

Multiple alignment

Adapted from the courses of the Bonsai team,

CRIStAL UMR 9189

Sylvain.legrand@univ-lille.fr

Introduction

Definitions

- Input: K sequences

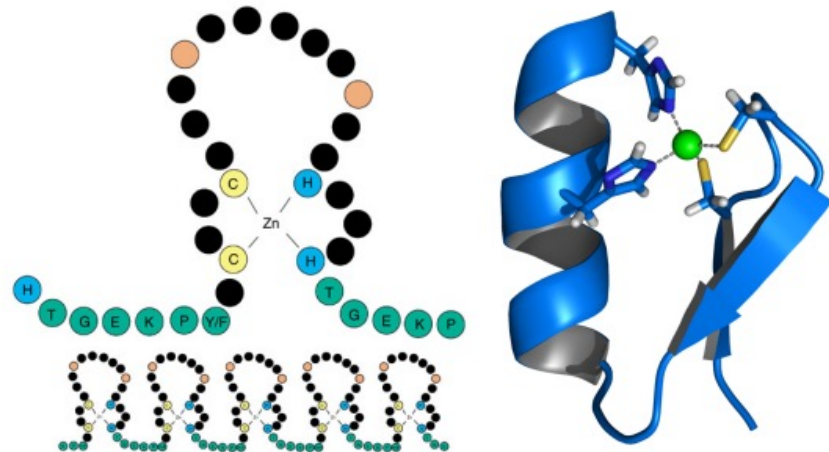
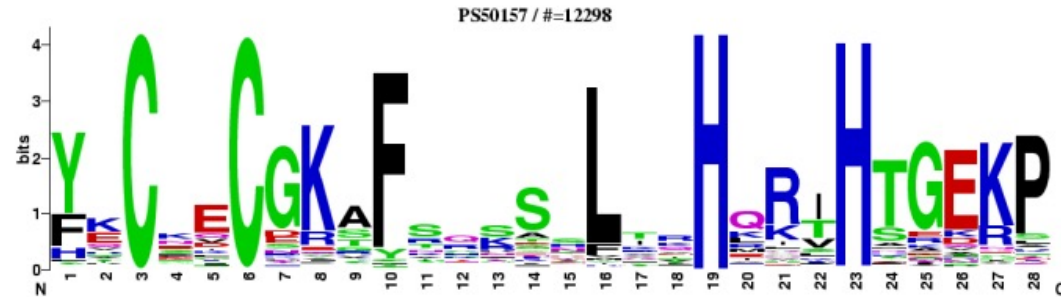
```
C A T G C G A G T A G T A G
C A T G G T A G T A G
C C T G G A G T A C G T A G
C A T G A G C G T A G
```

- Output: a table containing the k sequences, with indels

```
C A T G C G A G T A - G T A G
C A T G - - - G T A - G T A G
C C T G - G A G T A C G T A G
C A T G - - A G - - C G T A G
```

Zinc finger pattern (C2H2-type)

TTY1_HUMAN	YVCPFDG	CNKKFAQSTNLKSH	ILT--H
YKQ8_CAEEL	YKCT--V	CRKDISSESRLTH	MFQHH
BASO_HUMAN	FQCD--I	CKKTFKNACSVKI	HHKN-MH
ZG2-9_XENL	FVCT--V	CGKTYKYKHGLNTH	LHS--H
P43_XENBO	LKCSVPG	CKRSFRKKRALRI	HVSE--H
IKAR_MOUSE	FECN--M	CGYHSQDRYEFSS	HITRGEH
TRA1_CAEEL	YKCEFAD	CEKAFSNASDRAK	HQNR-TH
ZN10_HUMAN	YKCN--Q	CGIIFSQNSPFIV	HQIA--H
XFIN_XENLA	FRCS--E	CSRSFTHNSDLTA	HMRK--H
TF3A_BUFAM	CKCETEN	CNLAFTTASNML	HFKR-AH
ZG58_XENLA	FVCT--E	CNLSFAGLANLRSH	QHL--H
P43_XENBO	YRC SYED	CQTVSPTWTALQT	HLKK--H
TSH_DROME	FRCV--W	CKQSFPPTLEALT	THMKDSKH
ZN76_HUMAN	FRCGYKG	CGRLYTTAHHLKV	HERA--H
TF3A_BUFAM	YRCPREN	CDRTYTTKFNLKSH	ILT-FH
SUHW_DROAN	YACK--I	CGKDFTRS YHLKR	HQKYSSC
ZN76_HUMAN	YTCPEPH	CGRGFTSATNYKN	HVRI--H
SRYC_DROME	FKCN--Y	CPRDFTNFPNWLK	HTRR-RH
EVI1_HUMAN	YRCK--Y	CDRSFSISSNLQR	HVRN-IH



