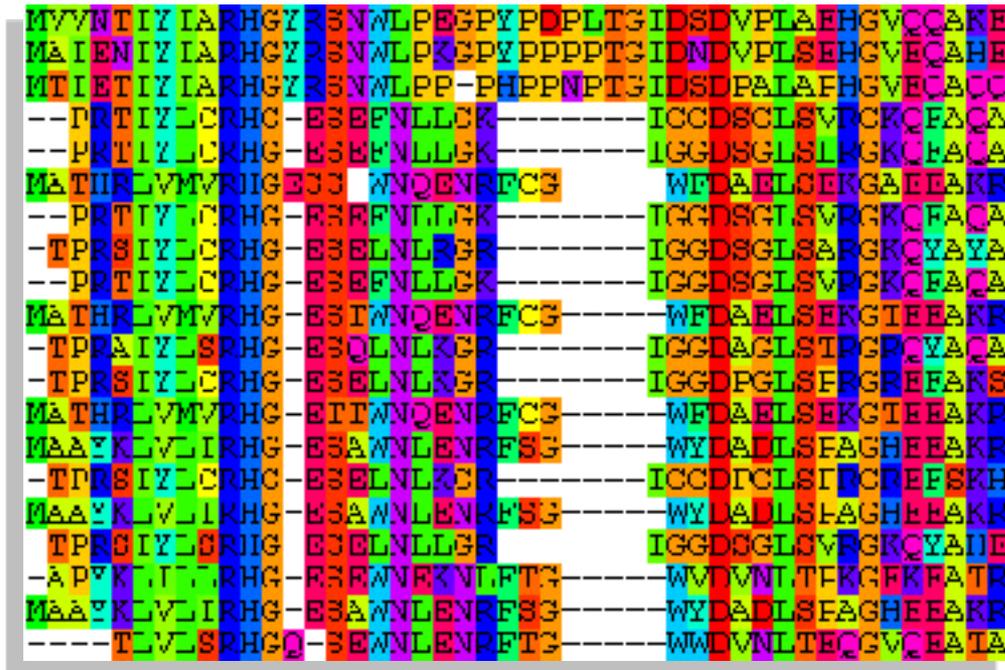


Bioinformatique pour la Biologie Structurale : Quelles informations tirer de l'analyse d'une séquence ?



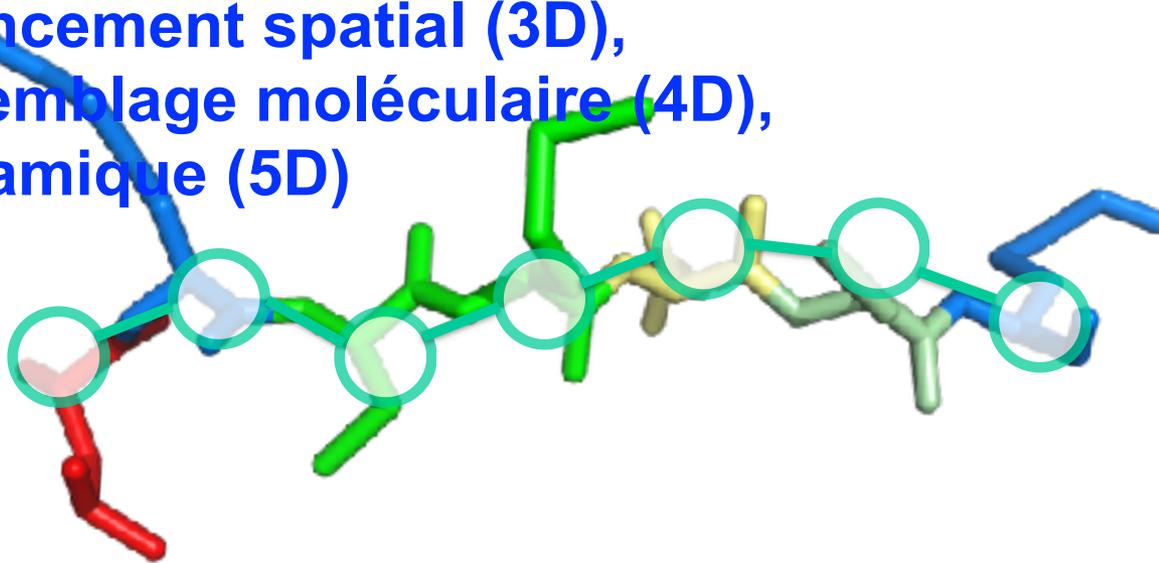
Marie-Hélène Le Du
CEA Saclay, Institut Joliot
91190 Gif sur Yvette Cedex
marie-helene.ledu@cea.fr

Qu'est-ce que la structure d'une macromolécule ?

- Séquence (1D),
- Structures secondaires (2D),
- Agencement spatial (3D),
- Assemblages moléculaires (4D),
- Dynamique (5D)

Qu'est-ce que la structure d'une macromolécule ?

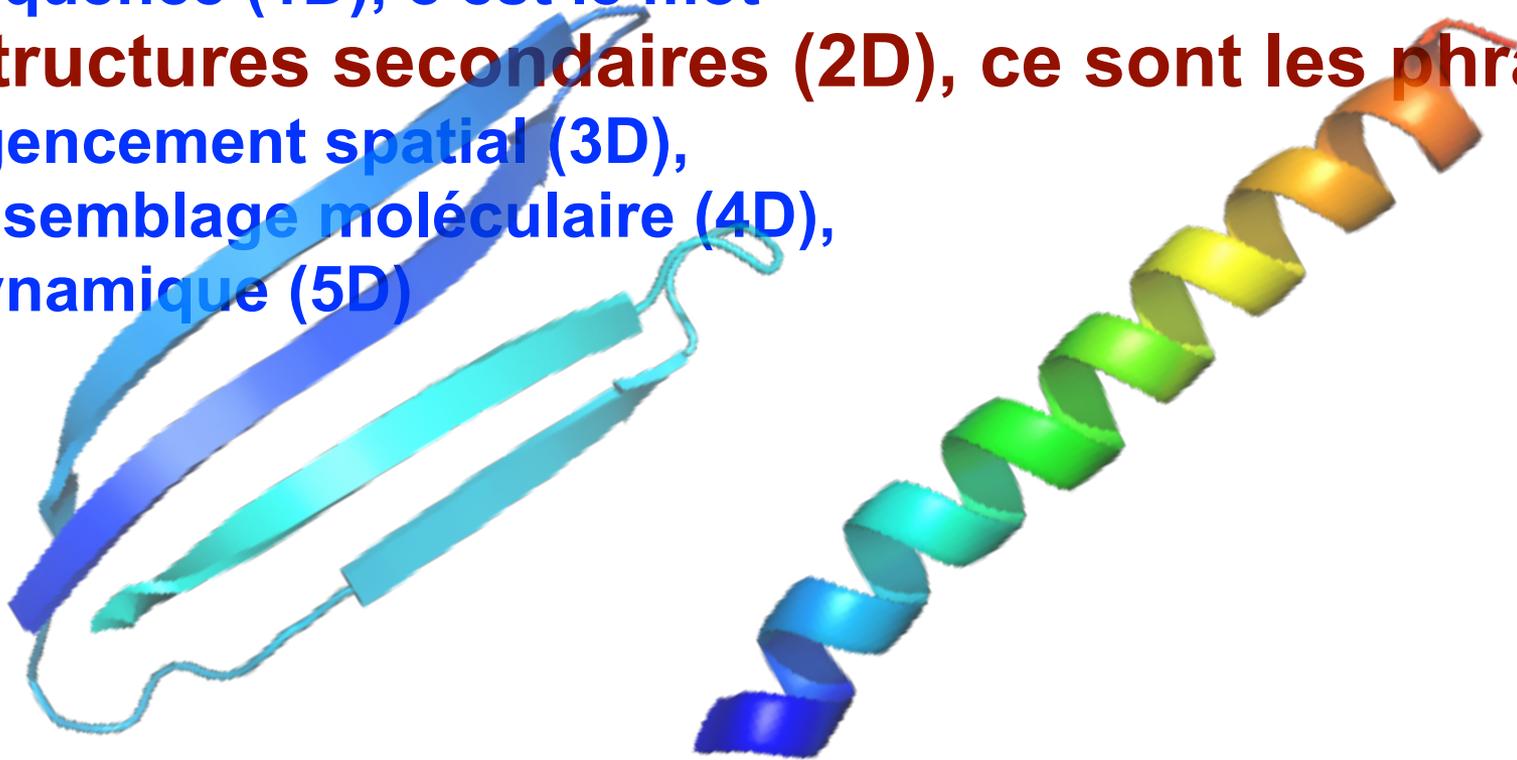
- **Séquence (1D), c'est le mot**
- Structures secondaires (2D),
- Agencement spatial (3D),
- Assemblage moléculaire (4D),
- Dynamique (5D)



Ordre des acides aminés => constitue une chaîne polypeptidique avec un squelette, l'enchaînement peptique, et des chaînes latérales qui dépendent des acides aminés

Qu'est-ce que la structure d'une macromolécule ?

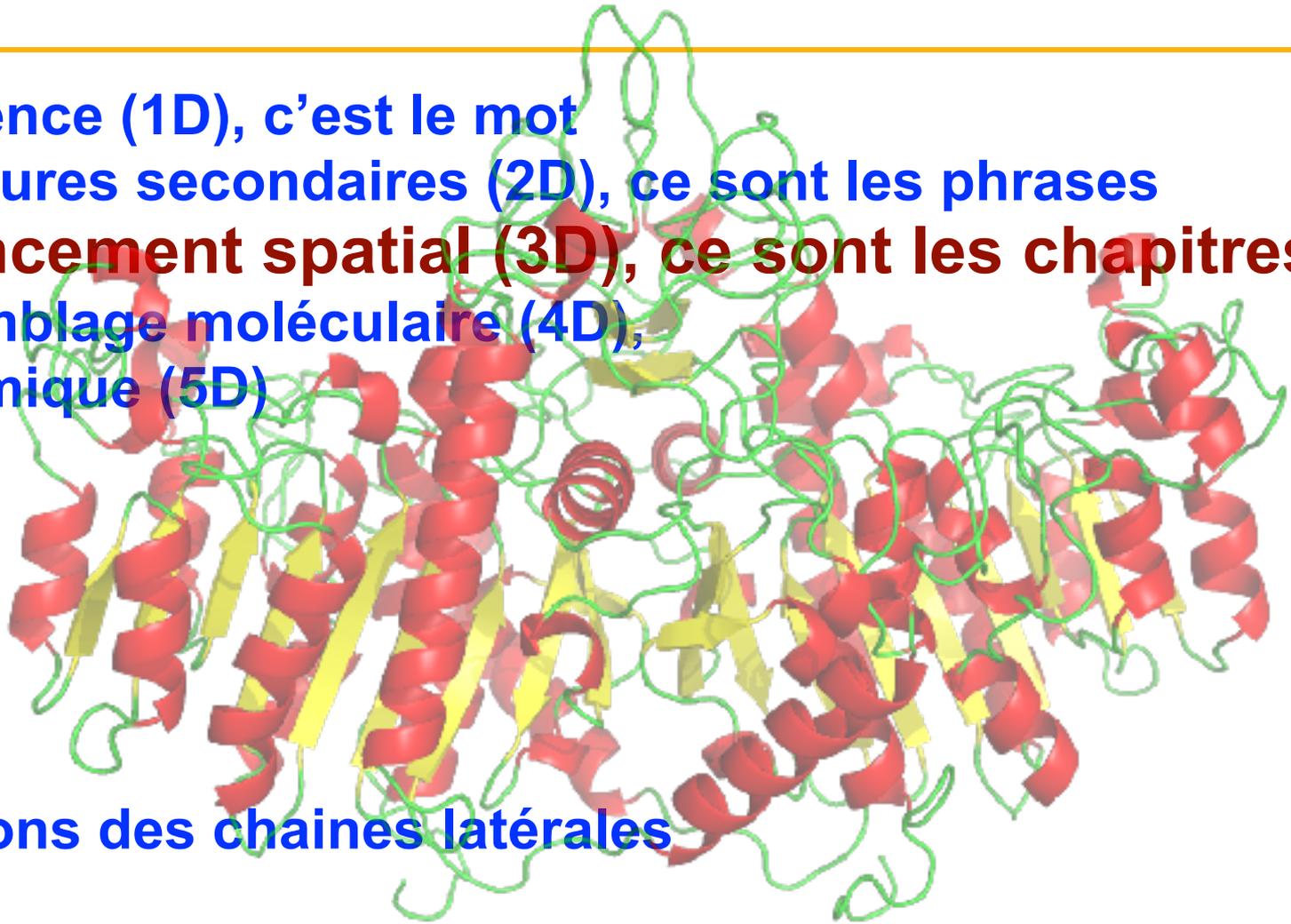
- Séquence (1D), c'est le mot
- **Structures secondaires (2D), ce sont les phrases**
- Agencement spatial (3D),
- Assemblage moléculaire (4D),
- Dynamique (5D)



**Interactions du squelette;
exemples: feuillets bêta, hélices alpha**

Qu'est-ce que la structure d'une macromolécule ?

