

Romain GAUTIER

Romain GAUTIER

Rechercher des similitudes de séquence

Recherche de motifs fonctionnels: PROSITE, PFAM , BLOCKS,, PRINTS...

- Profil d'hydropathie, accessibilité au solvant
- Prédiction de structure secondaire
- Prédiction des hélices transmembranaires
- Prédiction des boucles (coudes)
- Modélisation par homologie
- Positionnement des chaînes latérales

Lorsque l'on compare une séquence contre une banque par BLAST, un grand nombre de séquences similaires peuvent être données en sortie. Cela nous permet de poser des hypothèses sur notre séquence inconnue (Fonction, ...):

query:

Mais, comment comparer toutes ces séquences et étudier leur relations ?

Alignement multiple

Alignement multiple


[illegible]

Un alignement multiple se conçoit à partir de 3 séquences et peut concerner jusqu'à plusieurs centaines de séquences.

Entrée : k séquences (pas forcément la même longueur)

Sortie: un tableau contenant les k séquences, avec des indels

Sortie: un tableau contenant les k séquences, avec des indels



Démarche

Alignement 2 à 2

2 séquences quelconques

Détecter une similarité **syntactique**

Il y a-t-il une **fonction commune** ?

Alignement multiple

Famille de séquences avec la même **fonction**

A quelle **conservation syntactique** cela correspond-il ?

Alignement 2 à 2

2 séquences quelconques

Détecter une similarité syntaxique

Il y a-t-il une fonction commune ?

Alignement multiple

Famille de séquences avec la même fonction

A quelle conservation syntaxique cela correspond-il ?

Exemple: Motif doigt de Zinc (JS Varré H Touzot)

Voici un alignement multiple de séquences protéiques:

```
TTY1_HUMAN  YCYPFGDGNKKFAGTMLEKHIL--TH- 25
TF1B_BUFAF  YRCPRNCRDITYTTKFLKELHIL--TPR 26
ENT7_HUMAN  YTCPEHRCURGPTATNENHVR--IH- 25
ZNF6_HUMAN  FRCVQYQCDRLITTAHLEKHYER--AK- 25
P44_XENBO  YRCYSIEDQYTVSPNTALQTLHK--KH- 25
TSH_BROME  FRCVW--CKQSPPTLEALTTHMKDEK-- 25
ZF1X_XENLA  FRCSE--CRRSTFSDSLTAMR--KH- 23
EVL1_HUMAN  YRCYV--CRRSPISINLQKQVNR--IH- 24
TRA1_CAEEL  YECFADCEKAFSNASIRAKQNR--TH- 26
TF1A_BUFAF  CCCTENCLNAPLATFSLRFRFKH--AH- 26
SVC1_BROME  FPCNY--CRRDTFVFNRLGLTER--RH- 24
Z02_9_XENL  FVCVT--COKTKYKHLNLTALIS--H- 23
Z058_XENLA  FVCVE--CHLSFPAIANLRSHQLK--H- 23
Y4B8_CAEEL  YRCYV--CKKDISSEKLLTDMPE-QM- 25
BOM0_HUMAN  FQCDI--CKTKYFNACSVKLIHKN-MH- 24
SUNW_BROAN  YACKI--CKDKPTFSLRGLRFGY--SSC 25
ZB10_HUMAN  YKQNG--CKIIPGNSIFVGI--AK- 23
P43_XENBO  LKCSVPQCKRFRFKKRLRHYH--KH- 25
IKAR_MOUSE  FRCNM--CQYHSGDYRFSESHITRGHE- 25
```