modélisation : motif **Prosite**

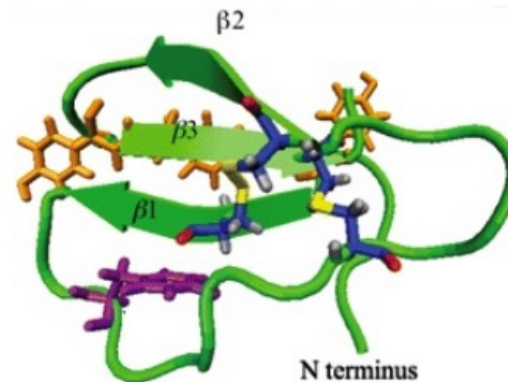
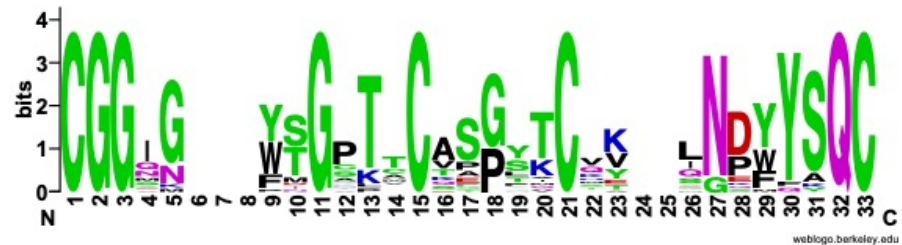
C-x(2,4)-C-x(3)-[LIVMFYWC]-x(8)-H-x(3,5)-H

Cellulose binding site

```

HWGQCGGI---GYSGCKTCTSGTTQYNSNDYYSQCL
HYGQCGGI---GYSGPTVCASGTTQVLNPYYSQCL
QWGQCGGI---GYTGSTTCASPYTCHVLNPYYSCY
VWGQCGGQ---NWSGPTCCASGSTCVYSNDYYSQCL
LYGQCGGA---GWTGPTTCQAPGTCKVQNWYSQCL
IWGQCGGN---GWTGATTCASGLKCEKINDWYYQCV
VWGQCGGN---GWTGPTTCASGSTCVKQNDFYSQCL
DWAQCGGN---GWTGPTTCVSPYTCTKQNDWYSQCL
QWGQCGGQ---NYSGPTTCKSPFTCKKINDFYSCQ
RWQQCGGI---GFTGPTTCEEPYICTKLNDWYSQCL
HWAQCGGI---GFSGPTTCEPEYTCAKDHDYISQCV
LYEQCGGI---GFDGVTCSEGLMCMKMGPPYYSQCR
VWAQCGGQ---NWSGTPCCTSGNKCVKLNDFYSCQ
PYGQCGGM---NYSGKTMCSPGFKCVELNEFFSQCD
AyyQCGGSKSAYPNGNLACATGSKCVKQNEYYSQCV
EYAACGGE---MFMGAKCKFGLVCYETSGKWSQCR
    
```