- Séquence (1D), c'est le mot
- Structures secondaires (2D), ce sont les phrases
- **Agencement spatial (3D), ce sont les chapitres**
- Assemblage moléculaire (4D),
- Dynamique (5D)

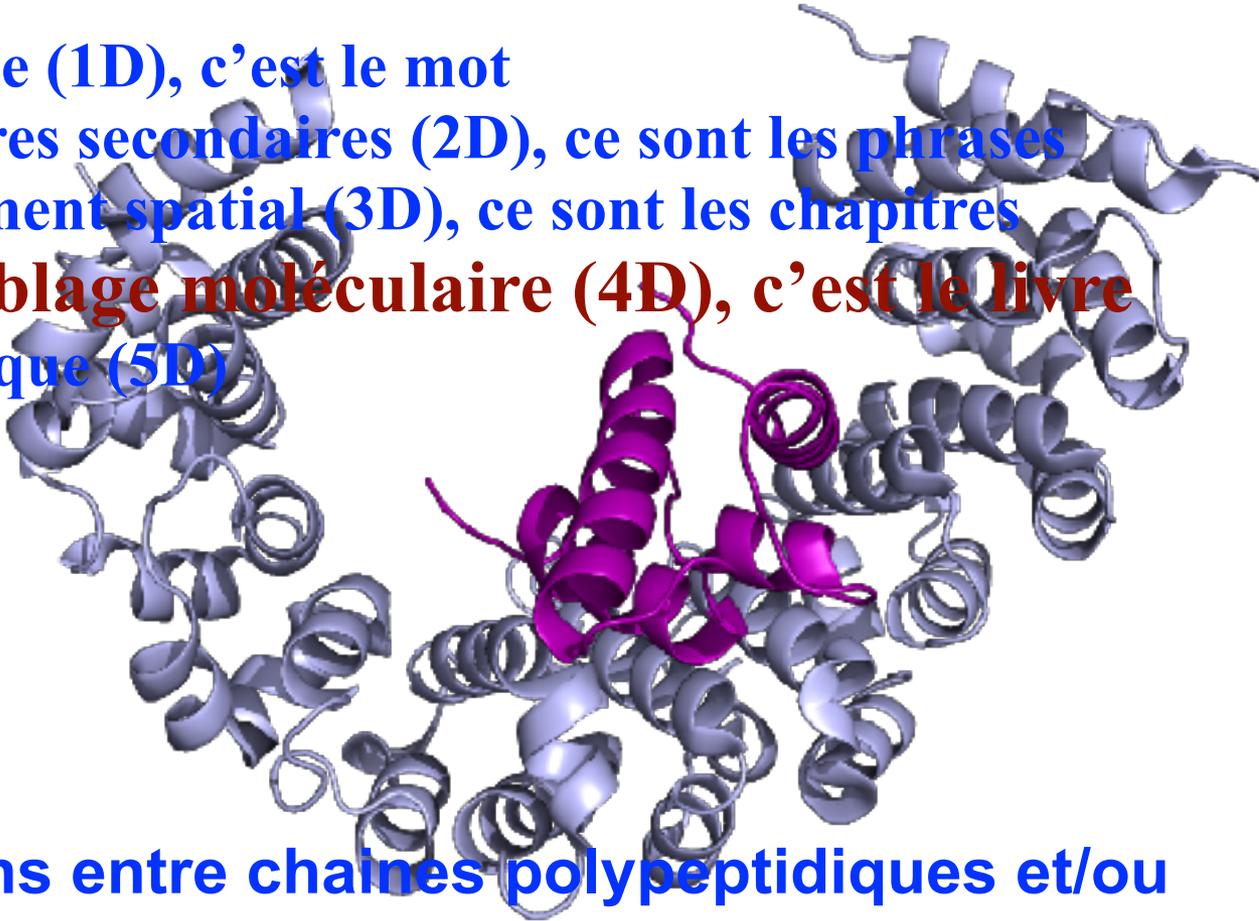


Interactions des chaînes latérales

Llinas et al., Protein Sci, 2006

Qu'est-ce que la structure d'une macromolécule ?

- Séquence (1D), c'est le mot
- Structures secondaires (2D), ce sont les phrases
- Agencement spatial (3D), ce sont les chapitres
- **Assemblage moléculaire (4D), c'est le livre**
- Dynamique (5D)

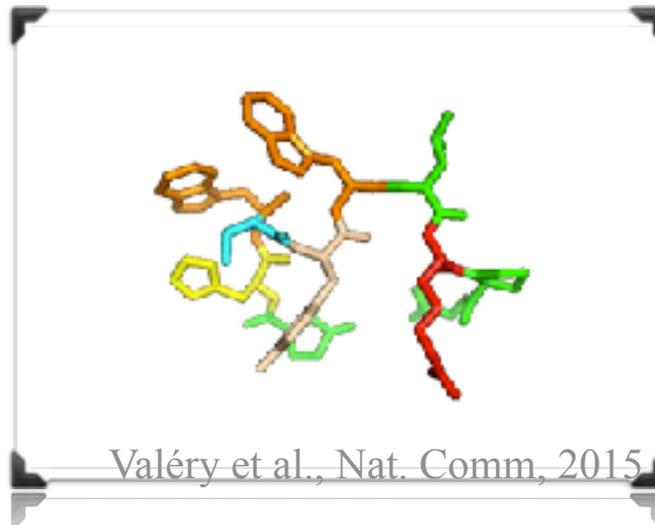


Interactions entre chaînes polypeptidiques et/ou nucléotidiques

Barrault et al., PNAS, 2012

Qu'est-ce que la structure d'une macromolécule ?

- Séquence (1D), c'est le mot
- Structures secondaires (2D), ce sont les phrases
- Agencement spatial (3D), ce sont les chapitres
- Assemblages moléculaires (4D),
- **Dynamique (5D), c'est la saga**



Interactions avec le milieu

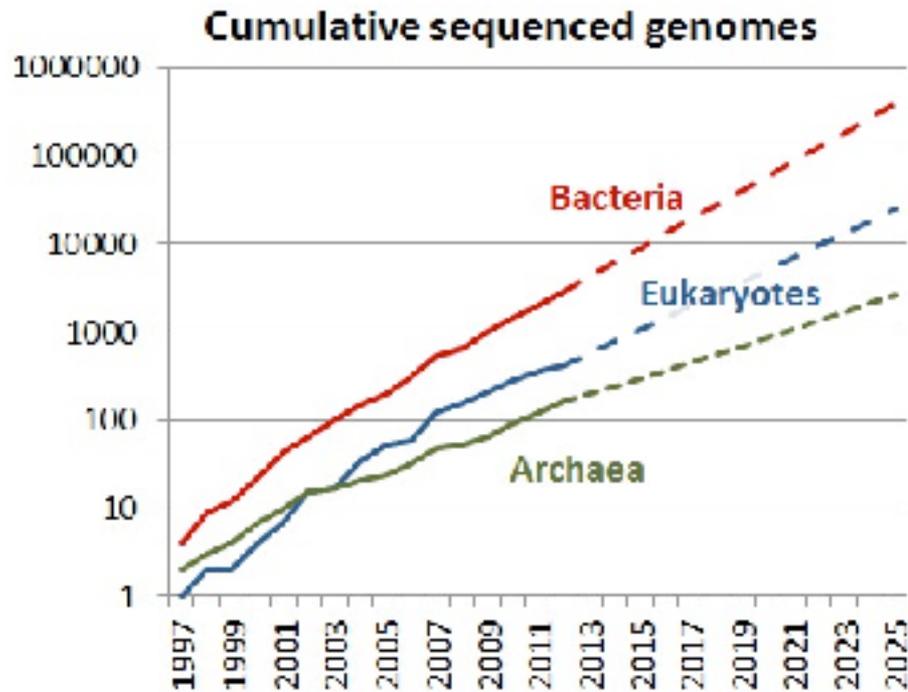
La Question du Repliement des Protéines

Dans les conditions physiologiques, une protéine adopte une structure unique

KVFGRCELAA AMKRHGLDNY
RGYSLGNWVC AAKFESNFNT
QATNRNTDGS TDYGILQINS
RWWCNDGRTP GSRNLCNIPC
SALLSSDITA SVNCAKKIVS
DGNGMNAWVA WRNRCKGTDV
QAWIRGCRL



Contexte de la Prédiction de la Structure des Protéines



Nous sommes à l'ère post-génomique

**=> Beaucoup de séquences, peu de fonctions,
peu de structures**

Contexte de la Prédiction de la Structure des Protéines

Structure:

Méthodes expérimentales :

Cristallographie

RMN

Cryo-EM

=> Quelques structures par jour à travers le monde

... incompatible avec le nombre de gènes séquencés

Forte demande en prédiction de structure:

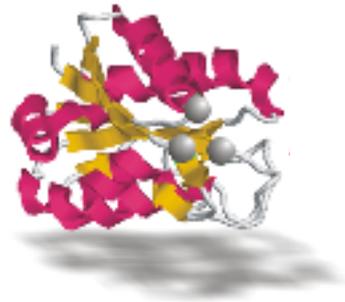
plus de 30,000 gènes chez l'homme.

nombreux génomes seront séquencés au cours des 10 prochaines années.

Pourquoi prédire la structure ?

```
KVFGRCRLAA AMKRHGLDNY  
IGYSLGNWVC AAKPRSNFNT  
QATNRHTDGS TDYGILOINS  
RWWNDGRET GSRNLCNIEC  
SALLSSDIEA SVMCAKRIYS  
DGNEMMAWVA WNRKCKGIDV  
QAWIRGRL
```

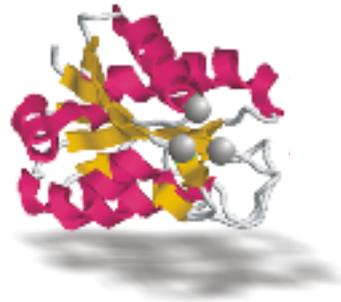
Structure prediction



Et avant une étude de biologie structurale: Pourquoi prédire la structure ?

```
KVFGRCRLAA AMKRHGLDNY  
IGYSLGNHVC AAKPRSNFNT  
QATNRHTDGS TDYGIHQINS  
RWWNDGRTP GSRNLCNTEC  
SALLSSDIEA SVMCAKRIYS  
DGNEMAMVA WNRKCKQIDV  
QAWIRGCRLL
```

