLNLDNIKNT...FAQLS.E.LHCDKL...HVDPENF.RL.LGDIL
MYHU EMKASeDLKKHGA.TVL.....TALGGILKKGHH..EAEIKPL.AQSHATK...HKIPVKYLEFISECII
MYOR EMKASaDLKKHGG.TVL.....TALGNILKKGQH..EAELKPL.AQSHATK...HKISIKFLEYISEAII
IGLOB T.GA...FATHATRIVSFLseVIALSGNTSNAAAV...NSLVSKL.GDDHKA...R.GVSAA.QF..GEFR
GPUGNI NNPK...LKAAHAaVIFKTI...CESATELRQKGHAVwdNNTLKR.L.GSIHLK...N.KITDP.HF.EVMKG
GPYL NNPd...LQAHAG.KVFKL..TYEAAIQLEVNGAVAs.DATLKS.L.GSVHVS...K.GVVDa.HF.PVVK
GGZLB Q.....PKALAM.TVL.....AAAQNIENLPAIL..PAVKKIaVKhCQAGVaaH.YPIVGQEL.LGAIK
```

```
XXXXXXXXXXXXX.XXXXXXXXXXX.XXXXXXXXXXXXXXXXXXXXXXXXXXXXXX.X
HAHU VT.LAA.H..LPAEFTPA..VHASLdkFLASV.STVLTS..KY..R
HAOR VV.LAR.H..CPGEFTPS..AHAAMDKFLSKV.ATVLTS..KY..R
HADK VV.VAI.H..HPAALTPE..VHASLdkFMCAV.GAVLTA..KY..R
HBHU VCVLAH.H..FGKEFTPP..VQAAYQKVAGV.ANALAH..KY..H
HBOR IVVLAH.H..FSKDFSPE..VQAawQKLvSGV.AHALGH..KY..H
HBDK IIVLAA.H..FTKDFTPe..CQAawQKLVRVv.AHALAR..KY..H
MYHU QV.LQSKHPgDFGADAQGA.MNKALELFRKDM.ASNYKELGFQ..G
MYOR HV.LQSKHSaDFGADAQAA.MGKALELFRNDM.AAKYKEFGFQ..G
IGLOB TA.LVA.Y..LQANVSWGDnVAaAWNKA.LDN.TFAIVV..PR..L
GPUGNI ALLGTIKeA.IKENWSDE..MGQawTEAYNQLVATIKAE..MK..E
GPYL AILKTIKeV.VGDKWSEE..LNTawTIAYDELAIiIKKE..MKdaA
GGZLB EVLGDAAT..DDILDAWGK.AYGVIAADVFIQVEADLYAQ..AV..E
```

Meta-MEME Alignment of RH Sequences

```
.....XXXXXXXXX.....XXXXXXXXX
HTLV-II ldt.....apCLFSDGSPqkaayvlwdqtilqqd.....itplpshethsaqkgELLALICG
SRV-I lnn.....alLVFTDGSSStgmaaytladtti.....kfqtnlnsaqlvELQALIAV
RSV pvp.....gpTVFTDASSsthkgvvvwwregprw.....eikeiadlgasvqqlEARAVAMA
HIV-II ipg.....aeTFYTDGSCnrqskegkagyvtdrg.....kdkvkkleqttnqqaELEAFAMA
MoMLV pda.....dhtWYTDGSSllqegqrkagaavttet.....ewiwakaldagtsagraELIALTQA
Ingi pre.....hyKLWTDGSSvslgeklgaaallhrnntl.....icapktgagelscsyraECVALEIG
CAMV pee.....kLIIEYTDASDDdywggmlkaikinegtntelicryasgsfkaaeknyhsndkETLAVINT
17.6 ftk.....kftLTTDASDdvalgavlsqdgghplsyi.....srtlnheheinystiekELLAIWVA
MAUP fnnstnlqepsdsrLLYRKGSwvnirfaay.....lysklseEKHGLVPK
HBV rpg.....lcQVFADATPtgwglvmghqrmr.....gtfsaplpihtaELLAACFA
Copia fen.....kiIGYVDSDWagseidrktstgylfkmfdf.nlicwntkrqnsvaassteaEYMALFEA
E.coli mlk.....qvEIFTDGSClgnppggygailryrg.....rektfsagytrttnnrmELMAAIVA

x.....XXXXXXXXX.....
HTLV-II Lraak.....pwpsLNIFLDSKYLikylhslaigafllgtsahqtlqaalp.....
SRV-I Lsafp.....nqpLNIYTDsAYLahsiplletvaqikhisetaklflqcqq.....
RSV Lllwp.....tptTNVVTDSAFVakmllkmgqegvpstaaafiledal.....
HIV-II Ltids.....gpkVNIIVDSQYVmgisasqpteseskivnqie.....
MoMLV Lkmae.....gkklNVYTDsRYAfatahihgeiyrrrglltsegkeiknkdeil.....