extrait de Prosite, entrée PS00562



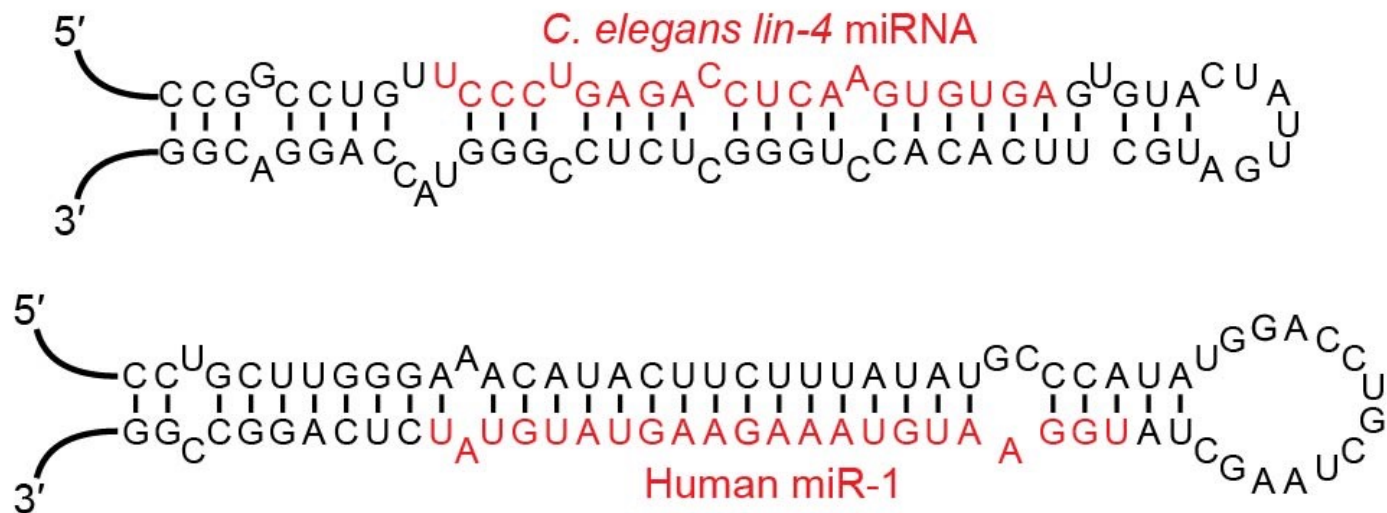
C-G-G-x(4,7)-G-x(3)-**C**-x(5)-**C**-x(3,5)-[NHG]-x-[FYWM]-x(2)-Q-**C**



the 4 cysteines are involved in disulfide bonds (SS-bond)

RNA structure

- We have a family of RNAs possessing the same secondary structure
- For a given structure pairing:
 - if a base mutates in the RNA structure, the base that matches it must mutate too... → **compensatory mutation**



RNA structure

- Example

G	A	G	C	C	C	A	G	U	U	C
	A	G	G	A	C	U	C	U	U	C
A	A	U	C	A	C	C	C	G	A	U

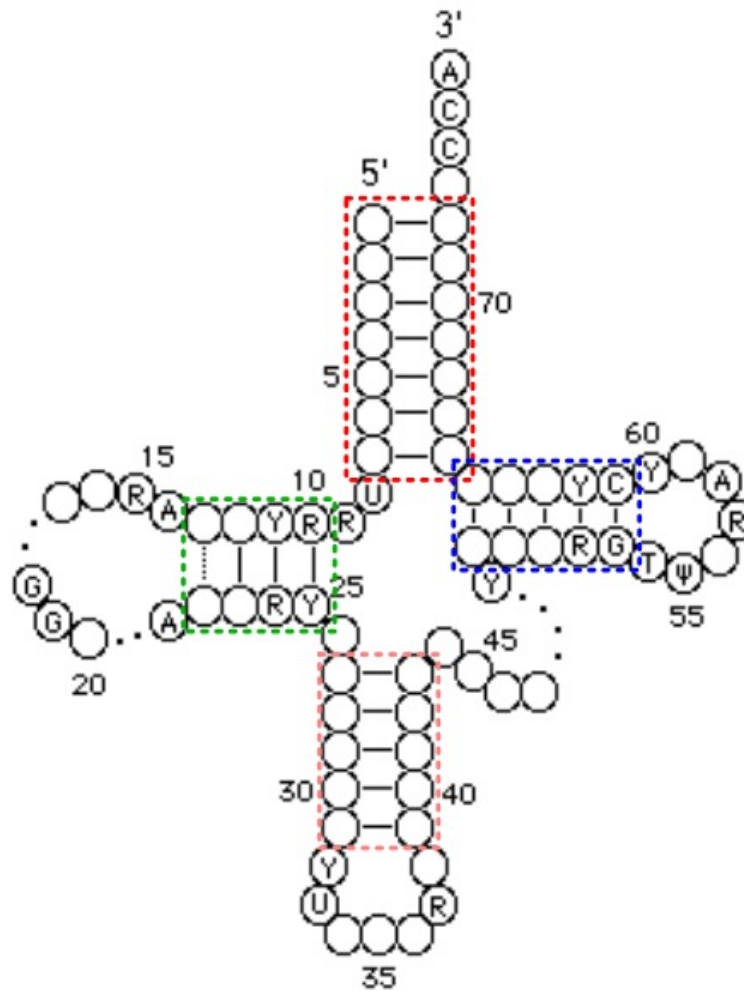
- Step 1: **multiple alignment** construction