Structure prediction



► Délimitation des
domaines structurés

Méthodes actuelles de prédiction de structure 3D

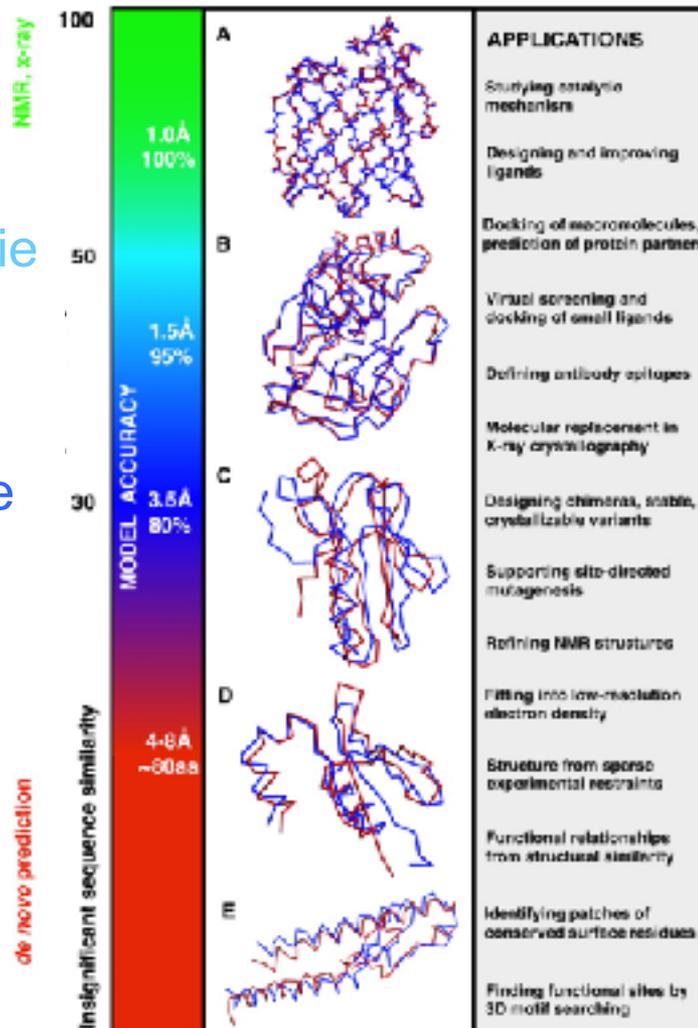
Expérimentales

Modélisation par homologie ou comparative

Beginning of twilight zone

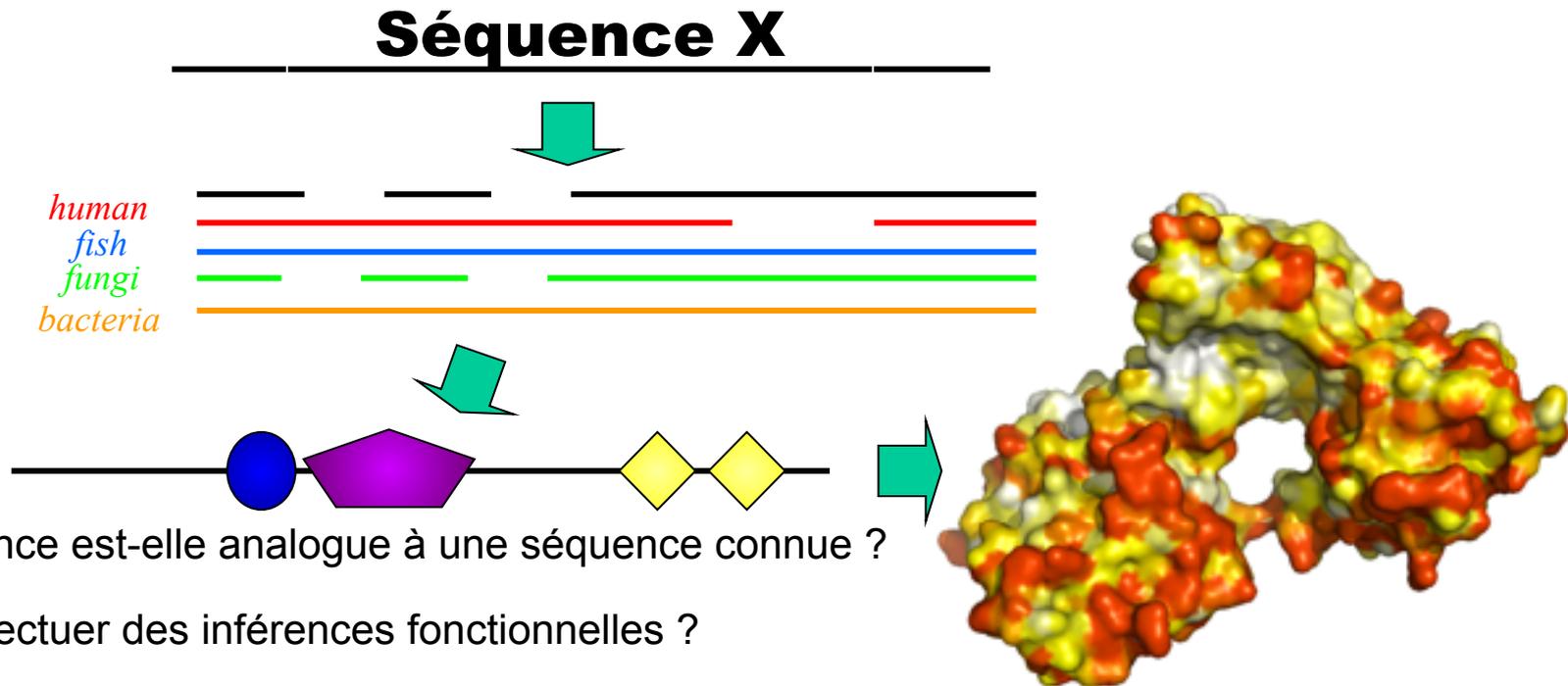
Profile-profile required

Ab initio



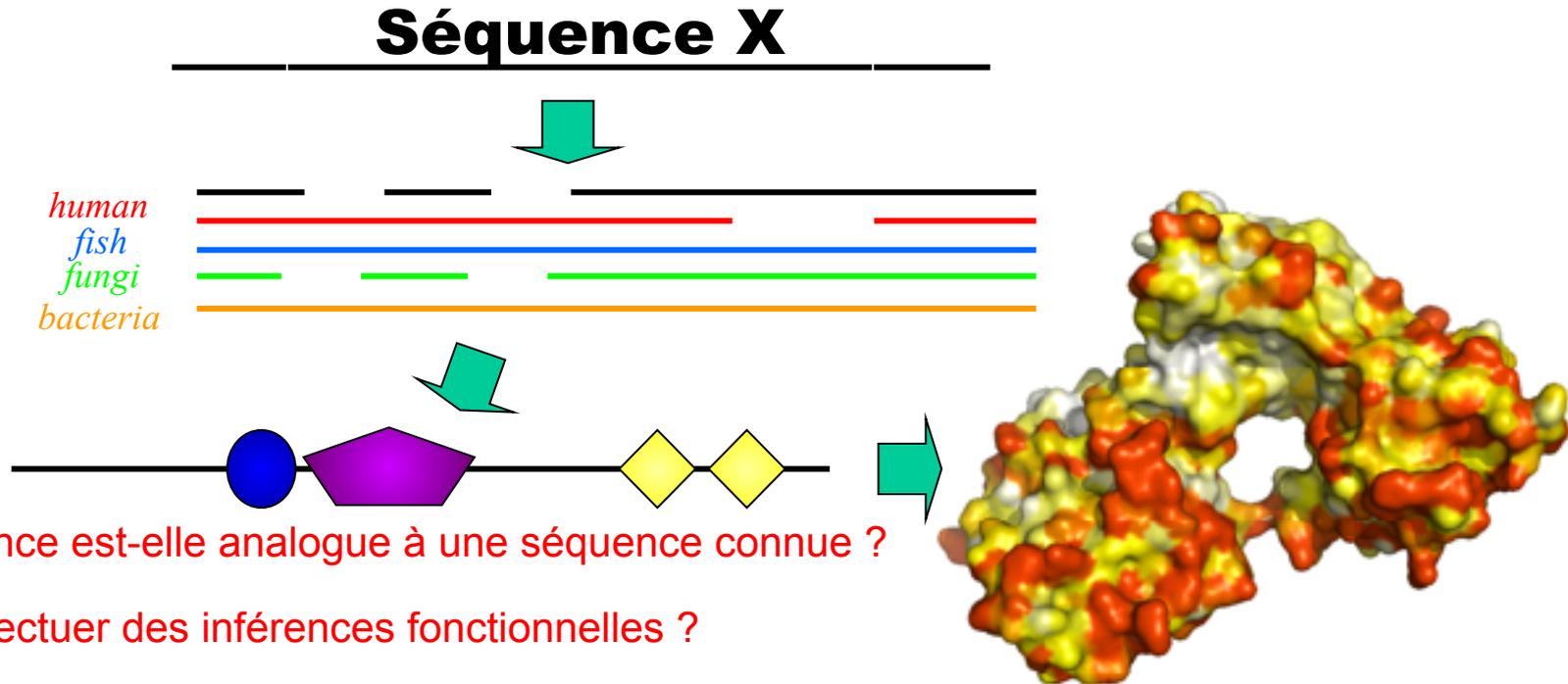
Baker & Sali, *Science* **294**, 2001, pp. 93-96

Quelles informations peut-on extraire d'une analyse de séquence ?



- 1 - Ma séquence est-elle analogue à une séquence connue ?
- 2 - Puis-je effectuer des inférences fonctionnelles ?
- 3 - Comment découper ma séquence en modules fonctionnels repliés ?
 - Quelles sont les régions fonctionnelles et structurales ?
- 4 - Peut-on prédire son organisation structurale ?
- 5 - Quelles sont les régions soumises à des pressions de sélection particulières ?

Quelles informations peut-on extraire d'une analyse de séquence ?



- 1 - Ma séquence est-elle analogue à une séquence connue ?
- 2 - Puis-je effectuer des inférences fonctionnelles ?
- 3 - Comment découper ma séquence en modules fonctionnels repliés?
 - Quelles sont les régions fonctionnelles et structurales ?
- 4 - Peut-on prédire son organisation structurale ?
- 5 - Quelles sont les régions soumises à des pressions de sélection particulières ?

Récupérer la séquence correcte de ma protéine

Type d'analyse	Approche	URL	Commentaires
Choix de la séquence	mots-clés, nom de protéine, fragment de séquence	http:// www.uniprot.org/	Sites généraux d'accès aux bases de données de protéomique / génomique,...
		http:// www.expasy.org/	
Question biologique	Gene Ontology (GO)	http:// geneontology.org/	Gene functions
	CAZY, Cazypedia	https:// www.cazypedia.org/ index.php/Main_Page	Bases de données spécifiques glycobiologie
Choix de l'organisme	Bibliographie	Pubmed, Isi web of science, Google Scholar,...	

Expasy

Categories

proteomics

protein sequences and identification

proteomics databases

function analysis

sequence data, analysis and motifs

protein modifications

protein structure

protein interactions

visualization of data

genomics

structure analysis

systems biology

evolutionary biology

population genetics

transcriptomics

metabolics

imaging

IT infrastructure

network connectivity

glycomics

Resources A-Z

Links/Documentation

Databases

-  **UniProtKB** - functional information on proteins - [\[more\]](#)
-  **UniProtKB/Swiss-Prot** - protein sequence database - [\[more\]](#)
-  **STRING** - protein-protein interactions - [\[more\]](#)
-  **SWISS-MODEL Repository** - protein structure homology models - [\[more\]](#)
-  **PROSITE** - protein domains and families - [\[more\]](#)
-  **ViralZone** - portal to viral UniProtKB entries - [\[more\]](#)
-  **neXtProt** - human proteins - [\[more\]](#)

-  **EMBLnet services** - bioinformatics tools, databases and courses - [\[more\]](#)
-  **ENZYME** - enzyme nomenclature - [\[more\]](#)
-  **GlycoCan** - International glycan structure repository - [\[more\]](#)
-  **GPDB** - gene and protease/enzyme - [\[more\]](#)
-  **HAMAP** - UniProtKB family classification and annotation - [\[more\]](#)
-  **MathGO** - protein-glycosaminoglycan interactions - [\[more\]](#)
-  **MetNetK** - Metabolic Network Repository & Analysis - [\[more\]](#)
-  **MSAPROtein** - MS/MS deconvolution - [\[more\]](#)
-  **MyHits** - protein domains database and tools - [\[more\]](#)
-  **PaADb** - protein abundance database - [\[more\]](#)
-  **Proline** - Popular science articles (in French) - [\[more\]](#)
-  **Protein Model Portal** - structural information for a protein - [\[more\]](#)
-  **Protein Spotlight** - Informally written reviews on proteins - [\[more\]](#)
-  **Reax** - expert curated resource of biochemical reactions - [\[more\]](#)
-  **SugarBind** - pathogen sugar-binding - [\[more\]](#)
-  **SWISS-3D-PAGE** - proteins on 2-D and SDS-PAGE maps - [\[more\]](#)
-  **SwissEcolcstere** - bioisomers for small molecules - [\[more\]](#)
-  **SwissLipids** - knowledge resource for lipid biology - [\[more\]](#)
-  **SwissPfam** - database of Pfam/Trifam events - [\[more\]](#)
-  **SwissEcolcstere** - non natural amino acid diastereois - [\[more\]](#)
-  **SwissVar** - variants in UniProtKB entries - [\[more\]](#)
-  **TCS** - interaction specificity in two-component systems - [\[more\]](#)
-  **UniCarbDB** - glycan mass and structural data - [\[more\]](#)
-  **UniCarbKB** - curated glycan database - [\[more\]](#)
-  **UniParc (UniProt sequence archive)** - UniProt sequence archive - [\[more\]](#)
-  **UniPathway** - metabolic pathways for the UniProtKB - [\[more\]](#)
-  **UniRef (UniProt sequence clusters)** - UniProtKB sequence clusters - [\[more\]](#)
-  **VenomZone** - portable venom proteins UniProtKB entries - [\[more\]](#)
-  **World-SDPAGE Consortium** - set of SDPAGE resources - [\[more\]](#)
-  **World-SDPAGE Repository** - gel-based proteomics data - [\[more\]](#)

Tools

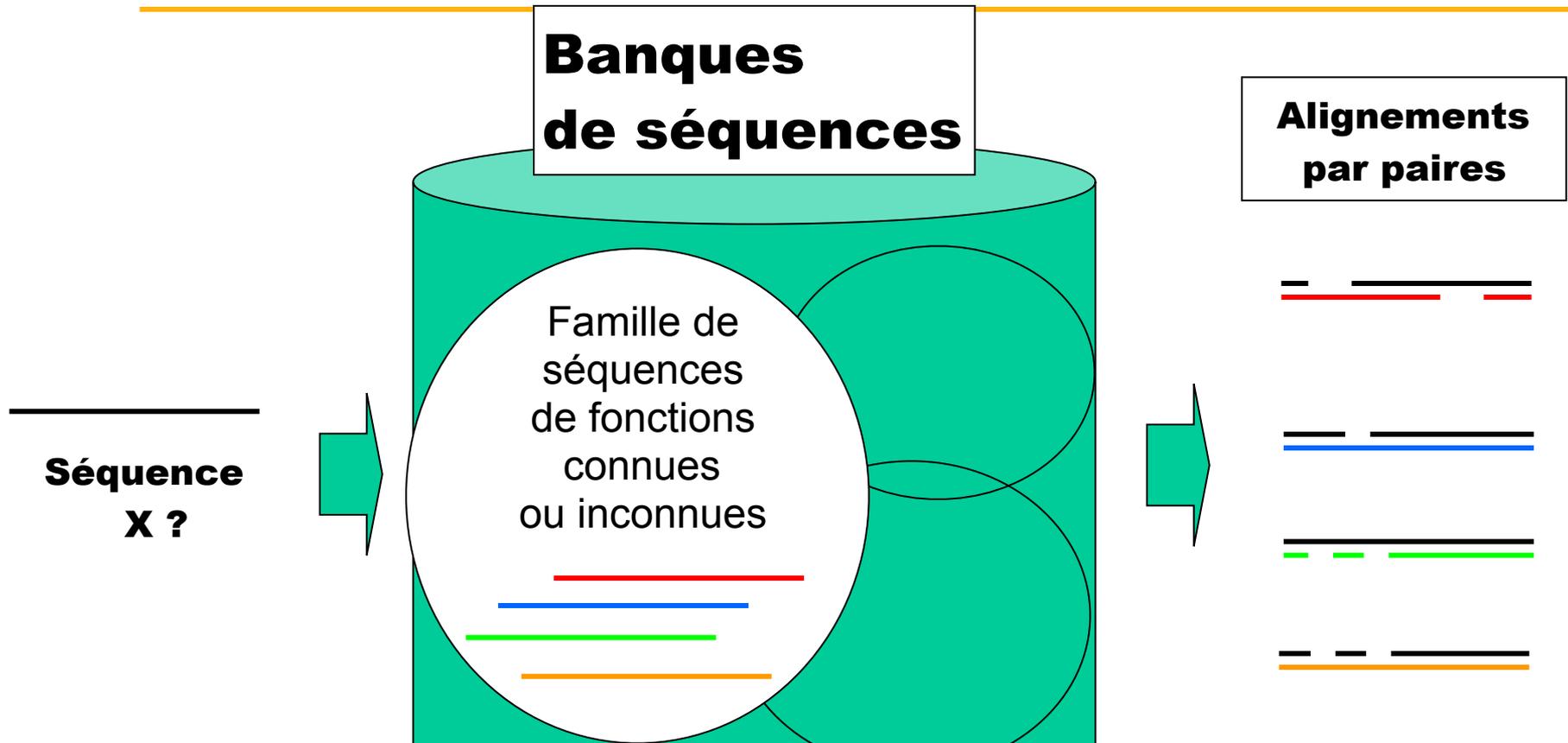
-  **SWISS-MODEL Workspace** - structure homology-modelling - [\[more\]](#)
-  **SwissDock** - protein ligand docking server - [\[more\]](#)
-  **ZIP** - Prediction of leucine zipper domains - [\[more\]](#)
-  **3dS** - Trifam-defined patterns in protein sequences - [\[more\]](#)
-  **AACompIdent** - protein identification by aa composition - [\[more\]](#)
-  **AACompSim** - amino acid composition comparison - [\[more\]](#)
-  **Agadir** - Prediction of the helical content of peptides - [\[more\]](#)
-  **ALF** - simulation of genome evolution - [\[more\]](#)
-  **Alignment tools** - Four tools for multiple alignments - [\[more\]](#)
-  **APESP** - Advanced Protein Secondary Structure Prediction - [\[more\]](#)
-  **Arasight** - Molecular modeling software - [\[more\]](#)
-  **bigPI** - predict GPI modification sites - [\[more\]](#)
-  **Biochemical Pathways** - Biochemical Pathways - [\[more\]](#)
-  **BLAST** - sequence similarity search - [\[more\]](#)
-  **BLAST (UniProt)** - BLAST search on the UniProt web site - [\[more\]](#)
-  **BLAST - NCBI** - Biological sequence similarity search - [\[more\]](#)
-  **BLAST - EMBL** - BLAST search on protein sequence databases - [\[more\]](#)
-  **blast2fasta** - Blast to Fasta conversion - [\[more\]](#)
-  **blastseq** - MSA pretty printer - [\[more\]](#)
-  **CPESP** - Protein secondary structure prediction - [\[more\]](#)
-  **ChroP** - chloroplast transit peptides & cleavage sites - [\[more\]](#)
-  **Click2Drug** - Directory of computational drug design tools - [\[more\]](#)
-  **ClustalC (UniProt)** - Align two or more protein sequences - [\[more\]](#)
-  **ClustalW** - Multiple sequence alignment - [\[more\]](#)
-  **ClustalW - MBL** - Multiple sequence alignment program - [\[more\]](#)
-  **ClustalW2** - Multiple sequence alignment program - [\[more\]](#)
-  **Coiled-Coils prediction** - Prediction of coiled coils regions - [\[more\]](#)
-  **CCILS** - Prediction of Coiled Coil Regions in Proteins - [\[more\]](#)
-  **ColorSeq** - Color Protein Sequence - [\[more\]](#)
-  **Compute pI/MW** - theoretical pI and MW computation - [\[more\]](#)
-  **CPIHomo** - Protein homology modeling - [\[more\]](#)
-  **CSS-Pfam** - Prediction of palmitoylation sites in proteins - [\[more\]](#)
-  **CAS-TMHier** - Prediction of transmembrane regions - [\[more\]](#)
-  **Decrease redundancy** - Sequence redundancy reduction - [\[more\]](#)
-  **DALIGN** - Local multiple sequence alignment - [\[more\]](#)
-  **DistoGlyc** - GlyNAc O-glycosylation sites in Databases - [\[more\]](#)
-  **DisEMBL** - Prediction of disordered protein regions - [\[more\]](#)
-  **DUP-SVM** - Domain linker predictor - [\[more\]](#)