Voici un alignement multiple de séquences protéiques:

```

T7Y1_NEMO      YVCFDQKQKKFAGSTLNKSHL--TH-
T7B1_BUFAM     YKRECDNRDTH2TFLNKLSHL--TFR-
Z777_NEMO      YK7FCEHCGRGT7T7YKXNHRV--LH-
Z777_BUFAM     YK7FCEHCGRGT7T7YKXNHRV--LH-
P441_NEMO      YKSCYEDQ7QVSP7TAL7H7LKH--KH-
TSL_DROME      PRCWV--CRQ5PT7LEL7HMKDKSH-
EV11_NEMO      XHJLH1_KFCLSA--C88ST7NSL7HDL7HAR--KH-
EV11_NEMO      YK7ACKY--CDR8SF7L7NKH7VRVH--LH-
TAR1_CAEEL     YK7FCEYK7FAS7NAS7RKH7NRH--TH-
TAR1_BUFAM     YK7FCEYK7FAS7NAS7RKH7NRH--TH-
Z029_DROME     YK7FCEYK7FAS7NAS7RKH7NRH--TH-
Z029_NEMO      YK7FCEYK7FAS7NAS7RKH7NRH--TH-
Z029_XENL       FVCTV--CDPT7K7YKH7GL7NL7HLS--RH-
Z058_NEMO      FVCTV--CH7LS7FAG7L7NL7SH7QL--RH-
YQ08_CAEEL     YK7CTV--CH7LS7S7S7L7L7H7P7QH--QH-
BASO_NEMO      FQCDI--C77F7AC7M7AC7S7I7H7K7N--KH-
BASO_BUFAM      FQCDI--C77F7AC7M7AC7S7I7H7K7N--KH-
Z100_NEMO      YK7CQI--C77F7AC7M7AC7S7I7H7K7N--KH-
P413_NEMO      LK7S7V7C7P7C7R7F7R7K7L7H7H7S--RH-
TKAR_MOUSE     F6C3M--C7Y7H7D7R7F7F7S7H7T7R7G7-

```

Motif doigt de Zinc (JS Varré H Touzet)

Regardons les résidus totalement conservés:

```

TY11_HUMAN  YVCFDGCNCKFQAQTNLSKSL--TH_25
TF3B_MOUSE  YRCPRNCDCRTYTTKPNLSKSL--TFH_26
ZNF77_HUMAN  YCFEPHRCCLGTAFTNYYKNRV--IKH_25
ZNF76_HUMAN  FRCYQYGCQGLYLTAMBLAYR--AM_25
P44_XENRO   YRCSYECQTVFPTFALQWLE--KH_25
ZNF_KRNO     FRCVW--CQSPFTLEALTITNEDKSD_25
XFLN1_XENLA  FRCSE--CGRSPTNINSLTAMR--KH_23
KV11_HUMAN  YRCKY--CGRSFSISNMLGRHYN--IKH_24
TRAI_CAEAL  YCEFPADCAKAFSNADRAKNKNS--TH_26
TF3A_MOUSE  CKCTENCTGCLTAFTNLSHLYR--AR_26
SRYC_BROME  FPCNY--CRPFTFPPNLAATL--RH_26
ZIG2_XENL  FVCTV--CGKTKYKYGKGLTFLRS--H_23
ZIG5_XENLA  FVCTE--CNLSFAPLAKLSHGLG--H_23
YGOB_CAEAL  YCCTV--CKKLSISSELTLEFQ--QHR_25
BAGO_HUMAN  FQCDI--CKKTFPACNACVSLIKHN--H_24
SDHM_DROAN  YAKCI--CKDPTFSYMLKRGKILRY--SSC_25
ZNF10_HUMAN  YCKWQ--CGLISGNSFPIVYGL--AM_23
P43_XENRO   FSKSVPGCKRFRKCKKALLRYS--EH_25
IKAR_MOUSE  FECNN--CQYSGQRVTFSSSHTRGEM--25

```

On peut établir une carte par position des résidus conservées :

Weblogo : <http://weblogo.berkeley.edu/logo.cgi>

Regardons les résidus totalement conservés:

TY11_HUMAN YVPCFDDKGGKFAFQNLKSLIL--TF
 ZP38_BUFAN YVPCFDDKGGKTFNFKLNLIL--TFP
 ZN77_HUMAN YVPCFDDKGGRTFATYVNLNVR--IK
 ZN77_MOUSE YVPCFDDKGGRTFATYVNLNVR--IK
 P44_KEMBO YVPCFDDKGGRTFATYVNLNVR--IK
 ZN3_MOUSE FCWCV--CGKQFFPTLALTHMKDKSK
 XFIN_XENLA FCWCV--CGRSFTNSHDLTAMDE--XU
 EV11_HUMAN YKCKY--CGRSFSLNSNLGRVNV--IT
 TRAL_CAEL YKCKY--CGRSFSLNSNLGRVNV--IT
 ZN1_MOUSE YKCKY--CGRSFSLNSNLGRVNV--IT
 SKYC_DROME FCKMY--CGRTKFKKGLNLTNR--IS
 Z02-9_XENL FCKTY--CGRTKFKKGLNLTNR--IS
 Z058_XENLA FVCTE--CHKDISFGLNSLSRSG--QH
 X08L_CAEL YKCKY--CHKDISFGLNSLSRSG--QH
 ZN1_MOUSE YKCKY--CHKDISFGLNSLSRSG--QH
 SDRK_DROME YACKI--CHKDISFGLNSLSRSG--QH
 ZN10_HUMAN YKCKY--CGIIFSPLSPFVIVQI--AS
 P43_KEMBO LKSLVPCQKSPFRKKALRIWS--EK
 IKAR_MOUSE FCNIN--CQYHSCDRTSSFTTITR--