G	A	G	C	–	C	C	A	G	U	U	C
–	A	G	G	A	C	–	U	C	U	U	C
A	A	U	C	A	C	C	C	G	A	U	–

- Step 2: **correlated positions** detection

G	A	G	C	–	C	C	A	G	U	U	C
–	A	G	G	A	C	–	U	C	U	U	C
A	A	U	C	A	C	C	C	G	A	U	–

tRNA structure



tRNA structure

GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAAGCAGGGGGTTCAGATCCCGCTATTCTCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTTGCAACGCAGGAGGTCTGCGGTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTTGCAATGGCATGCAAGAGGTTCAGCGGTTCGATCCCGCTTAGCTCCACCA
GGGGAATTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAAGCAGGGGGTTCAGATCCCGCTATTCTCCA---
GGGGCTTATAGCTCAGTC-GGTAGAGCACCTGCCCTTTTGCAAGGCAGATGTCAGGGGTTCGATTCCCCTAGGCTCCA---
GGGGGTATAGCTCAGTT-GGTAGAGCGCTGCCCTTTTGCAAGGCAGGAAGTTCAGCGGTTCGATTCCGCTTACCCCA---
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTTGCAACGCAGGAGGTCTGCGGTTCGATCCCGCATAGCTCCACCA
GGGGCTATAGCTCAGCT-GGGAGAGCGCCTGCTTTTGCAACGCAGGAGGTCTGCGGTTCGATCCCGCATAGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAAGCAGGGGT-CGTCGGTTCGATCCCGTCTGCCCTCCACCA
GGGGCCATAGCTCAGCT-GGGAGAGCGCCTGCTTTTGCAACGCAGGAGGTTCAGAGTTCGATCCTCCTTGGCTCCACCA
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAAGCAGGGGT-CGTCGGTTCGATCCCGTCTGCCCTCCACCA
GGGGCCATAGCTCAGCTGGGGAGAGCGCCTGCCCTTGCAACGCAGGAGGTCAACGGTTCGATCCCGTTTGGCTCCA---
GGGGGCATAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAAGCAGGGGT-CGTCGGTTCGATCCCGTCTGCCCTCCACCA
GGGGCATTAGCTCAGCT-GGGAGAGCGCCTGCTTTTGCAACGCAGGAGGTTCAGCGGTTCGATCCCGCTATTCTCCACCA
GGGGCCATAGCTCAGTT-GGTAGAGCGCCTGCTTTTGCAAGCAGGTGT-CGTCGGTTCGAATCCGTCTGGCTCCACCA
GGGGCCGTAGCTCAGCTGGG-AGAGCACCTGCTTTTGCAAGCAGGGGGTTCGGAGTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GGGAGAGCACCTGCTTTTGCAAGCAGGGGGTTCGTCGGTTCGATCCCGTCCGGCTCCACCA
GGGGCCGTAGCTCAGCT-GG-AGAGCACCTGCTTTTGCAAGCAGGGGGTTCGTCGGTTCGATCCCGTCCGGCTCCACCA

Implementation

- **Pairwise alignment**

2 sequences → detect a syntactic similarity → Is there a common function?

- **Multiple alignment**

Family of sequences with the same function → To which syntactic conservation does this correspond?

- **Sum of pairs**

$$SP(m_i) = \sum_{1 \leq j < k \leq n} s(m_i^j, m_i^k)$$

m_i = la i -ème colonne de l'alignement

m_i^j = j -ème aa dans la colonne i

Score of a multiple alignment

Scoring system:

$$s(x,x)=1, \quad s(x,y)=-1, \quad s(x,-)=s(-,x)=-2, \quad s(-,-)=0$$

A	A	C	G	T	A	C	G	A	T	A
A	-	C	G	T	A	-	A	A	T	G
G	T	C	G	T	A	-	-	T	T	A

(1-2)	1	-2	1	1	1	1	-2	-1	1	1	-1
(1-3)	-1	-1	1	1	1	1	-2	-1	-1	1	1
(2-3)	-1	-2	1	1	1	1	0	-2	-1	1	-1
	=	=	=	=	=	=	=	=	=	=	=
	-1	-5	3	3	3	3	-4	-4	-1	3	-1