Récupérer la séquence correcte de ma protéine

Type d'analyse	Approche	URL	Commentaires
Choix de la séquence	mots-clés, nom de protéine, fragment de séquence	http:// www.uniprot.org/	Sites généraux d'accès aux bases de données de protéomique / génomique,...
		http:// www.expasy.org/	
Question biologique	Gene Ontology (GO)	http:// geneontology.org/	Gene functions
	CAZY, Cazypedia	https:// www.cazypedia.org/ index.php/Main_Page	Bases de données spécifiques glycobiologie
Choix de l'organisme	Bibliographie	Pubmed. Isi web of science. Google Scholar...	

Recherche d'homologues

Comparaison séquence/séquence



OUTILS : Blast (<http://www.ncbi.nlm.nih.gov/BLAST/>)

Basic Local Alignment Search Tool

Méthodes actuelles de prédiction de structure 3D

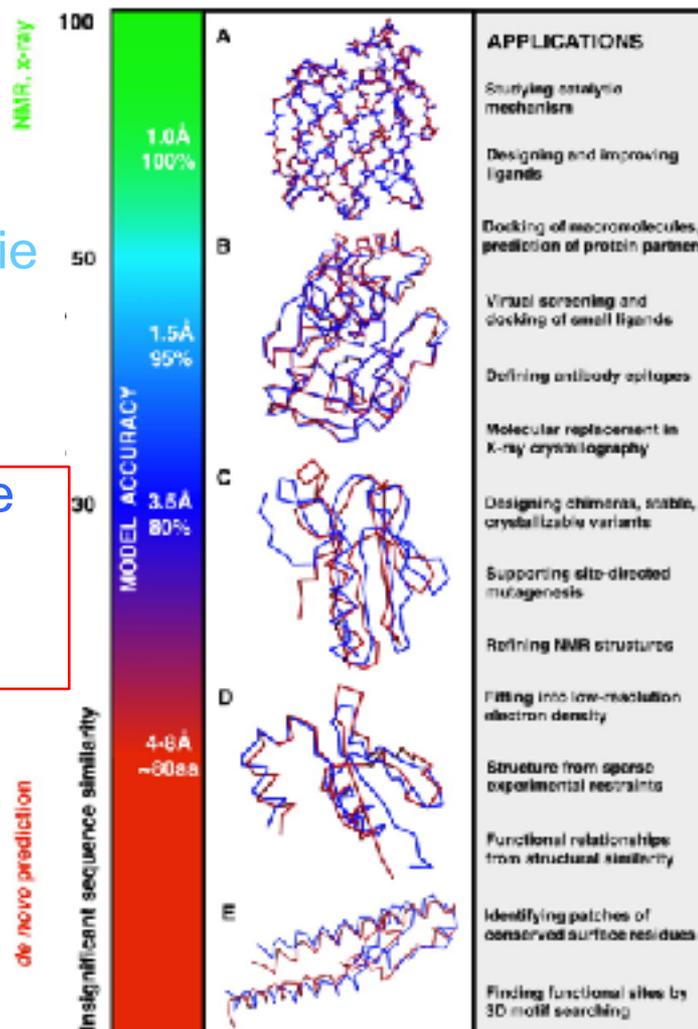
Expérimentales

Modélisation par homologie
ou comparative
(blast suffisant)

Beginning of twilight zone

Profile-profile required

Ab initio

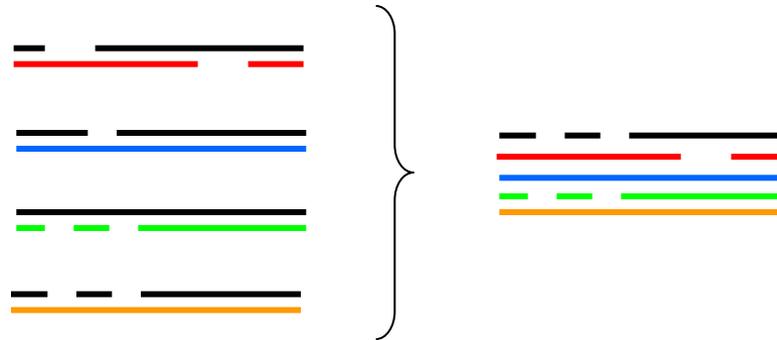


Baker & Sali, *Science* **294**, 2001, pp. 93-96

Combiner les alignements entre paires → alignements multiples

Importance de **l'histoire évolutive** d'une famille de protéines
“pairwise alignment whispers... multiple alignment shouts out loud.” (Hubbard et al., 1996)

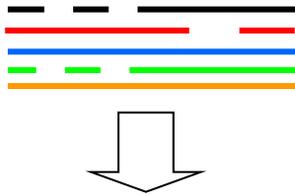
Alignements par paires



OUTILS : * Clustalw (<http://www.ebi.ac.uk/Tools/clustalw/index.html/>)
*** MUSCLE (<http://www.drive5.com/muscle/>)
*** MAFFT (<http://mafft.cbrc.jp/alignment/server/>)

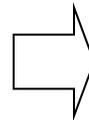
Extraction des informations des alignements multiples → Notion de profil

Détection des séquences reliées à la séquence cible.



Alignement Multiple

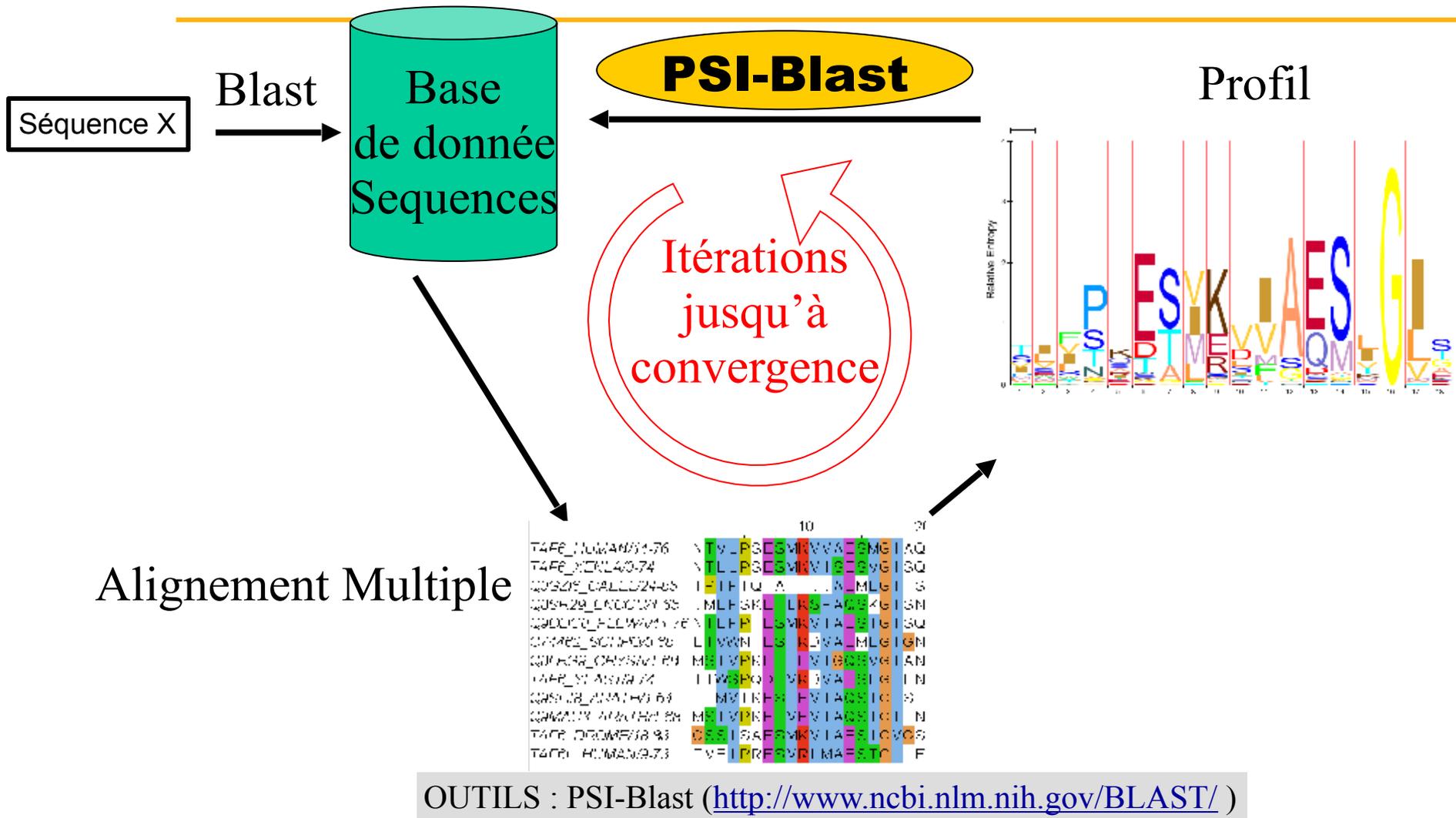
	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15	16	17	18	19	20		
T4FE_HUMAN/1-76	N	T	L	P	S	G	M	V	V	A	E	S	M	G	I	A	Q					
T4FE_NORL40-74	N	T	L	P	S	G	M	V	V	A	E	S	M	G	I	A	Q					
Q09205_DALLD24-65	L	F	T	T	Q	A																
Q09429_DROD12Y-55	L	M	L	F	S	K	L	L	K	S	-	A	G	S	K	G	I	G	N			
Q9D010_FELMVMY-76	N	T	L	P	S	G	M	V	V	A	E	S	M	G	I	A	Q					
Q14482_SGTHFV0-90	L	T	V	W	N	L	S	K	L	V	A	M	L	G	I	G	N					
Q14482_CHYSAV1-84	M	S	L	V	P	K																
T4FE_YEAST/1-74	L	T	V	W	N	L	S	K	L	V	A	M	L	G	I	G	N					
Q09429_DROD12Y-55	L	M	L	F	S	K	L	L	K	S	-	A	G	S	K	G	I	G	N			
Q14482_SGTHFV0-90	L	T	V	W	N	L	S	K	L	V	A	M	L	G	I	G	N					
Q14482_CHYSAV1-84	M	S	L	V	P	K																
T4FE_DROME/1-83	C	S	S	L	S	A	F	R	M	K	V	I	A	F	S	L	C	V	G			
T4FE_HUMAN/9-73	T	V	E	L	P	R	F	S	M	V	I	M	A	F	S	T	C	F				



Profil ou PSSM (position specific scoring matrix)



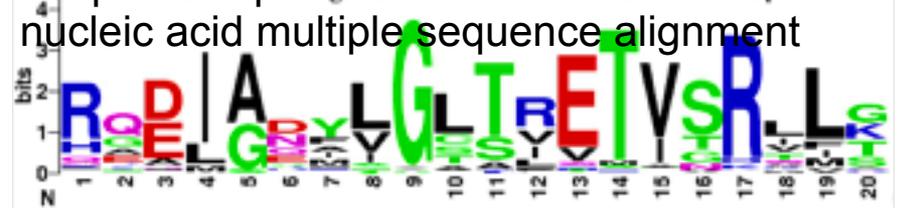
PSI-BLAST : méthode itérative pour construire les profils