On peut établir une carte par position des résidus conservées :
WebLogo : <http://weblogo.berkeley.edu/logo.cgi>

[illegible]

TYTL_BIFAN YVCFPQDHCNAGQNSLHSHLT-
 T273_HIFAN YVGFQDHCNCTITFNFKSLHSHLT-
 Z877_HIFAN YVGFQDHCNCTGTSATNFKNVR-
 T287_HIFAN FRCYQDHCNCTITASHLHVR-AN
 P44_XENID YKCYDCHVYSPWPTATQTLK-AN
 T273_XENID YKCYDCHVYSPWPTATQTLK-AN
 T287_XENID YKCYDCHVYSPWPTATQTLK-AN
 E711_HIFAN YKCYDCHVYSPWPTATQTLK-AN
 T287_CAEK YKCYDCHVYSPWPTATQTLK-AN
 T273_BIFAN CXCXCHNCHVYSPWPTATQTLK-AN
 T287_BIFAN YKCYDCHVYSPWPTATQTLK-AN
 T287_XENID YKCYDCHVYSPWPTATQTLK-AN
 Y058_CAEK YKCYDCHVYSPWPTATQTLK-AN
 N430_HIFAN YKCYDCHVYSPWPTATQTLK-AN
 Z877_BIFAN YKCYDCHVYSPWPTATQTLK-AN
 P44_XENID YKCYDCHVYSPWPTATQTLK-AN
 IAKR_MOUSE FRCYDCHVYSPWPTATQTLK-AN

L'alignement multiple permet de mieux comprendre l'organisation d'un motif particulier et permet de modéliser le motif : (Prosit)

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

[illegible]

A partir de l'alignement multiple, on peut déterminer la séquence consensus:

On attribut à chaque position l'acide aminé ou le nucléotide qui est le plus souvent retrouvé (totalement conservé ou partiellement).

[illegible]

Alignement multiple: principe

L'approche la plus courante consiste à aligner progressivement des paires de séquences:

Alignements itératifs et progressifs

Algorithme Needleman & Wunsch 2 à 2 traditionnel répété n fois

L'approche la plus courante consiste à aligner progressivement des paires de séquences:

Alignements itératifs et progressifs

Algorithme Needleman & Wunsch 2 à 2 traditionnel répété n fois

Alignement multiple: principe

L'approche la plus courante consiste à aligner progressivement des paires de séquences.

Les différentes approches se distinguent par:

- La façon de choisir la paire initiale de séquences

Pour progresser dans l'alignement, les programmes vont:

- Soit aligner chaque séquence les unes après les autres à un alignement unique enrichi à chaque étape
- Soit créer des sous-familles de séquences d'abord alignées au sein de ces familles puis entre les familles.

La méthode de pondération des alignements individuels des paires de séquences et des alignements cumulés.

Alignement basé sur un arbre

Idée: reconstruire l'alignement multiple à partir d'un **arbre guide** (clusters)

- feuilles : séquences
- noeuds : alignements

Partir des feuilles puis remonter dans l'arbre

Utilisation de la technique de **profil alignement** → produire un seul alignement multiple avec deux.

CLUST_{er} + **AL**ignement **CLUSTAL**

Thompson et al. 1994

ClustalW EBI <http://www.ebi.ac.uk/cluster/index.html>

ClustalW est l'un des programmes les plus utilisés pour l'alignement progressif.