Ingi Lqrlkwlpl...ryrstpsrLSIFSDSLMltalqtgplavtdpilrrlwrll.....
CAMV Ikkfsiy.....ltpvhFLIRTDNTHfksfvnlnykgdsklgrnir.....
17.6 Tktfrhy.....llgrhFEISSDHQPLswlyrmkdpnkskltrwr.....
MAUP Flek.....lreINFALDKVDVteidsklsrlmkfsvsaaydevgtlalkslfkfrnseres
HBV Rsrss.....gaNIIGTDNSVvlrkytsfpwllgcaanwilrgtsfvvypsa.....
Copia VrealwklfltsiniklenpIKIYEDNQGCisiannp schkrakhidiky.....
E.coli Lealk.....ehceVILSTDSQYVrqgitqwihnwkkrgwktadkpkvknvd.....

.....XXXXXXXXX.....
HTLV-II .....pllqgktylhhvrshtnlpdpistfNEYTDSLILapl.....
SRV-I .....liynrsipfyighvrahsglpgpiahgNQKADLTKtvasn.....
RSV .....sqrsamaavlhrshsevpdfftegnDVADSQATfqay.....
HIV-II .....emikkeaiyvavwpahkgiggNQEVVDHLVSqgirqvl.....
MoMLV .....allkalflpkrlsiihcpghqkghsaeargNRMADQAARkaaitetpdtstll...
Ingi .....lqvqrrkirirlqfvfdhcgvkrNEVCDEMAKkaadlpql.....
CAMV .....wqawlshysfdvehikgtDNHFADFLSRefnkvn.....
17.6 .....vklsefddikyikgkeNCVADALSRIkleety.....
MAUP ikasfkqlrengkiaefsearrlwfeilkirldlfnasSLACD DLLShlqdrresi.....
HBV .....lnpaddpsrgrlglrpllrpfrpttgrtSLYADSPSVpshlpdrvh.....
Copia .....hfareqvqnnvicleyipteNQLADIFTKplpaarfve.....
E.coli .....lwqrldaalgqhikwewkghaghpNERCDELARaaamnptledtgyqvev
```


Alignment Scores

| Globins | 1 | 2 | 3 | 4 | 5 | Total |
|----------------|----------|----------|----------|----------|----------|--------------|
| META-MEME | 12 | 11 | 11* | 11*' | 12* | 57 |
| HMMER | 11 | 12*' | 11* | 10*' | 11*' | 55 |
| AMULT | 12 | 12 | 12 | 12 | 12 | 60 |
| ASSEMBLE | 12 | 11 | 12 | 12 | 12 | 59 |
| CLUSTAL V | 12 | 11 | 12 | 12 | 12 | 59 |
| DFALIGN | 12 | 12 | 12 | 12 | 12 | 60 |
| GENALIGN | 11*' | 12 | 12 | 10' | 11 | 56 |
| MULTAL | 12 | 11 | 12 | 12 | 12 | 59 |
| MACAW | 9 | 11 | 9 | 8 | 8 | 45 |
| PIMA | 12 | 12 | 12 | 12 | 12 | 60 |
| PRALIGN | 8 | 8* | 9* | 8* | 10 | 43 |

| Kinases | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | Total |
|----------------|----------|----------|----------|----------|----------|----------|----------|----------|--------------|
| META-MEME | 12 | 11 | 10 | 12 | 12 | 12 | 12 | 11 | 92 |
| HMMER | 12 | 12* | 12* | 12 | 12 | 12 | 12 | 12 | 96 |
| AMULT | 12 | 10 | 11 | 12 | 12 | 12 | 12 | 12 | 93 |
| ASSEMBLE | 11 | 7 | 10 | 12 | 12 | 12 | 12 | 12* | 87 |
| CLUSTALV | 12 | 11 | 11* | 12 | 12 | 12 | 12 | 12* | 94 |
| DFALIGN | 12 | 12 | 12 | 12 | 12 | 12 | 12 | 12 | 96 |
| GENALIGN | 12' | 9* | 10 | 12 | 12 | 12 | 12* | 11* | 90 |
| MULTAL | 12 | 9* | 10* | 12 | 12 | 12* | 12 | 12 | 91 |
| MACAW | 8 | 0 | 9 | 12 | 12 | 10 | 12 | 0 | 63 |
| PIMA | 12 | 11 | 11 | 12 | 12 | 12 | 12 | 12 | 94 |
| PRALIGN | 12 | 10* | 6* | 4 | 9* | 9* | 4 | 4 | 58 |

| Proteases | 1 | 2 | 3 | Total |
|------------------|----------|----------|----------|--------------|
| META-MEME | 12 | 4*' | 9* | 25 |
| HMMER | 11 | 8 | 6* | 25 |
| AMULT | 11 | 7 | 10 | 28 |
| ASSEMBLE | — | — | — | 0 |
| CLUSTALV | 12 | 9* | 6* | 27 |
| DFALIGN | 12 | 12* | 12 | 36 |
| GENALIGN | 11 | 8*' | 7* | 26 |
| MULTAL | 10 | 7* | 9* | 26 |
| MACAW | 12 | 4 | 8 | 24 |
| PIMA | 12 | 5* | 5* | 22 |
| PRALIGN | 8* | 4* | 8* | 20 |

| RHs | 1 | 2 | 3 | 4 | Total |
|------------|----------|----------|----------|----------|--------------|
| META-MEME | 11 | 9 | 11 | 12 | 43 |
| HMMER | 11 | 10* | 5*' | 7* | 33 |
| AMULT | 11 | 9* | 8* | 7* | 35 |
| ASSEMBLE | — | — | — | — | 0 |
| CLUSTALV | 12 | 9 | 9* | 9* | 39 |