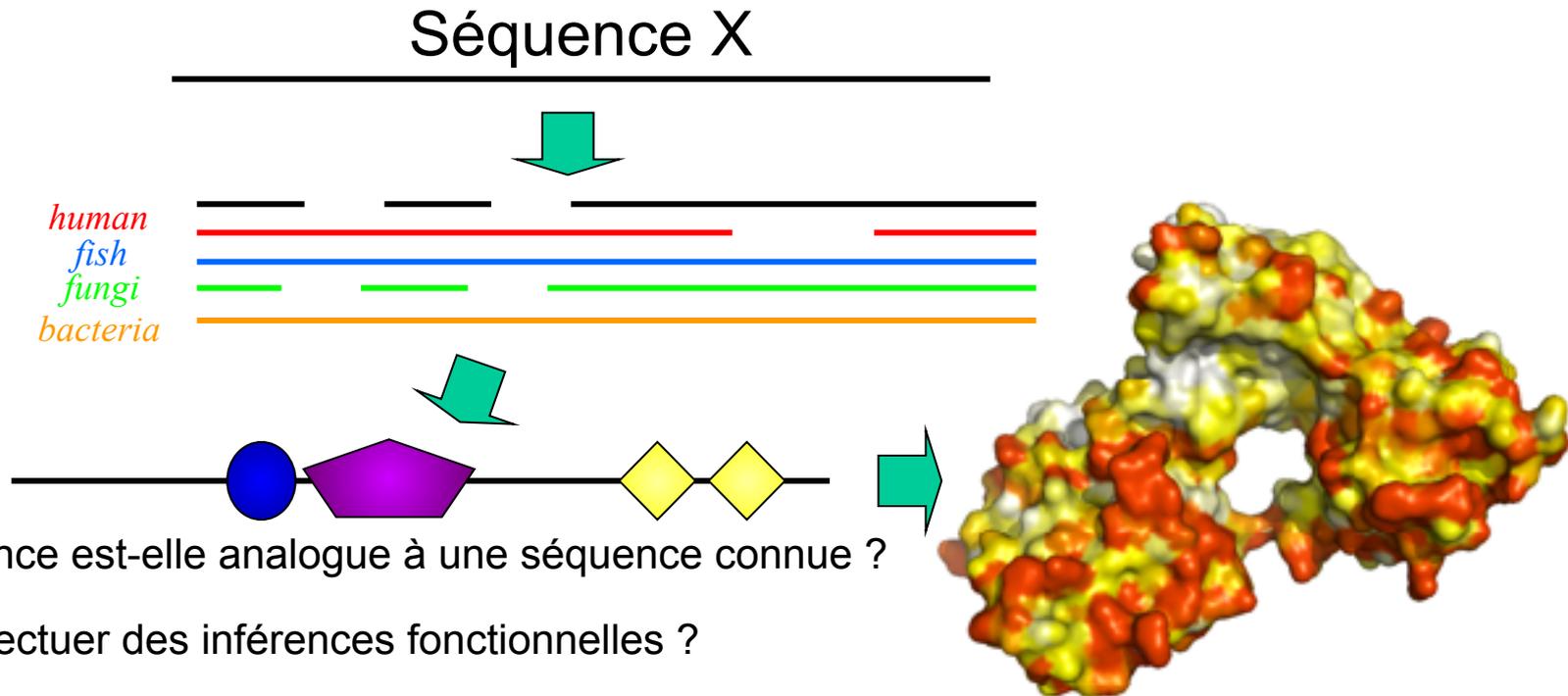
OUTILS : PSI-Blast (<http://www.ncbi.nlm.nih.gov/BLAST/>)

Visualisation des alignements multiples

Logiciel	URL	
Jalview	http://www.jalview.org/	Logiciel pour l'édition et la visualisation d'alignement de séquence, analyse phylogénique, analyse structurale
Geneious	http://www.geneious.com/	Logiciel pour l'édition et la visualisation d'alignement de séquence, analyse phylogénique, analyse structurale (payant mais free trial)
Alscript	http://www.compbio.dundee.ac.uk/software.html	Format multiple sequence alignments in PostScript for publications
ESPrict	http://esprict.ibcp.fr/ESPrict/ESPrict/	Easy Sequencing in PostScript': met en valeur les similarités et l'information de structures secondaires à partir d'alignements de séquences pour l'analyse et la publication
WebLogo	http://weblogo.berkeley.edu/logo.cgi	Graphical representation of an amino acid or nucleic acid multiple sequence alignment



Quelles informations peut-on extraire d'une analyse de séquence ?



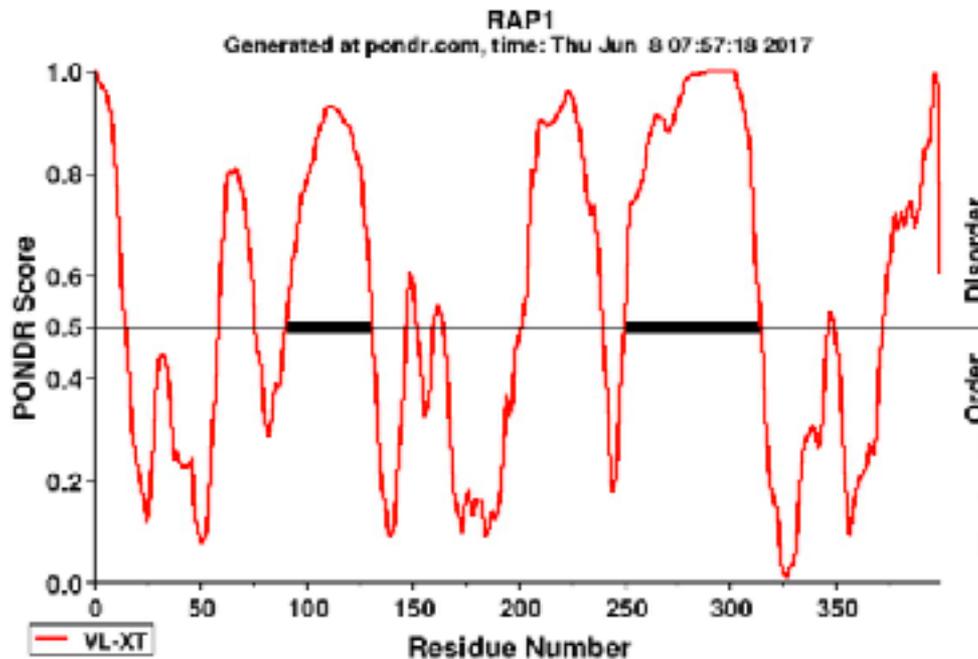
- 1 - Ma séquence est-elle analogue à une séquence connue ?
- 2 - Puis-je effectuer des inférences fonctionnelles ?
- 3 - Comment découper ma séquence en modules fonctionnels repliés?
 - Quelles sont les régions fonctionnelles et structurales ?
- 4 - Peut-on prédire son organisation structurale ?
- 5 - Quelles sont les régions soumises à des pressions de sélection particulières ?

Recherche des régions ordonnées / désordonnées

PONDR[®]

Predictor of Natural Disordered Regions :

<http://www.pondr.com/>



Estimation des régions
désordonnées ou
ordonnées

Molecular Kinetics
(www.molecularkinetics.com;
main@molecularkinetics.com)

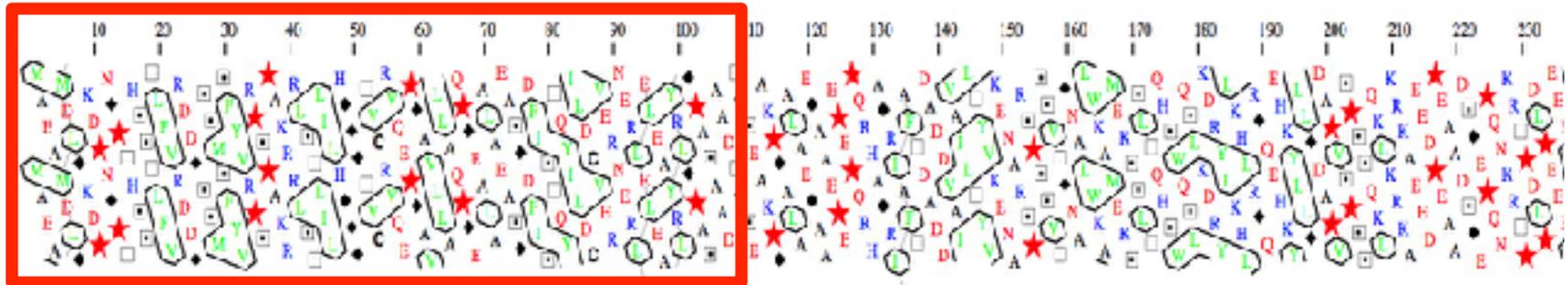
PONDNRs are typically feedforward neural networks that use sequence attributes taken over windows of 9 to 21 amino acids. These attributes, such as the fractional composition of particular amino acids, hydrophathy, or sequence complexity, are averaged over these windows and the values are used to train the neural network during predictor construction.

Recherche des régions ordonnées / désordonnées

HCA 1.0.2

<http://mobyli.rpbs.univ-paris-diderot.fr/cgi-bin/portal.py#forms::HCA>

Hydrophobic Cluster Analysis.



Domaine replié

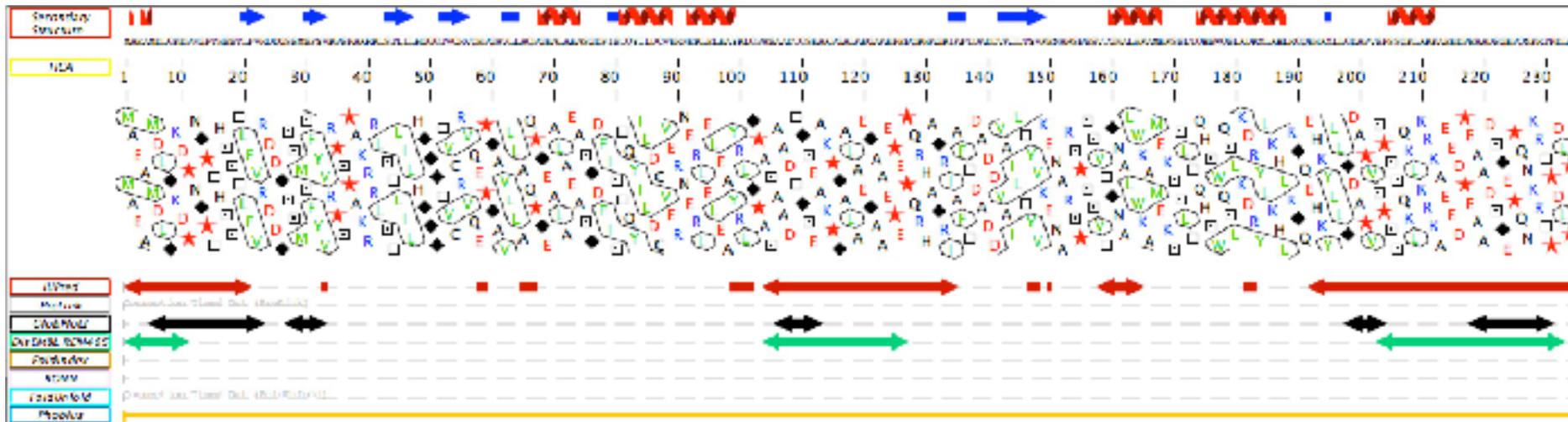
Deciphering protein sequence information through hydrophobic cluster analysis (HCA): current status and perspectives. Cell Mol Life Sci. 1997 Aug;53(8):621-45. Review.

Callebaut I, Labesse G, Durand P, Poupon A, Canard L, Chomilier J, Henrissat B, Mornon JP.

Recherche des régions ordonnées / désordonnées



MeDor (Metaserver of Disorder) : <http://www.vazymolo.org/MeDor/index.html>



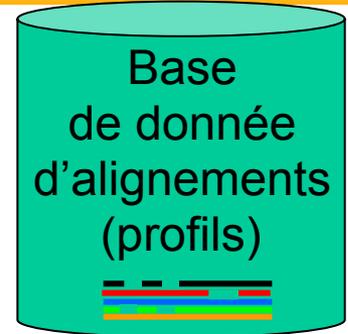
"MeDor: a metaserver for predicting protein disorder." Philippe Lieutaud, Bruno Canard and Sonia Longhi. BMC Genomics. 2008 Sep 16;9 Suppl 2:S25.

Exploitation des bibliothèques d'alignements multiples

→ Identification rapide de la composition en domaines

Séquence X

Existe-t-il des régions conservées ?