Etape 1: Alignements globaux 2 à 2

Etape 2: Regroupements des alignements (clusters), construction arbre guide

Etape 3: Alignement multiple obtenu par combinaisons des alignements 2 à 2 (profils)

ClustalW exemple

4 séquences

s1	cgatgagtcattgtgactg
s2	cgagccattgtagctactg
s3	cgaccattgtgactacctg
s4	cgatgagtcactgtgactg

Jeu de score:

- Indel= -2
- Substitution=-1
- Identité= 1

ClustalW étape 1: Calcul des scores

Les alignements de toutes les paires de séquences sont réalisés puis le programme génère une matrice de distances décrivant leur taux de similitude.

s1	cgatgagctcattgt-g-actg	s2	cgagccattgtactga-ctg
s2		s3	
	cga-g--cattgtgactgactg		cga--ccattgtgactgactg
s1	cgatgagctc-tgactg	s2	cga-g--ccattgtgactgactg
		s3	
s3	cagacca-ttgtgactgactg	s4	cgatgagctcattgt-g-actg
s1	cgatgagctcattgtgactg	s3	cgaccattgtgactgactg
		s4	
s4	cgatgagctcactgtgactg	s4	cgatgagctcactgtgactg

Tableau des scores d'alignement:

	s1	s2	s3	S4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	

Pour N séquences:
 $N(N-1)/2$ calculs

ClustalW étape 2: Construction arbre

A l'aide de la matrice de scores, un arbre guide (ou dendrogramme) est ensuite construit par un algorithme dit de neighbor-joining:

	s1	s2	s3	s4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	

s1

s2

s3

s4

Regroupement des séquences suivant leur similitude à partir de la matrice des scores 2 à 2.

ClustalW étape 2: Construction arbre

A l'aide de la matrice de scores, un arbre guide (ou dendogramme) est ensuite construit par un algorithme dit de neighbor-joining:

	s1	s2	s3	s4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	

s2
s3
s1
s4

ClustalW étape 2: Construction arbre

A l'aide de la matrice de scores, un arbre guide (ou dendogramme) est ensuite construit par un algorithme dit de neighbor-joining:

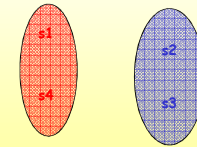
	s1	s2	s3	s4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	

s2
s3
s1
s4

ClustalW étape 2: Construction arbre

A l'aide de la matrice de scores, un arbre guide (ou dendogramme) est ensuite construit par un algorithme dit de neighbor-joining:

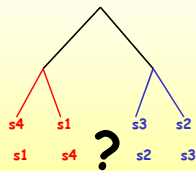
	s1	s2	s3	s4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	



ClustalW étape 2: Construction arbre

A l'aide de la matrice de scores, un arbre guide (ou dendogramme) est ensuite construit par un algorithme dit de neighbor-joining:

	s1	s2	s3	s4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	



ClustalW étape 2: Construction arbre

A l'aide de la matrice de scores, un arbre guide (ou dendogramme) est ensuite construit par un algorithme dit de neighbor-joining:

	s1	s2	s3	s4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	

Alignement s1s2 est plus proche que s1s3

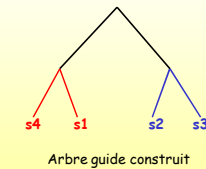
Alignement s4s2 est plus proche que s4s3

Matrice de score -> matrice de distance

ClustalW étape 2: Construction arbre

A l'aide de la matrice de scores, un arbre guide (ou dendogramme) est ensuite construit par un algorithme dit de neighbor-joining:

	s1	s2	s3	s4
s1		2	0	17
s2	2		14	0
s3	0	14		-1
s4	17	0	-1	



Etape 3: Construction alignement final

ClustalW aligne les séquences en se servant de l'arbre guide: chaque paire de séquences situées sur une même branche extérieure de l'arbre est alignée par programmation dynamique.

Les alignements partiels permettent de constituer des **profils**, représentés sous forme de tableau dans lequel sont données pour chaque position la fréquence observée de chaque lettre.

L'algorithme aligne ensuite les profils associés par un même nœud de l'arbre. Cet alignement de séquences puis de profils se poursuit de façon récursive jusqu'à l'alignement final complet depuis les branches de l'arbre vers la racine.

