

# Recherche de similarités faibles entre séquences homologues

# Limitation des comparaisons deux à deux (BLAST, FASTA, ...)

Seq A	CGRRLLILFMLATCGECDTDSSE ... HICCIKQCDVQDIIRVCC
	:: : ::: :: :
Insuline	CGSHLVEALYLVCGERGFFYTP ... EQCCTSICSLYQLENYCN
	::: : : : :: : :
Seq B	YQSHLLIVLLAITLECFPSDRK ... KRQWISIFDLQTLRPMTA

Comparaisons 2 à 2:

Insuline / Seq A : 25% d'identité

Insuline / Seq B : 25% d'identité

# Insulin gene family: sequence alignment

	B-chain	A-chain
INSL4	ELRG <b>CG</b> PRFGKHL <b>LSY</b> CPMPEKTFTTTPGG... [x] 58	...SGRHRFDPF <b>CC</b> EVI <b>C</b> DDGT <b>SV</b> KL <b>CT</b>
INSL3	REKL <b>CG</b> HHFVRA <b>LVRV</b> CGGPRWSTEA... [x] 51	...AAATNPARY <b>CC</b> LSG <b>CT</b> QQD <b>LL</b> TL <b>CPY</b>
RLN1	VIKL <b>CG</b> RELVRAQIA <b>IC</b> GMSTWS... [x] 109	...PYVALFEK <b>CC</b> LIG <b>CT</b> KR <b>SL</b> AKY <b>C</b>
BBXA	VHTY <b>CG</b> RHLARTLAD <b>LC</b> WEAGVD... [x] 25	...GIVDE <b>CC</b> LRP <b>CS</b> VDV <b>LL</b> SY <b>C</b>
BBXB	ARTY <b>CG</b> RHLADTLAD <b>LC</b> F--GVE... [x] 23	...GVVDE <b>CC</b> FRP <b>CT</b> LDV <b>LL</b> SY <b>CG</b>
BBXC	SQFY <b>CG</b> DFLARTMSI <b>LC</b> WPDMP... [x] 25	...GIVDE <b>CC</b> YRP <b>CT</b> TDV <b>LL</b> KLY <b>CD</b> KQI
BBXD	GHIY <b>CG</b> RYLAYKMAD <b>LC</b> WRAGFE... [x] 25	...GIADE <b>CC</b> LQP <b>CT</b> NDV <b>LL</b> SY <b>C</b>
LIRP	VARY <b>CG</b> EKLSNALK <b>LV</b> CRGNYNTMF... [x] 58	...GVFDE <b>CC</b> RKS <b>CS</b> ISE <b>LQ</b> TY <b>CG</b> R
MIPI	RRGV <b>CG</b> SALADLVDF <b>AC</b> SSSNQPAMV... [x] 29	...QGTTNIVCE <b>CC</b> MKP <b>CT</b> LSE <b>LR</b> QY <b>CP</b>
MIPII	PRGI <b>CG</b> SNLAGFRAF <b>IC</b> SNQNSPSMV... [x] 44	...QRTTNLVCE <b>CC</b> FNY <b>CT</b> PDV <b>V</b> RKY <b>CY</b>
MIPIII	PRGL <b>CG</b> STLANMVQ <b>WL</b> CSTYTTSSKV... [x] 30	...ESRPSIVCE <b>CC</b> FNQ <b>CT</b> VQ <b>EL</b> LAY <b>C</b>
MIPV	PRGI <b>CG</b> SDLADLRAF <b>IC</b> SRRNQPAMV... [x] 44	...QRTTNLVCE <b>CC</b> YNV <b>CT</b> VDV <b>F</b> YEY <b>CY</b>
MIPVII	PRGL <b>CG</b> NRLARAHAN <b>LC</b> FLLRNTYPDIFPR... [x] 86	..EVMAEPSLVCD <b>CC</b> YNE <b>CS</b> VRK <b>L</b> ATY <b>C</b>
ILP	AEYL <b>CG</b> STLADVLS <b>FV</b> CGNRGYNSQP... [x] 31	...GLVEE <b>CC</b> YNV <b>CD</b> YSQ <b>LE</b> SY <b>CN</b> PYS
INS	NQHL <b>CG</b> SHLVEALY <b>LV</b> CGERGFYTPKT... [x] 35	...GIVEQ <b>CC</b> T <b>SI</b> CSLYQ <b>LE</b> NY <b>CN</b>
IGF1	PETL <b>CG</b> AELVDALQ <b>FV</b> CGDRGFYF... [x] 12	...GIVDE <b>CC</b> FRS <b>CD</b> LR <b>LE</b> MY <b>CA</b> PLK
IGF2	SETL <b>CG</b> GELVDTLQ <b>FV</b> CGDRGFYF... [x] 12	...GIVEE <b>CC</b> FRS <b>CD</b> LAL <b>LE</b> TY <b>CA</b> TPA
	* .                    . *	**    *                .    *