= -1

Score of a multiple alignment

Scoring system:

$$s(x,x)=1, \quad s(x,y)=-1, \quad s(x,-)=s(-,x)=-2, \quad s(-,-)=0$$

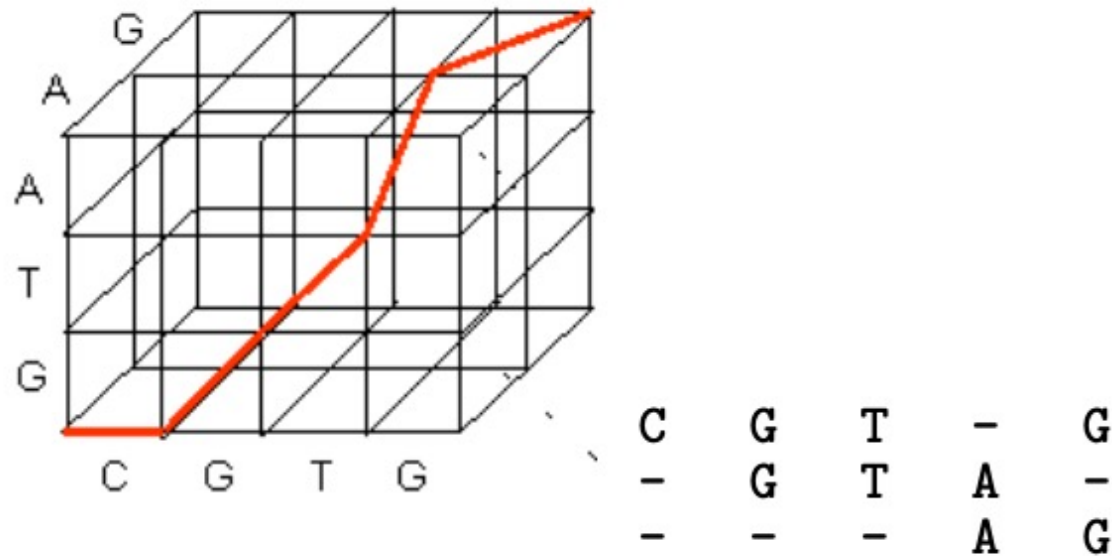
Alternatively... (identical)

A	A	C	G	T	A	C	G	A	T	A
A	-	C	G	T	A	-	A	A	T	G
G	T	C	G	T	A	-	-	T	T	A

(1-2)	1	-2	1	1	1	1	-2	-1	1	1	-1	=	1
(1-3)	-1	-1	1	1	1	1	-2	-1	-1	1	1	=	0
(2-3)	-1	-2	1	1	1	1	0	-2	-1	1	-1	=	-2
												=	-1

Exact algorithm: dynamic programming

- **Pairwise alignment** \rightarrow path in a **2-dimensions matrix**
- **Multiple alignment** of n sequences \rightarrow path in a **n -dimensions matrix**



But, impossible to use in practice... Will be **too time consuming**...

- **Definition:** algorithm using simple rules to reduce the search space for solutions (but not necessarily giving the best solution)
- Examples: Clustal, Dialign, Muscle, Multalin, T-coffee...
→ as many programs that can produce different alignments !
- The algorithms will be seen next semester...

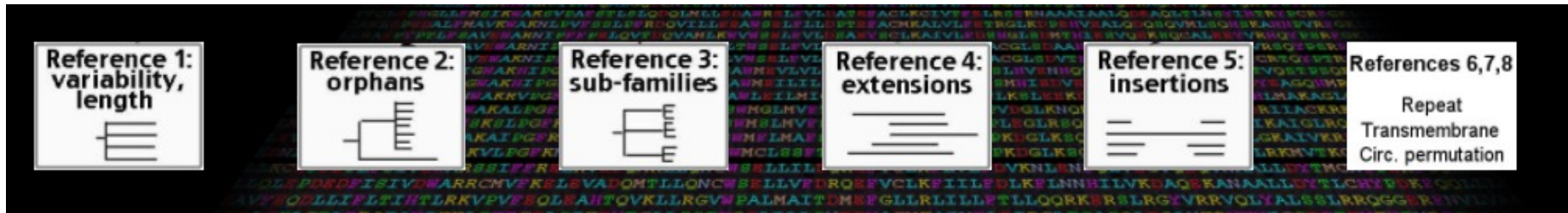
Méthode	Idée	Stratégie
MSA DCA	Extension de l'algorithme de Needleman et Wunsch	Simultanée
Clustal PIMA PILEUP MULTALIGN Dialign	Ajout successif de séquences ou groupes de séquences	Progressive
Saga/Coffee PRRN HMMT MUSCLE MA-FFT	Réalignement lors de l'ajout successif de séquences ou groupes de séquences	Itérative

What method should I use?