| DFALIGN | 12 | 12 | 10 | 12 | 46 |
| GENALIGN | 12*' | 7 | 8*' | 9*' | 36 |
| MULTAL | 11* | 11* | 9* | 10 | 41 |
| MACAW | 7 | 5 | 7 | 3 | 22 |
| PIMA | 10 | 9 | 8* | 11* | 38 |
| PRALIGN | 9 | 8* | 6* | 3 | 26 |

Scores represent the number of sequences in which each motif was correctly aligned. Scores for non-HMM methods are from (McClure et al. 1994). A * indicates that the motif was correctly aligned in two or more misaligned subsets of the test sequences. A ' indicates that a gap was inserted into the motif.

Results

| | Globins | Kinases | Proteases | RHs | Total |
|------------------|-----------|-----------|-----------|-----------|------------|
| DFALIGN | 60 | 96 | 36 | 46 | 238 |
| CLUSTALV | 59 | 94 | 27 | 39 | 219 |
| META-MEME | 57 | 92 | 25 | 43 | 217 |
| MULTAL | 59 | 91 | 26 | 41 | 217 |
| AMULT | 60 | 93 | 28 | 35 | 216 |
| PIMA | 60 | 94 | 22 | 38 | 214 |
| HMMER | 55 | 96 | 25 | 33 | 209 |
| GENALIGN | 56 | 90 | 26 | 36 | 208 |
| MACAW | 45 | 63 | 24 | 22 | 154 |
| PRALIGN | 43 | 58 | 20 | 26 | 147 |
| ASSEMBLE | 59 | 87 | 0 | 0 | 146 |

- Meta-MEME ranks third overall.
- MEME discovers 19 of the 20 motifs in the test set.
- Only DFALIGN (Feng & Doolittle, 1987) significantly outperforms other methods.

Discussion

- Meta-MEME's focus on the common core allows accurate models to be trained from fewer sequences than are required by standard HMMs.
- By focusing its models on highly conserved regions of the training set, Meta-MEME effectively ignores noisy portions of the data, thereby allowing the software to properly align distant homologs.
- Meta-MEME alignments could be extended by running another alignment program on the unaligned, non-motif regions.
- A server for MEME and MAST, as well as the source code for Meta-MEME, is available at <http://www.sdsc.edu/MEME>.

References

- Bailey, T. L. and C. P. Elkan. "Fitting a mixture model by expectation-maximization to discover motifs in biopolymers." *2nd ISMB*, 1994.
- Chothia, C. and A. M. Lesk. "The relation between the divergence of sequence and structure in proteins." *EMBO*, 5:823-826, 1986.
- Eddy, S. R. "Multiple Alignment Using Hidden Markov Models." *3rd ISMB*, 114-120, 1995.
- Feng, D.-F. and R. F. Doolittle. "Progressive sequence alignment as a prerequisite to correct phylogenetic trees." *J. Mol. Evol.*, 25:351-360, 1987.
- Grundy, W. N., T. L. Bailey, C. P. Elkan and M. E. Baker. "Meta-MEME: Motif-based hidden Markov models of protein families." *CABIOS*, 1997. To appear.
- Grundy, W. N., T. L. Bailey, C. P. Elkan and M. E. Baker. "Hidden Markov model analysis of motifs in steroid dehydrogenases and their homologs." *BBRC*, 231(3):760-766, 1997.
- Krogh, A., M. Brown, I. Mian, K. Sjolander and D. Haussler. "Hidden Markov Models in Computational Biology: Applications to Protein Modeling." *JMB*, 234:1501-1531, 1994.
- McClure, M. A., T. K. Vasi and W. M. Fitch. "Comparative analysis of multiple protein-sequence alignment methods." *Mol. Bio. Evol.*, 11(4):571-592, 1994.
- McClure, M. A. and R. Raman. "Parametrization studies of hidden Markov models representing highly divergent protein sequences." *28th HICSS*, 184-193, 1995.
- Sonnhammer, E. and D. Kahn. "The modular arrangement of proteins as inferred from analysis of homology." *Protein Science*, 3:482-492, 1994.