<http://smart.embl-heidelberg.de/>



Identification de domaines

Recherche par architecture analogue

Interface avec les cartes interactions



<http://pfam.xfam.org/>



CDD, conserved domain database au NCBI

Exemple hRAP1

>gi|8102033|gb|AAF72711.1|AF262988_1 TRF2-interacting telomeric RAP1 protein [Homo sapiens]

MAEAMDLGKDPNGPTHSSTLFVRDDGSSMSFYVRPSPAKRRLSTLILHGGGTVCRVQEPGAVLLAQPGEA
LAEASGDFISTQYILDCVERNERLELEAYRLGPASAADTGSEAKPGALAEGAAEPEPQRHAGRIAFTDAD
DVAILITYVKENARSPSSVTGNALWKAMEKSSLTQHSWQSLKDRYLKHLRGQEHKYLLGDAPVSPSSQKLK
RKAEDPEAADSGEPQNKRTPDLPREEYVKEEQENEEAVKKMLVEATREFEEVVVDESPPDFEIHITMC
DDDPPTPEEDSETQPDEEEEEEEEEKVSQPEVGAAIKIIRQLMEKFNLDLSTVTQAFLKNSGELEATSAFL
ASGQRADGYPIWSRQDDIDLQKDEDTREALVKKFGAQNVARRIEFRKK

Exemple hRAP1

1/ Uniprot

Family & Domains¹

Domains and Repeats

Feature key	Position(s)	Description	Actions	Graphical view	Length
Domain ¹	78 - 101	BRCT	Add BLAST		24
Domain ¹	128 - 188	Myb-like	Add BLAST		61

Motif

Feature key	Position(s)	Description	Actions	Graphical view	Length
Motif ¹	383 - 399	Nuclear localization signal Sequence analysis	Add BLAST		17

Compositional bias

Feature key	Position(s)	Description	Actions	Graphical view	Length
Compositional bias ¹	214 - 304	Asp/Glu-rich (addic)	Add BLAST		91

>gi|8102033|gb|AAF72711.1|AF262988_1 TRF2-interacting telomeric RAP1 protein [Homo sapiens]

MAEAMD LGKDPNGP THSSTLFVRDDGSSMSFYVRPSPAKRRLSTLILHGGGTVCRVQEPGAVLLAQPGEA
LAEASGDFISTQYILDCVERNERLELEAYRLGPASAADTGSEAKPGALAE GAAEPEPQRHAGRIAF TDAD
DVAILTYVKENARSPSSVTGNALWKAMEKSSLTQHSWQSLKDRYLKHLRGQEHKYL LGDAPVSPSSQKLK
RKA EEDPEAADSGEPQNKRT PDLPEEEYVKEEIQENEEAVKKMLVEATREFEEVVVDESPPDFEIHITMC
DDDPPTPEEDSETQPDEEEEEEEEEKVSQPEVGAAIKIIRQLMEKFNLDLSTVTQAF LKNSGELEATS AFL
ASGQRADGYPIWSRQDDIDLQKDD E D TREALVKKFGAQNVARRIEF RKK

Exemple hRAP1

1b/ Expasy => protparam

Number of amino acids: 399

Molecular weight: 44259.9

Theoretical pI: 4.64

Amino acid composition:

Ala (A)	40	10.0%
Arg (R)	23	5.8%
Asn (N)	9	2.3%
Asp (D)	30	7.5%
Cys (C)	3	0.8%
Gln (Q)	18	4.5%
Glu (E)	51	12.8%
Gly (G)	23	5.8%
His (H)	7	1.8%
Ile (I)	14	3.5%
Leu (L)	32	8.0%
Lys (K)	25	6.3%
Met (M)	7	1.8%
Phe (F)	11	2.8%
Pro (P)	25	6.3%
Ser (S)	30	7.5%
Thr (T)	19	4.8%
Trp (W)	3	0.8%
Tyr (Y)	8	2.0%
Val (V)	21	5.3%
Pyl (O)	0	0.0%
Sec (U)	0	0.0%

Total number of negatively charged residues (Asp + Glu): 81

Total number of positively charged residues (Arg + Lys): 48

Atomic composition:

Carbon	C	1917
Hydrogen	H	3033
Nitrogen	N	537
Oxygen	O	646
Sulfur	S	10

Formula: C₁₉₁₇H₃₀₃₃N₅₃₇O₆₄₆S₁₀

Total number of atoms: 6143

Extinction coefficients:

Extinction coefficients are in units of M⁻¹ cm⁻¹, at 280 nm measured in water.

Ext. coefficient 28545

Abs 0.1% (=1 g/l) 0.645, assuming all pairs of Cys residues form cystines

Ext. coefficient 28420

Abs 0.1% (=1 g/l) 0.642, assuming all Cys residues are reduced

Estimated half-life:

The N-terminal of the sequence considered is M (Met).

The estimated half-life is: 30 hours (mammalian reticulocytes, in vitro).

>20 hours (yeast, in vivo).

>10 hours (Escherichia coli, in vivo).

Instability index:

The instability index (II) is computed to be 57.42

This classifies the protein as unstable.

Aliphatic index: 70.25

Grand average of hydropathicity (GRAVY): -0.763

Exemple hRAP1

```
>gi|8102033|gb|AAF72711.1|AF262988_1 TRF2-interacting telomeric RAP1
protein [Homo sapiens]
MAEAMD LGKDPNGP THSSTL FVRDDGSSMSFYVRPSPAKRRLSTLILHGGGTVCRVQEPGAVLLAQPGEA
LAEASGDFISTQYILDCVERNERLELEAYRLGPASAADTGSEAKPGALAE GAAEPEPQRHAGRIAFTDAD
DVAILTYVKENARSPSSVTGNALWKAMEKSSLTQHSWQSLKDRYLKHLRGQEHKYLLGDAPVSPSSQKLK
RKA EEDPEAADSGEPQNKRT PDLPEEEYVKEEIQENEEAVKKMLVEATREFEEVVVDESPPDFEIHITMC
DDDPPTPEEDSETQPDEEEEEEEEEKVSQPEVGAAIKIIRQLMEKFNLDLSTVTQAFLKNSGELEATSAFL
ASGQRADGYPIWSRQDDIDLQKDD EDTREALVKKF GAQNVARRIEFRKK
```

En première approche => Découpage en 3 domaines

A confirmer par une protéolyse ménagée

Méthodes actuelles de prédiction de structure 3D

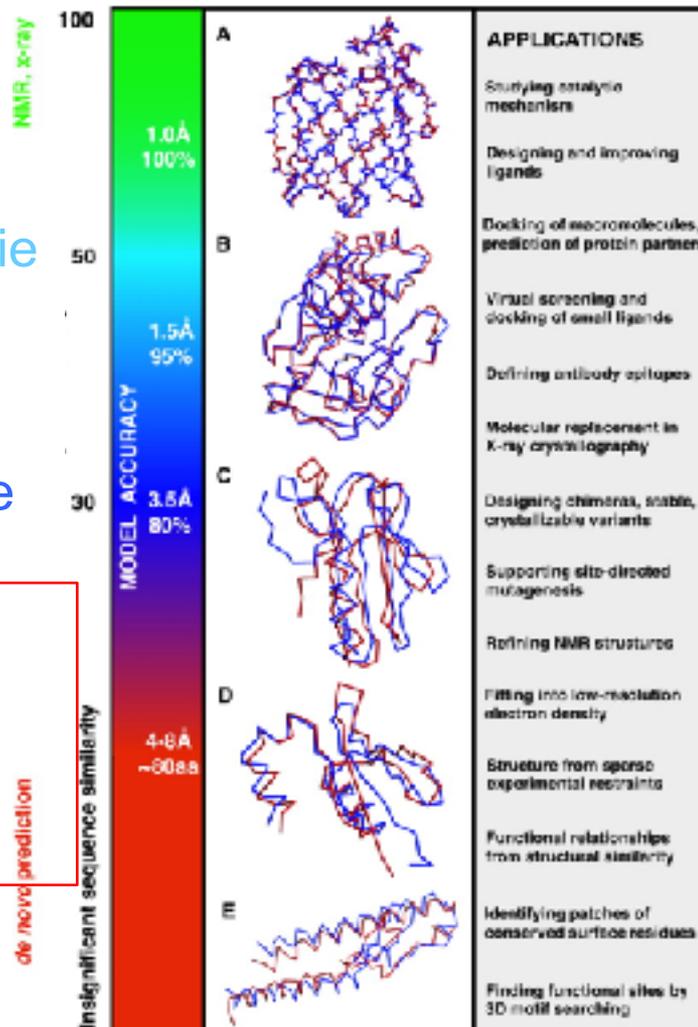
Expérimentales

Modélisation par homologie
ou comparative
(blast suffisant)

Beginning of twilight zone

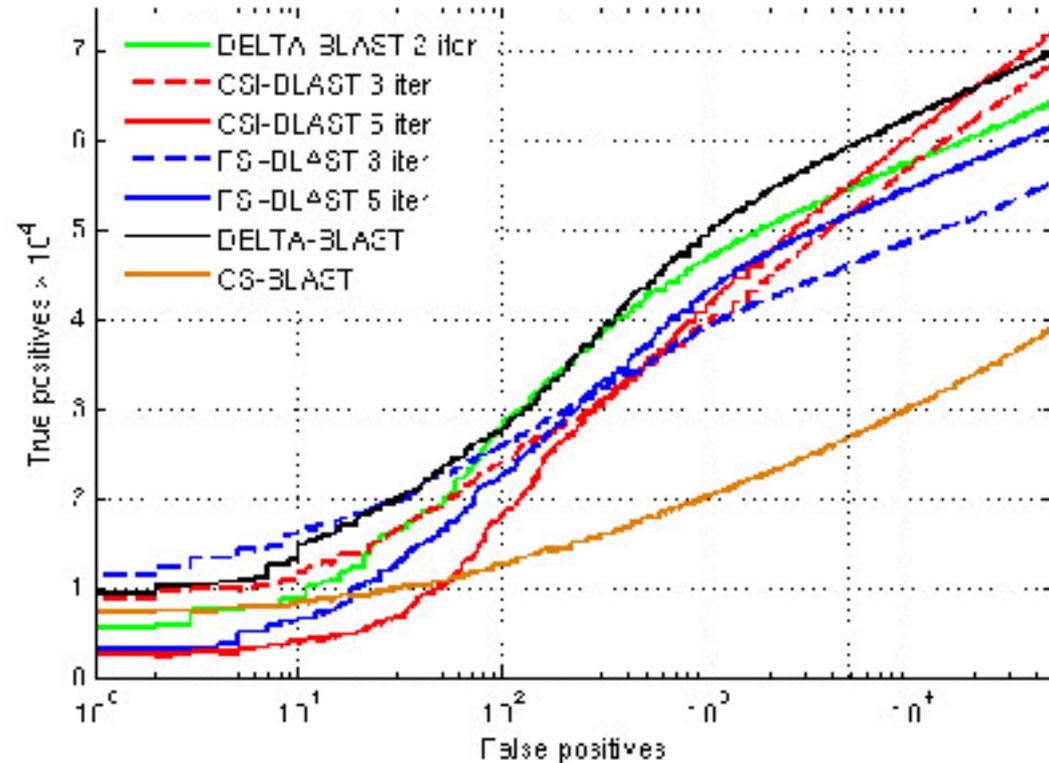
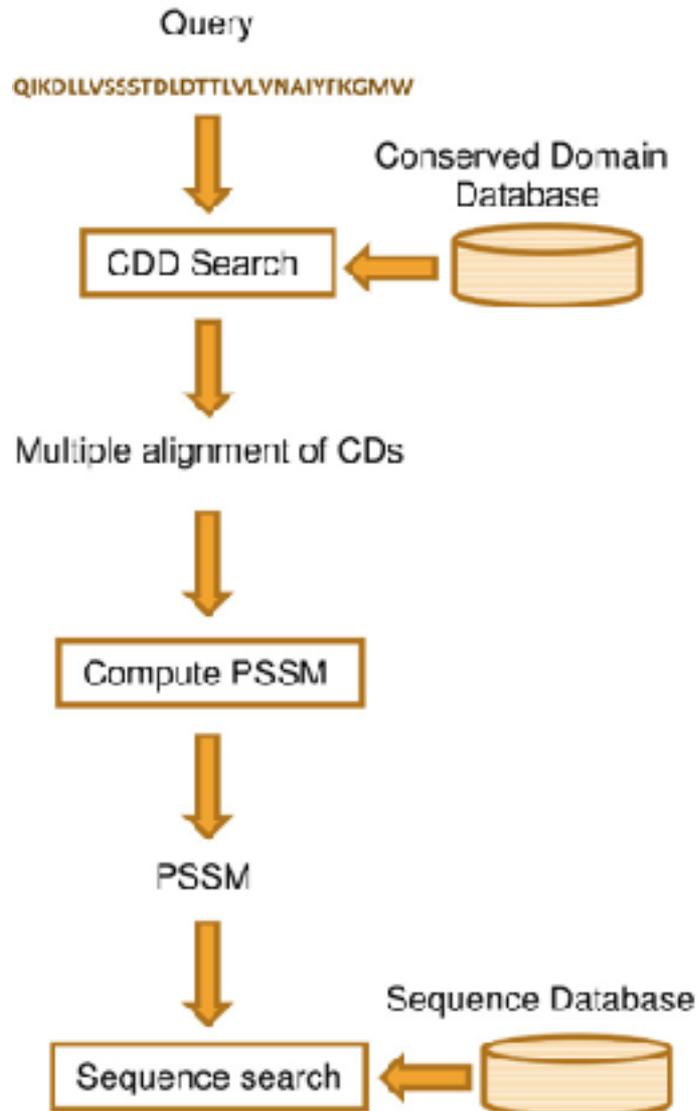
Profile-profile required