- It **depends on the type of sequences** to be aligned...
- **The more divergent the sequences, the less reliable the result**
- When the identity percent is **superior to 35%** → all methods are **satisfying**
- **Clustal** tends to allow **fewer gaps than Dialign2**
 - local similarity → Dialign
 - Global similarity → Clustal

What method should I use?

- **BaliBase**: multiple alignment database for benchmark
 - More than 150 protein families
 - Alignments based on secondary structure of proteins



- For References 1, 2 and 3: Clustal>Dialign
- For references 4 and 5: Dialign > Clustal

Examples of alignments

Examples

- Example from Cédric Notredame

```
GARFIELD THE LAST FAT CAT
GARFIELD THE FAT CAT
GARFIELD THE VERY FAST CAT
THE FAT CAT
```

Alignement fourni par Clustal

```
seq1      GARFIELDTHELASTFA-TCAT
seq2      ----GARFIELDTHEFA-TCAT
seq3      GARFIELDTHEVERYFASTCAT
seq4      -----THEFA-TCAT
```

Alignement fourni par Dialign2

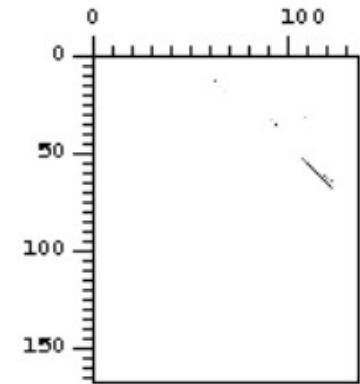
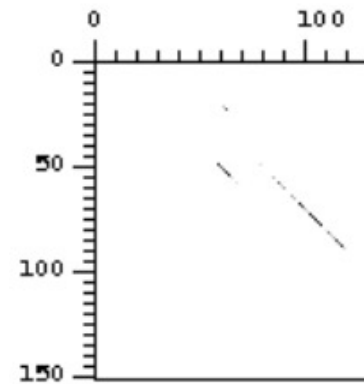
```
seq1      GARFIELD THE LAST FA-T CAT
seq2      GARFIELD THE ---- FA-T CAT
seq3      GARFIELD THE VERY FAST CAT
seq4      ----- THE ---- FA-T CAT
```

Examples

- **Helix – loop – helix domain**

- 5 sequences
- different lengths
- Local similarity

1)	HEN1-Human	133
2)	CBF1-Yeast	351
3)	HES5-Mouse	167
4)	IN04-Yeast	151
5)	ESC1-Yeast	413