Etape 3: Construction alignement final

L'alignement et création des profils:

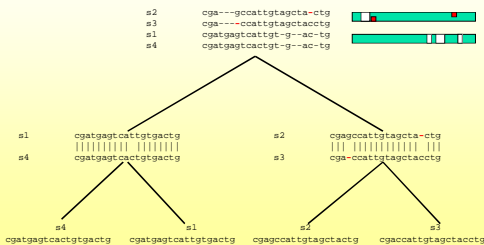


Etape 3: Construction alignement final



Etape 3: Construction alignement final

Once a gap, always a gap



ClustalW

ClustalW est optimisé pour les protéines:

Pondération des séquences en fonction de leur sur/sous représentation

Adaptation des matrices de similitudes au fil de l'algorithme en fonction de la divergence des séquences à aligner
 Blosom 80 pour aligner des séquences proches
 Blosom 50 pour aligner des séquences distantes

Pénalités de gaps spécifiques à chaque résidu.

Par exemple, les Glycines sont davantage susceptible d'avoir un gap que les Valines.

Pénalités de gaps réduites dans les régions hydrophiles

Encourage la formation de gaps dans des boucles plutôt que dans des régions structurées.

Pénalités de gaps augmentées dans le voisinage d'autres gaps

Evite la formation de petits gaps voisins, au profit de longs gaps.

Autre méthode

A partir des alignements locaux

Idée: repérer des similitudes locales fortes entre les séquences (les diagonales du dotplot par exemple)

Incorporer les diagonales dans l'alignement multiple

Conséquence: les gaps inter-diagonales sont moins importants

DIagonal + ALIGNement

DIALIGN

DIALIGN Morgenstern et al. 1996

DIALIGN sur Pasteur: <http://bioweb.pasteur.fr/seqanal/interfaces/dialign2-simple.html>

Alignements des paires avec optimisation des poids des diagonales

Tri des diagonales selon leur poids et leur chevauchement

DIALIGN : DNA and protein sequence alignment based on segment-to-segment comparison ([Morgenstern, Dress, Werner](#))

Reconstruction gloutonne
Insertion des diagonales par poids décroissants
Vérification de la consistance avec les diagonales déjà introduites

Recommencer ...

DIALIGN exemple

3 séquences:
YIAFLFAWDD
LACFIFGS
SWEDFMFAED

Etape 1 : Détection des diagonales dans les paires de séquences:

DIALIGN exemple

3 séquences:
YIAFLFAWDD
LACFIFGS
SWEDFMFAED

Etape 1 : Détection des diagonales dans les paires de séquences:

DIALIGN exemple

3 séquences:
YIAFLFAWDD
LACFIFGS
SWEDFMFAED

Etape 1 : Détection des diagonales dans les paires de séquences:

DIALIGN exemple

3 séquences:
YIAFLFAWDD
LACFIFGS
SWEDFMFAED

Etape 2 : Sélection d'un ensemble cohérent de diagonales pour construire l'alignement:

Pas de croisement
Pas de chevauchements
score maximal

yIA-FLFAWd
-LAcFIFgs--
swedFMFAED-

CLUSTAL vs DIALIGN

Exemple (C. Notre-Dame)

GARFIELD THE LAST FAT CAT

GARFIELD THE FAT CAT

GARFIELD THE VERY FAST CAT

THE FAT CAT

CLUSTAL vs DIALIGN

Alignement fourni par ClustalW:

```
seq2xx1 GARFIELDTHE----FAT-CAT
seq4xx3 -----THE----FAT-CAT
seq1xx0 GARFIELDTHELASTFAT-CAT
seq3xx2 GARFIELDTHEVERYFASTCAT
```

Clustal 1ère version:

```
GARFIELDTHELASTFA-TCAT
---GARFIELDTHEFA-TCAT
GARFIELDTHEVERYFASTCAT
-----THEFA-TCAT
```

Alignement fourni par Dialign2:

```
seq1 1 GARFIELDTHELASTFA-TCAT
seq2 1 GARFIELDTHE----FA-TCAT
seq3 1 GARFIELDTHEVERYFASTCAT
seq4 1 -----THE----FA-TCAT
```

Quelles méthode utiliser

Cela dépend du type de séquences à aligner !