Ab initio



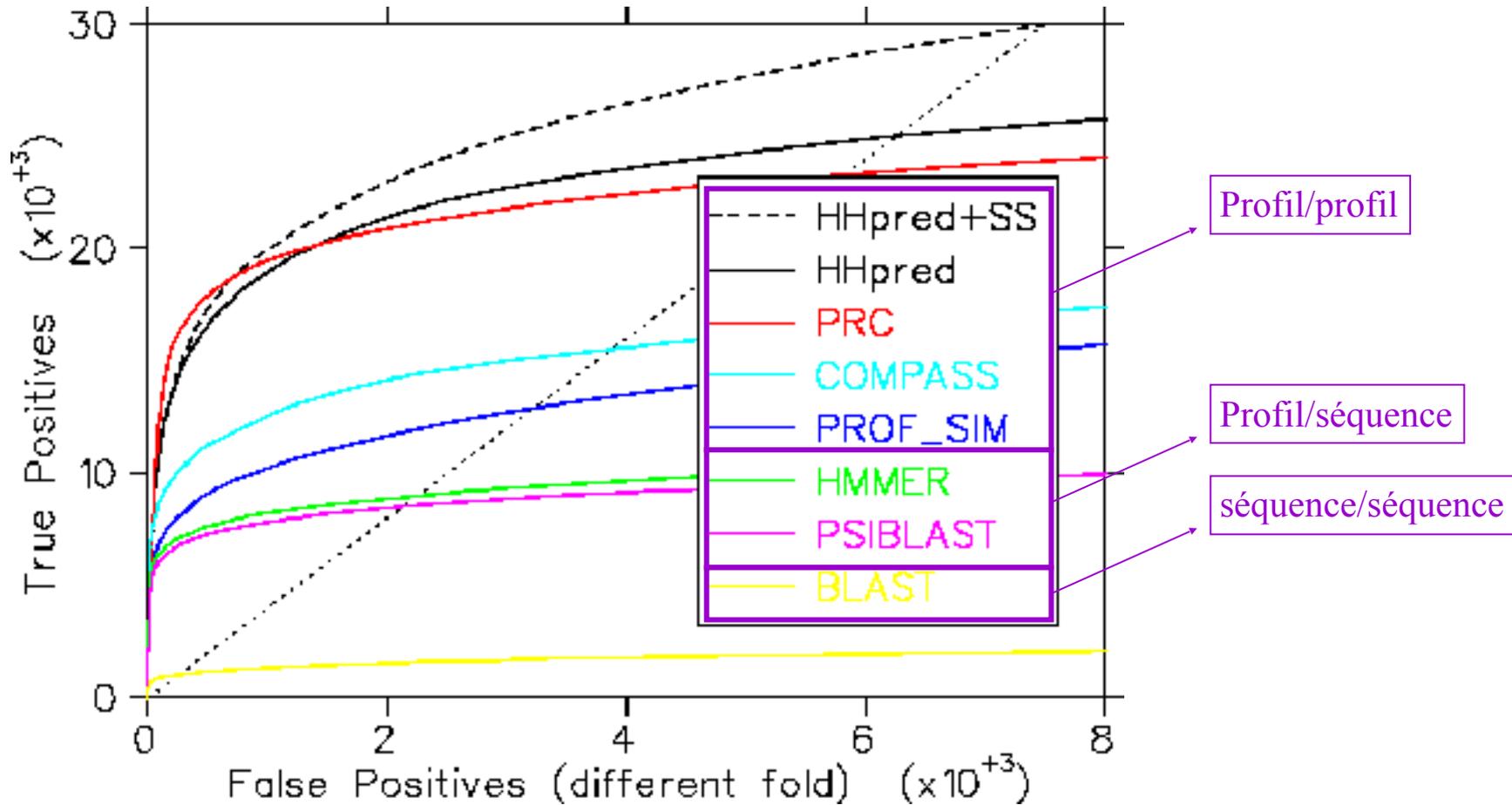
Baker & Sali, *Science* **294**, 2001, pp. 93-96

Delta-Blast : Exploitation des domaines pour l'amélioration des recherches d'homologies lointaines



Performances comparées

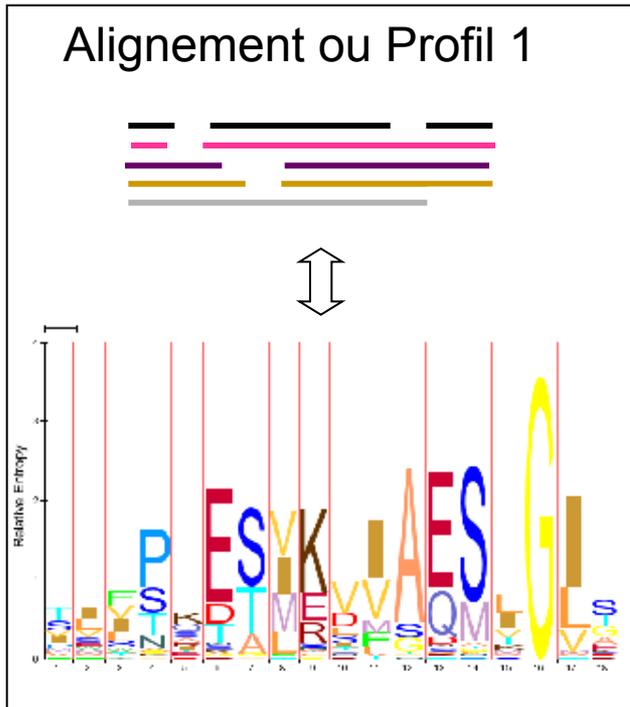
Gain très important en spécificité
grâce aux comparaisons profil/profil



J Söding. Bioinformatics 2005. 21(7):951-60

HHpred+SS : méthode comparaison profils HMM/HMM

Les développements méthodologiques pour l'analyse de séquence à très haute divergence

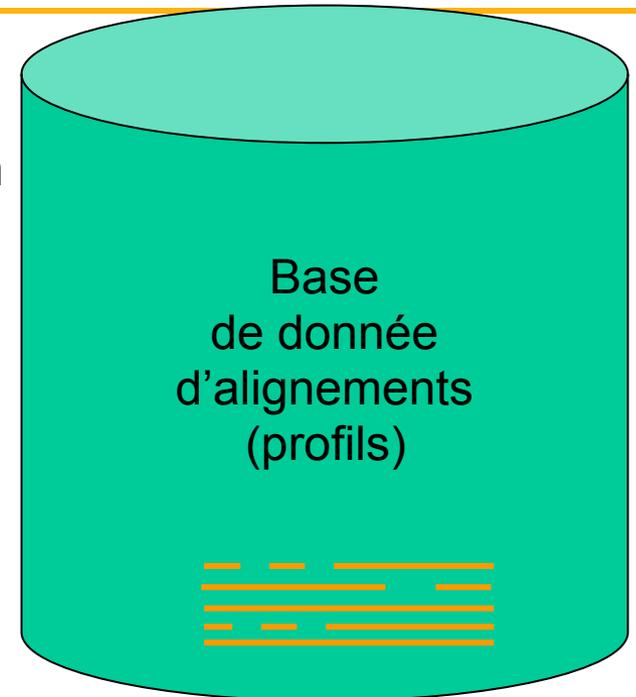


Outils de comparaison
jusqu'à 2005



ex : PSI-Blast

?



Exemples :

HHpred

Soding J, Bioinformatics, 2005

<http://toolkit.tuebingen.mpg.de/hhpred>

Très bon serveur pour l'analyse de séquences

<http://toolkit.tuebingen.mpg.de>

Soding J, Biegert A, Lupas AN. *Nucleic Acids Res.* (2005) 33(Web Server issue):W244-8

The screenshot displays the Bioinformatics Toolkit interface. On the left is a sidebar with a 'HOME' link, a logo, and a 'Show results of job:' section with a 'Show results' button. Below this is a 'Recent jobs:' section with 'Reset all' and 'Delete all' buttons, and a 'Clear job' button. A 'Bioinformatics Jobs' section contains a list of jobs with colored status indicators (blue, orange, green, red).

The main content area features the 'Bioinformatics Toolkit' header and 'Max Planck Institute for Developmental Biology' logo. A navigation bar includes 'Search', 'Alignment', 'Sequence Analysis', 'Zany Structures', 'Dary Structures', 'Classification', and 'Utils'. A menu of tools is visible, including 'CD-BLAST', 'HHpred', 'HHblits', 'HHsearch', 'HHMER2', 'PatternSearch', 'ProtBLAST', 'ProtELAST', 'PSI-BLAST', 'PSI-ELAST', and 'SimDist00'. A 'New beta HHsuite 3.0 available with better performance and new functionality!' announcement is present with a 'Download' link.

The 'HHpred - Homology detection & structure prediction by HMM-HMM comparison' tool is selected, with a 'Help' button. The 'Input' section contains a text area for 'Paste protein sequence or MSA', a 'Search with pancreaticoduodenal protein PIP4D.' button, and a 'Choose protein folder' button. A 'Select input format' dropdown is set to 'FASTA'. 'Reset form' and 'Submit job' buttons are at the bottom right of the input section.

The 'Search Options' section includes 'Select HMM databases (hold Ctrl to select several)', 'MSA Generation Method' (radio buttons for 'HMMs' and 'Ahoias'), 'Max. NSA Generation iterations' (dropdown), 'Scope secondary structure' (radio buttons for 'yes', 'no', and 'predicted vs predicted only'), 'Alignment mode' (radio buttons for 'local' and 'global'), and 'Realign with HAC' (checkbox).

Dissection d'un output HHpred



Switch : vue par profil

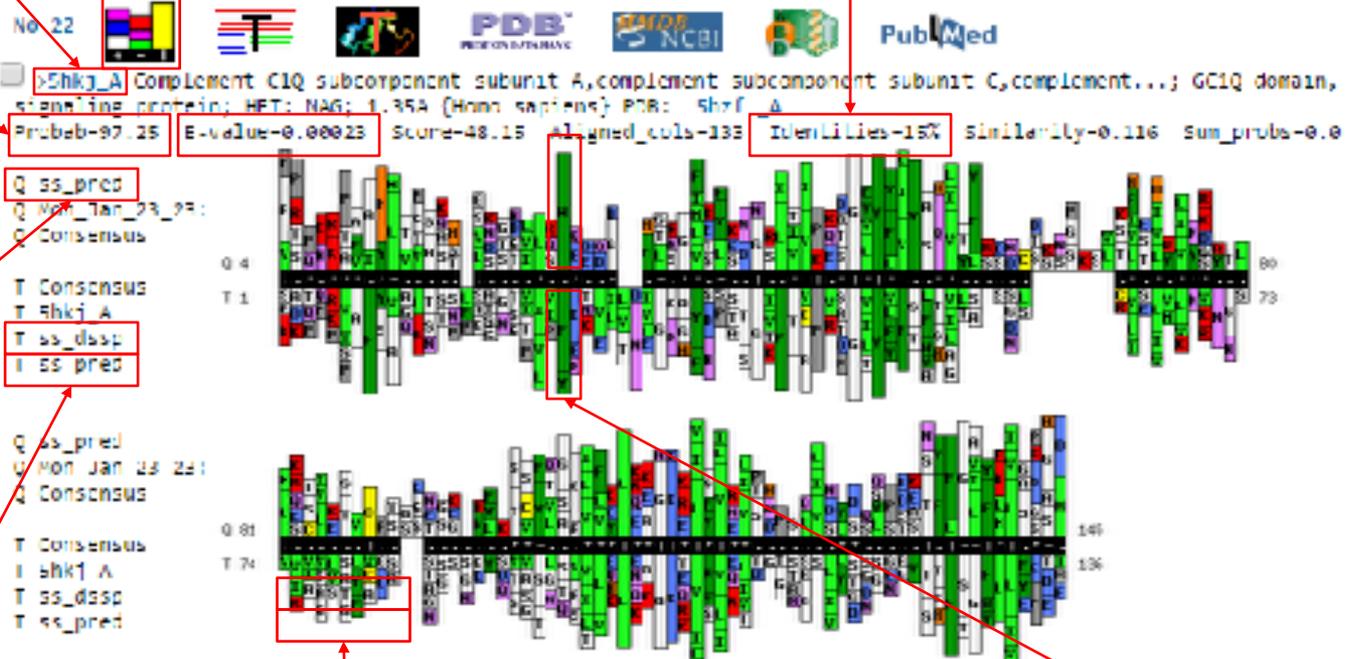
Match pdb

%identity Query vs Target

% Probabilité
fiabilité du match
Sécurité → P>90%
Possible → P>50%

Q=Query
ss_pred=prediction
struct. secondaire

T=Target
ss_pred=prediction
struct. secondaire
Ss_dssp=struct. sec.
observees dans pdb



E=Feuillet beta (e=moins fiable)
H=Hélices
C=désordonné

Majuscule > 50% conservé dans la famille
Minuscule = modérément conservé

Prédiction de Structures Secondaires

Quelques serveurs qui regroupent un ensemble d'outils de prédiction

Jpred

<http://www.compbio.dundee.ac.uk/jpred/>

The PredictProtein Server

<https://www.predictprotein.org/>

PBIL

<https://prabi.ibcp.fr/htm/index.php>

Considérations pratiques

- La modélisation par homologie permet d'obtenir des **modèles** de très bonne qualité
 - 1 Å de précision sur la position des atomes correspond à:
 - Structure *X-ray* à 2.5 Å de résolution et un facteur *R* de 25%
 - *RMN* avec 10 contraintes de distances inter-protons par résidus
- La modélisation par homologie réalisée avec des structures dont les identités de séquence sont >40% identité donne des résultats équivalents.

MAIS

- La qualité des modèles est dépendante de l'alignement de séquence
 - *importance d'un ajustement manuel éventuel de cet alignement*
- L'emploi de templates multiples améliore grandement la qualité des alignements, donc des modèles obtenus.