# Descripteurs de motifs dans des séquences

- √ Mot exact: e.g. EcoRI site de restriction GAATTC
- √ Consensus: e.g. TATA box: TATAAWR
- √ Expression régulière: e.g. PROSITE pattern des insulines  
C-C-{P}-x(2-4)-C-[STDNEKPI]-x(3)-[LIVMFS]-x(3)-C
- √ Profils :matrice position-spécifique de pondération des substitutions et indels

## Représentation d'un motif par une *matrice de pondération position-dépendante* (exemple)

- Site donneur d'épissage (vertébrés)
- fréquence (%) des 4 bases à chaque position
- log transformation  $\prod$  matrice de pondération

Base	Position								
	-3	-2	-1	+1	+2	+3	+4	+5	+6
A	<b>33</b>	<b>60</b>	8	0	0	<b>49</b>	<b>71</b>	6	15
C	<b>37</b>	13	4	0	0	3	7	5	19
G	18	14	<b>81</b>	<b>100</b>	0	<b>45</b>	12	<b>84</b>	20
T	12	13	7	0	<b>100</b>	3	9	5	<b>46</b>
Cons.	M	A	G	G	T	R	A	G	T

# Recherche d'un motif dans une séquence à l'aide d'un profil

- Calcul des scores de similarité en faisant glisser une fenêtre de la longueur du motif le long de la séquence. Exemple:

A	33	60	8	0	0	49	71	6	15
C	37	13	4	0	0	3	7	5	19
G	18	14	81	100	0	45	12	84	20
T	12	13	7	0	100	3	9	5	46

GAAAGGTGAGTCAT...

GAAAGGTGA

$$S=18+60+8+0+0+45+9+84+15=239$$

.AAAGGTGAG

$$S=33+60+8+100+0+3+12+6+20=242$$

..AAGGTGAGT

$$S=33+60+81+100+100+45+71+84+46=620$$

...AGGTGAGTC

$$S=33+14+81+0+0+49+12+5+19=213$$

....GGTGAGTCA

...etc

# Searching for distantly related homologues in sequence databases

- √ 1- search for homologues (e.g. BLAST)
  - √ 2- align homologues (e.g. CLUSTAL, MEME)
  - √ 3- compute a profile from the multiple alignment
  - √ 4- compare the profile to a sequence database (e.g. MAST, pfsearch)
- 
- √ pfsearch: <http://www.isrec.isb-sib.ch/profile/profile.html>
  - √ MEME/MAST: <http://meme.sdsc.edu/meme/website/>

# PSI-BLAST

- v Position-Specific Iterated BLAST
  - λ 1- classical BLAST search
  - λ 2- compute a profile with significant BLAST hits
  - λ 3- BLAST search based on the profile
  - λ 4 -repeat steps 2-3 up to convergence
  
- v More sensitive than Smith-Waterman
  
- v 40 times faster



# Comparison of a sequence to a database of protein motifs

- √ Databases: PROSITE, PFAM, PRODOM, ..., INTERPRO
- √ Search tools:
  - λ ProfileScan : <http://hits.isb-sib.ch/cgi-bin/PFSCAN>