BaliBase : base de données d'alignements multiples et de benchmarks
<http://www-bio3d-igbmc.u-strasbg.fr/balibase/>

Plusieurs références!

Ref1 : Séquences équidistantes avec différents niveaux de conservation

Ref2: Protéines homologues + une séquence orpheline

Ref3: Sous-groupes avec moins de 25% d'identité entre les groupes

Ref4: Extensions N/C terminales

Ref5: Insertions internes

Ref6: Répétitions internes

Ref7: Protéines transmembranaires

Ref8: Permutations de domaines

Ref1, 2 3: Préférer Clustal à Dialign

Ref4 et 5: Préférer Dialign à Clustal

Quelles méthode utiliser

Plus les séquences sont divergentes, moins le résultat est fiable.
Quand le taux d'identité est supérieur à 35%, toutes les méthodes sont satisfaisantes.

Twilight Zone : 10-20% d'identité

Aucune méthode n'assure un alignement avec plus de 50% de correction.

Clustal a tendance à autoriser moins de gaps que Dialign

Similitude locale : Dialign

Similitude globale : Clustal

Existe d'autres méthodes : MultiAlign, Tcoffee, etc...

Pas de méthode universelle

Pas de confiance aveugle vis-à-vis du résultat obtenu

Visualiser des alignements multiples

Il existe des outils simples pour visualiser ou modifier un alignement multiple:

SeaView : <http://pbil.univ-lyon1.fr/software/seaview.html>



BioEdit : <http://www.mbio.ncsu.edu/BioEdit/bioedit.html>

Autre type d'alignement: structuraux

On peut aussi vouloir aligner non pas des séquences mais des structures protéiques (Cela permet de mieux comprendre les zones clés, structurellement conservées).

TM-Align : <http://zhang.bioinformatics.ku.edu/TM-align/>
Permet d'aligner 2 structures au format PDB

Isuperpose : <http://bioserv.rpbs.jussieu.fr/cgi-bin/iSuperpose>
Permet d'aligner plusieurs structures (format PDB). Possibilité de fournir l'alignement ou pas.

...

Exercices

Exo 1: A partir du jeu de séquences [seq1](#), lancer un ClustalW et analyser l'alignement obtenu.

Regarder les différentes sorties proposées par ClustalW (lancer l'applet JalView)

Sauvegarder les différents fichiers de sortie.

Télécharger et installer SeaView et Njplot.

Regarder l'alignement avec SeaView et visualiser l'arbre guide avec NjPlot.

Exercices

Exo 2: Analyse des domaines SH3

Intro: SH3 (Src homology 3) domains are often indicative of a protein involved in signal transduction related to cytoskeletal organization. The SH3 domain has a characteristic fold which consists of five or six beta-strands arranged as two tightly packed anti-parallel beta sheets. The linker regions may contain short helices.

Voici 5 protéines à aligner: [seq2](#)

Ce sont des séquences courtes avec des similitudes faibles et diffuses (<25%).

Tester les différents alignements obtenus avec ClustalW, Dialign2, Multialign.

Proposer des régions ou résidus clés pour ce type de domaine.

Exercices

Exo 2: Les structures de ces protéines sont connues:

Voici pour chaque protéine les éléments de structures secondaires:

```
> laboA 58
HLFVALYDFVARGENTLILITGKELVLYVNNHGMCAQTRNGQMVPSNYITPVH
CSEEECCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
> lycsB 60
KUVLYALMDYEPQNDDELPHKEGDGMTIIHREDEDEIWWWARLNDEKGYVPRNLLGLY
CSEEECCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
> lpht 83
AEQYQFALYDYKKEEREDIDLHLGIDLTVNKGSLVALGFSQGEARPEIGWLNQYNETTGKRGDFFGTVEYIGREKISPP
CSEEECCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
> llnvA 52
MIQNFVYVYDSRDFVWKGPAKLLMKEGAVVIQNSDIKVVPRKAKIIRD
CSCSEEECCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
> lyie 60
PQNTATFQMDRVRKESGAAMQQQIVWYCTNLTPFGYAVESRAHPGSVQIYPVVALERIN
CSCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCCC
```

Voici les protéines alignées structurellement visualisées avec Pymol [fit](#)

Que peut-on dire au niveau de la structure ?
(par rapport à l'alignement ?)