Examples

- **Clustal** result (in pink, the domain that should be aligned)

```
-----MMLNSDTMELD-----LPPTHSETESG-----FSDCGGG
--MNSLANNNKLSTEDDEEHSARKRGYNEEQNYSEARKKQRDQGLLSQESNDGNIDSALLSEGATLKGTQSQYESG-----LTSNKDE
MSSYALPSMQPTPTSSIPLRQMSQPTTSAPSNSASSTPYSPQQVPLTHNSYPLSTPSSFQHGQTRLPPINCLAEPFNRPPWHSNSAAP
-----MAPSTVAVEMLSPEKNRLRKPVVEKMRRDR-----INSSIEQ
-----MTNDIKEIQTIQPGLSEIKEIKGELANVKKR-----

AGPD-----GAGPGG-----
KGSDDDEDASVAEAAVAATVNYTDLIQQQE-----DSSDAHTSNQTNANGEHKDSLNGERAITPSNEGVKPNTSLEGMTSSPMEST
ASSSPTSATLSTAHPVHTNAAQVAGSSSSSYVYSVPPTNSTTSQASAKHSAVPHRSSQFQSTTLTPSTTDSSSTDVSSSDSVSTSASS
LKLL-----

-----PGGGQARGPEPGEPRKD-----LQHLSREERRRRRRAT-----AKYRTA-----
QQSKNDMLIPLAEHARGPEHQDDEDNDDADID-----LKKDISMQPGRGRKPTTLATDEWKKQRKDS-----
NASNTVSVTSPASSSATPLPNQPSQQQFLVSKNDAFTTFVHSVHNTPMQQSMYVPQQQTSHSSGASYQNESANPPVQSPMQYSYSQGGP
-----LEQEFARHQPNKLEKAD-----ILEMAVSYLKHSKAFAAAAGPKSLHQDYSEG-----
-----KRRSKKINKLTDGQIR-----INHVSSEKKRRELERAIFFDELVAVVPDLQPQ-----

-----HATRERIRVEAFNLFA--ELRKLLPTLPP-----DKKLSKIEILR
-----HKEVERRRRRENINTAIN--VLSDLLPVRESSKAAILARAAEYIQKLKETDEANIEKWTQKLLSEQNASQ
FSYPQHKNQSFSASPIDPSMSYVYRAPESFSSINANVPYGRNEYLRRTSLVPNQPEYTGYPYTRNPELRTSHKLAERKRRKEIKELFDDLKDA
-----YSWCLQEAVQFLTTLHAASDTQMKLLYHFQRP-----APAAPAKEPPA
-----ESRSELIYILKSLSYLSWLYERNEKLR-----KQIIAKHEAKT

LAIC-----YISYLNHVLDV-----
LASANEKLQEELGNAYKEIEYMKRVLKKEGIEYEDMHTHKKQENERKSTRSDNPHEA
LPLDKSTKSSKWGLLTRAIQYIEQLKSEQVALEAYVKSLEENMQSNKEVTKGT----
PGAAPQPARSSAKAAAAAVSTSRQPACGLWRPW-----
GSSSSSDPVQEQQGNIRDLVPKELIWELGDGQSGQ-----
```


Examples

- **Dialign** result, first part

```
mml-----  
m-----NSLANNNKLS  
MA-----  
MT-----  
mssyalpsmqptptssiplrqmsqpttsapsnsasstpyspqqvplthnsyplstpssfqhggqtrlppinclaepfnrpqpwhsnaapaSSSPTSATLS  
-----NSDTMELD-----LPPTHSETESGFSDCGGGAGPDgagpgggpgggqarg-----  
TEDEEIHSARKRGYNEEQNYsearkkqrdqgllsqesndgnidsallsegatLKGTQSQYESGLTSNKDEKGSDDedasvaeaavaatvnytdliqgQED  
-----  
TAAHPVHTNAAQVAGSSSSYVYS-----VPPTNSTTSQAsakhsavphrssqfqsttltptst-----DSS  
-----PEPGEPGRK-----  
SSDAHTSNQTNANGEHKDSLNGERAITPSNEGVKP-----NTSLEGMTSSPMESTQQSKNdmliplaehdrg-----  
-----  
STDVSSSDSVSTSASSSNASNTVSVTSPASSSATPLPNQPSQQqflvskndaf ttfvhsvhNTPMQQSMYVPQQQTSHSSGasyqnesanppvqspmqys
```

Examples

- **Dialign** result, second part

```
-----DLQHL---SREERRRRRRATA-----K
-----PEHQqddednddadidlkkdismqpgrgrkPTTLAtt dew-KKQR-----
-----PSTVAVEMLSPEKN-----
-----NDIKEIQTIPGLSEIKEIKGELANVKKR---KRRSKKINKLTDG-----Q
ysqgqpfsyPQHK-----NQSFSASPIDPSMSYVYRAPESFSSINANvpyGRNEYLRRvtslvpnqpeytgpytrnpE
YRTAHATRERIRVEAFNLAFaelRKLLPTL----PPDKKLSKIEILRLAICYISYLNHvIdv-----
-KDSHKEVERRRRRENINTAINVLSDLLP-V----RESSKAA---ILARAAEYIqKLKETDEanieKWTlQKLLSEQNASQLASANEKLQEELGNaykeie
-RLRKPVVEKMRRDRINSSIEQLKLLLeqefarhQPNSKLEKADILEMAVSylKHSKAFaa----Aag-----P
IRINHVSSEKKRRELeraIFDELvAVVPDL----QPQESRSELIiYLKSLSYLSWLYERNE---KLRKQIIAKHEAKTGSSSSSDPVQEQNgnirdlvP
LRTSHKLAERKRRKEIKELFDDLKDALP-L----DKSTKSSKWGLLTRAiQYIEQLKSEQV---ALEAYVKSLEEnmqsnkevtkgt-----
-----
ymkrv1r-----KEGIEYEDMHTHkKqenerkstrsdnphea-----
KSLHQDYSEGYSwclQEAVQFLTLHAasdtqmkllYhfqrppapaapakeppapgaapqparssakaaaaavstsrqpacglwrpw
KELIWELGDGQSgq-----
-----
```

- **SH3 (Src homology 3) domain**

- Often indicative of a protein involved in signal transduction related to cytoskeletal organization
- 5 sequences
- short sequences
- low and diffuse similarity (<25%)

1aboA	P00520	57
1ycsB	P04637	60
1pht	P27986	80
1ihvA	P00383	49
1vie	P12497	51

- When we **aligned secondary structures manually**

```

1aboA    -NLFVALYDfvasgdntlsitkGEKLRVLgynhn-----
1ycsB    kGVIYALWDyepqnddelpmkeGDCMTIIhrede-----
1pht     gYQYRALYDykkereedidlhlGDILTVNkgslvalgfsd
1ihvA    -NFRVYYRDsrd-----pvwkGPAKLLWkg-----
1vie     -drvrkksga-----awqGQIVGWYctnlt-----

1aboA    -----gEWCEAQt--kngqGWVPSNYITPVN-----
1ycsB    -----deiEWWARl--ndkeGYVPRNLLGLYP-----
1pht     gqearpeeiGWLNGYnettgerGDFPGTYVEYIGrkkisp
1ihvA    -----eGAVVIQd--nsdiKVVPRRKAKIIRd-----
1vie     -----peGYAVESeahpgsvQIYPVAALERIN-----
  
```

- **Clustal** result

```

1aboA  -NLFV-ALYDFVASGDNTLSITKGEKLRV-----LGYNHNG
1ycsB  KGVIIY-ALWDYEPQNDDELPMKEGDCMTI-----IHREDED
1pht   -GYQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDGQ
1ihvA  -----NFRVYYRDSRD--PVWKGPAKLL-----WKGEG
1vie   -----DRVRKKSG--AAWQGQIVGW-----YCTNL
  
```

```

1aboA  -----EWCEA--QTKNGQGWVPSNYITPVN-----
1ycsB  EI-----EWWA--RLNDKEGYVPRNLLGLYP-----
1pht   EARPEEI GWLNGYNETTGERGDFPGTYVEYIGRKKISP
1ihvA  -----AVVIQ---DNSDIKVVPRRKAKIIRD-----
1vie   TP-----EGYAVESEAHPGSVQIYPVAALERIN-----
  
```


- Dialign** result

```

1aboA    n-LFVALYDFVASGDNTLSITKGEKLRVL-----
1ycsB    kgVIYALWDYEPQNDDELPMKEGDCMTIIhr----EDEDEI-----
1pht     gyQYRALYDYKKEREEDIDLHLGDILTVNKGSLVALGFSDgqearpeei
1ihvA    --NFRV---YYRDSRDPVWKGPAKLLWKGE GAVVIQDNSDI-----
1vie     -----DRVRKKSGaa-W-----QGQI-----
1aboA    ----GYNhn gEWCEAQTKN GQGWV-----PSNYIt p-----VN
1ycsB    -----EWWARLNDKEGYV-----PRNLLgLYP-----
1pht     gwln GYN-----ETTGERGDF-----PGTYV-EYigRKKIsp--
1ihvA    -----Kv-----V-----PRr-----KAKIIRd-
1vie     -----VGWYCTNLTPEGYAveseahPGSVQ-IYPv-AALERIN
  
```

Sylvain Legrand
Maître de Conférences
UMR CNRS 8198 EVO-ECO-PALEO
Evolution, Ecologie et Paléontologie
Université de Lille - Faculté des Sciences et Technologies
Bât SN2, bureau 208 - 59655 Villeneuve d'Ascq

sylvain.legrand@univ-lille.fr | <http://eep.univ-lille.fr/>
Tél. +33 (0)3 20 43 40 16