Exemple hRAP1

>gi|8102033|gb|AAF72711.1|AF262988_1 TRF2-interacting telomeric RAP1 protein [Homo sapiens]

MAEAMDLGKDPNGPHTHSSTLFVRRDDGSSMSFYVVRPSPAKRRLSTLILHGGGTVCRVQEPGAVLLAQPGEA
LAEASGDFISTQYILDCVERNERLELEAYRLGPASAADTGSEAKPGALAEGAAEPEPQRHAGRIAFTDAD
DVAILTYVKENARSPSSVTGNALWKAMEKSSLTQHSWQSLKDRYLKHLRGQEHKYLLGDAPVSPSSQKLK
RKAEDPEAADSGEPQNKRTPLPEEEYVKEEIQENEEAVKKMLVEATREFEEVVVDESPPDFEIHITMC
DDDPPTPEEDSETQPDEEEEEEEEEKVSQPEVGAAIKIIRQLMEKFNLDLSTVTOAFLKNSGELEATS AFL
ASGQRADGYPIWSRQDDIDLQKDDDETREALVKKFGAQNVARRIEFRKK

MAEAMDLGKDPNGPHTHSSTLFVRRDDGSSMSFYVVRPSPAKRRLSTLILHGGGTVCRVQEPGAVLLAQPGEA
LAEASGDFISTQYILDCVERNERLELEAYRLGPASAADTGS

GSEAKPGALAEGAAEPEPQRHAGRIAFTDAD
DVAILTYVKENARSPSSVTGNALWKAMEKSSLTQHSWQSLKDRYLKHLRGQEHKYLLGDAPV

PDEEEEEEEEEKVSQPEVGAAIKIIRQLMEKFNLDLSTVTOAFLKNSGELEATS AFL
ASGQRADGYPIWSRQDDIDLQKDDDETREALVKKFGAQNVARRIEFRKK

Exemple hRAP1

MAEAMDLGKDPNGP~~THSST~~ : zone d'incertitude

=> ajustement/amélioration des alignements

=> clonages multiples

=> protéolyse ménagée

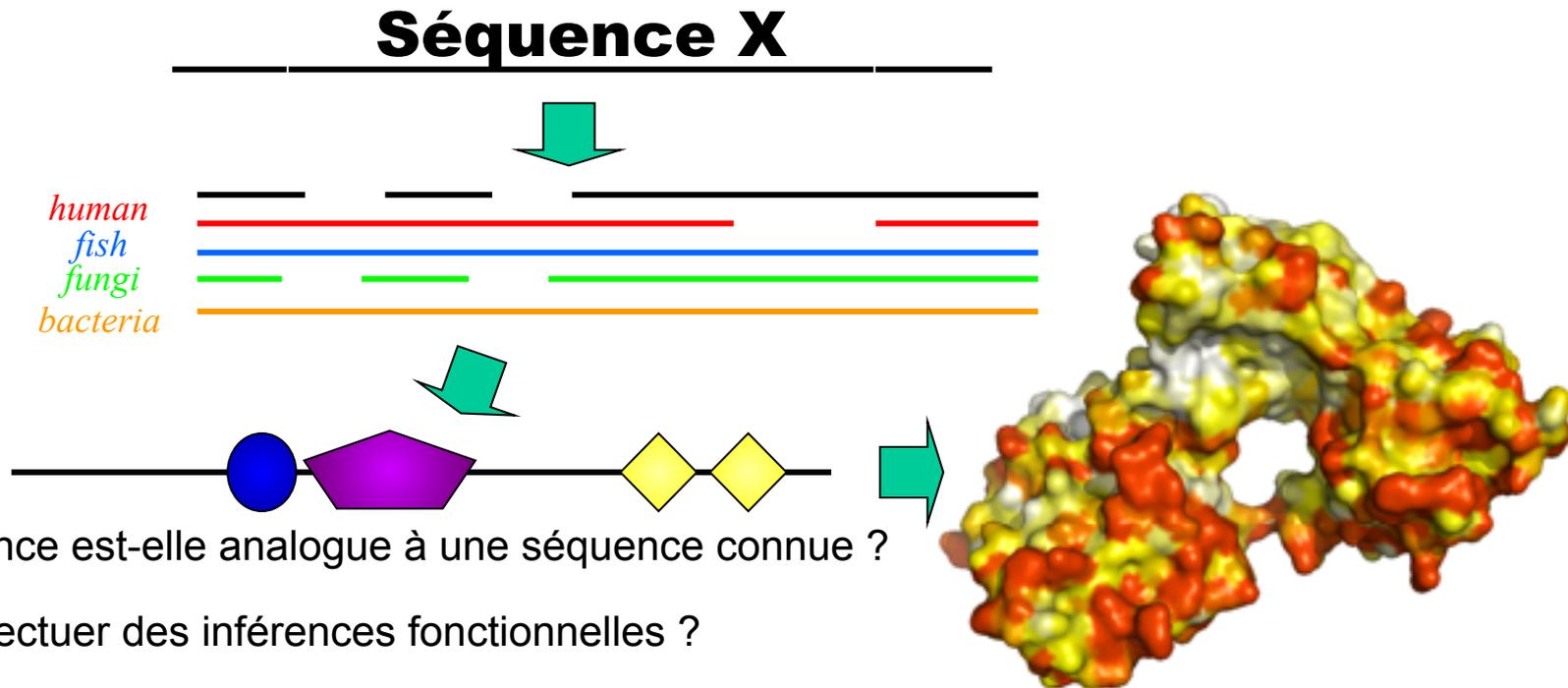
MAEAMDLGKDPNGP~~THSST~~LFVRRDDGSSMSFYVRPSPAKRRLSTLILHGGGTVCRVQEPGAVLLAQPGEA
LAEASGDFISTQYILDCVERNERLELEAYRLGPASAADTGS

GSEAKPGALAE~~GAAEPEPQRHAGRIAFTDAD~~
DVAILITYVKENARSPSSVTGNALWKAMEKSSLTQHSWQSLKDRYLKHLRGQEHKYLLGDAPV

PDEEEEEEEKVSQPEVGAAIKIIRQLMEKFNLDLSTVTQAF~~LKNSGELEATS~~SAFL
ASGQRADGYPIWSRQDDIDLQKDD~~EDTREALVKKFGAQN~~VARRIEFRKK

Et après la structure ?

Quelles informations peut-on extraire d'une analyse de séquence ?



- 1 - Ma séquence est-elle analogue à une séquence connue ?
- 2 - Puis-je effectuer des inférences fonctionnelles ?
- 3 - Comment découper ma séquence en modules fonctionnels repliés?
 - Quelles sont les régions fonctionnelles et structurales ?
- 4 - Peut-on prédire son organisation structurale ?
- 5 - Quelles sont les régions soumises à des pressions de sélection particulières ?

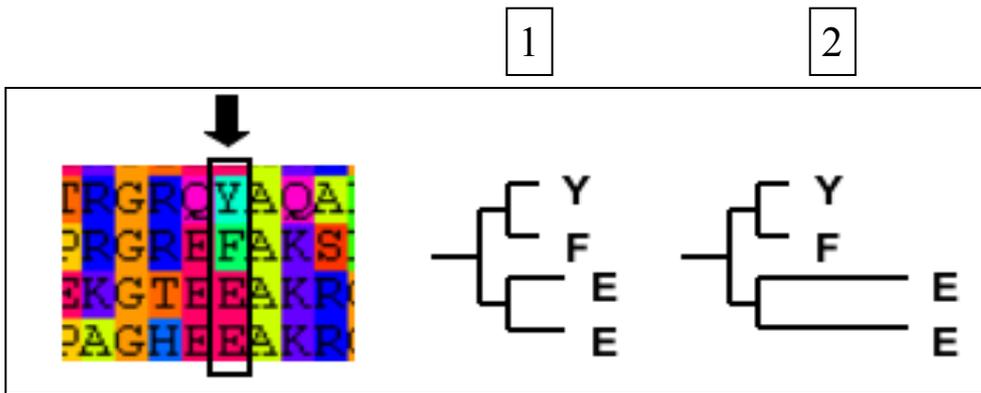
Analyse des conservations automatique

Une des meilleures sensibilités : ConSurf

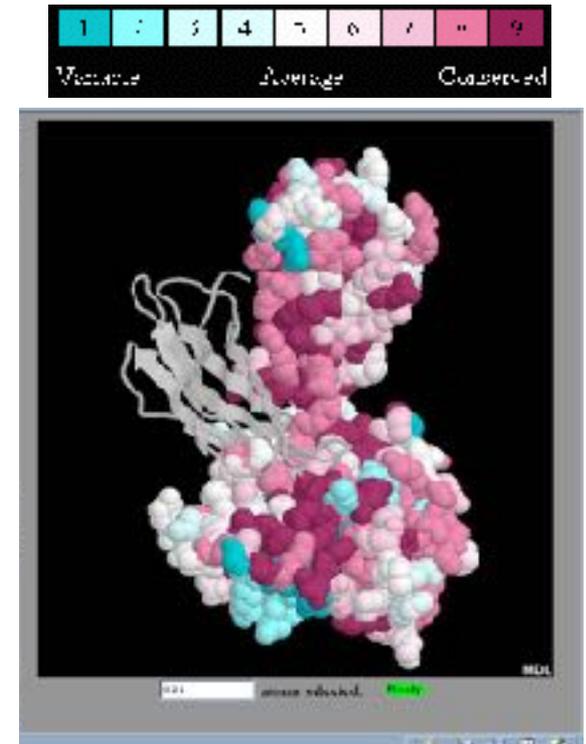
<http://consurf.tau.ac.il/>

N. Ben Tal, Nucleic Acids Research, 2005, Vol. 33, Web Server issue W299–W302

Méthode dont la sensibilité est accrue par une prise en compte des vitesses d'évolution à chaque position
(cf publi Rate4site Pupko et al Bioinfo 2002)



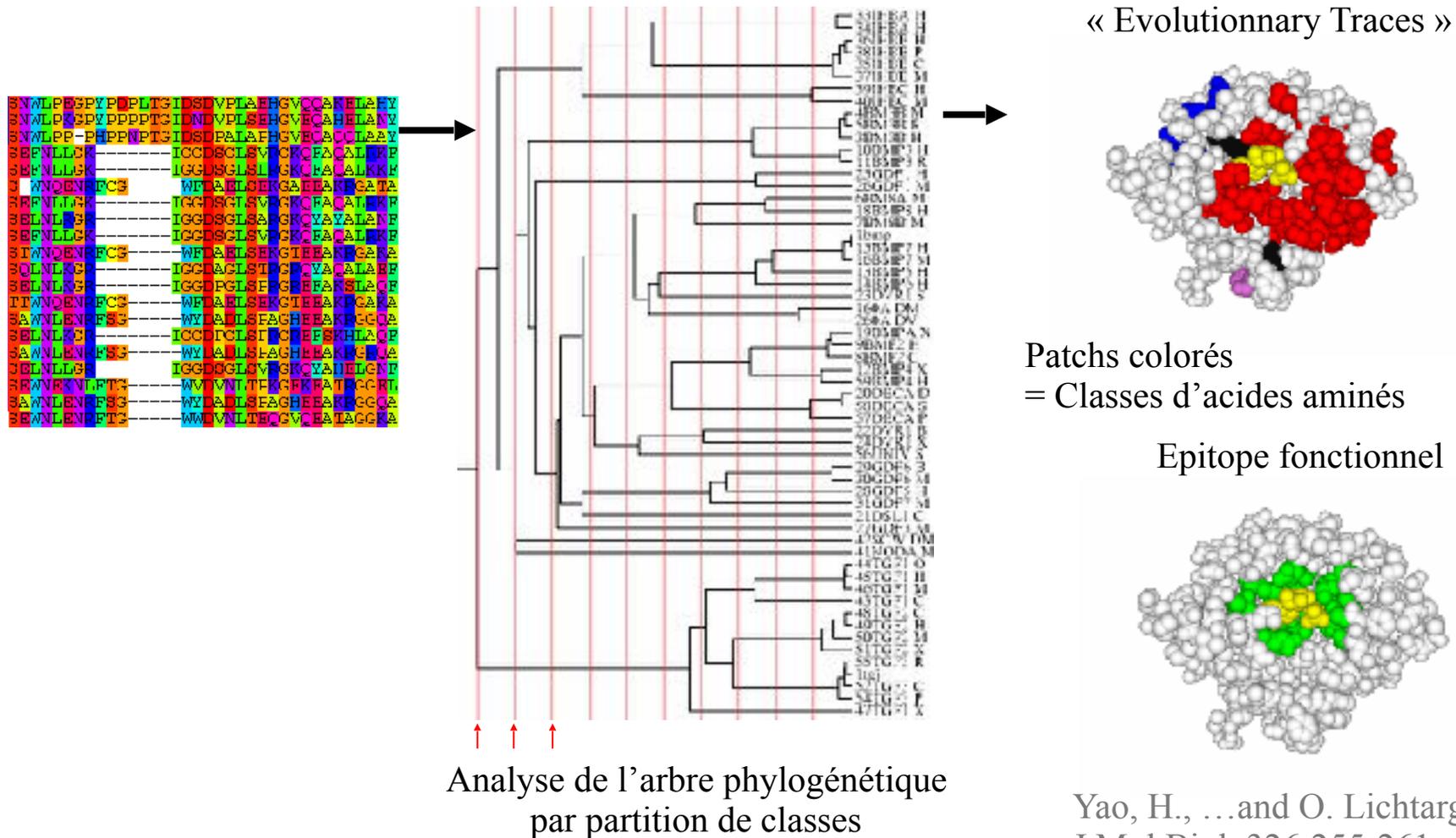
Pression de sélection plus importante sur le « E » dans le cas 2



Analyse des traces évolutives

<http://imgen.bcm.tmc.edu/molgen/labs/lichtarge/>

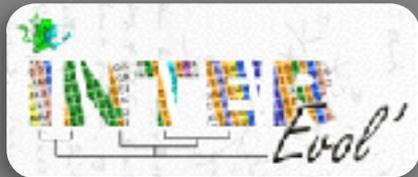
<http://www-cryst.bioc.cam.ac.uk/~jiye/evoltrace/evoltrace.html>



Yao, H., ...and O. Lichtarge. 2003.
J Mol Biol. 326:255-261.

Principe de la CoEvolution d'Interface :

Analyse Statistique → Extraction des caractéristiques → Score de Docking



Database of
Interfaces with
Evolutionary information

InterEvol:

~18,000 non-redundant interfaces
among which

~4,000 heteromeric interfaces

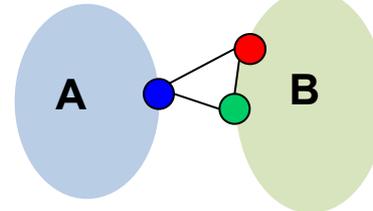
G. Faure, et al Nucleic Acids Res. (2012)

<http://biodev.cea.fr/interevol>

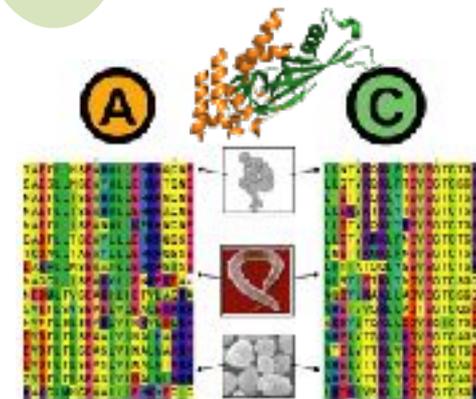
*J. Andreani, et al , Plos Comp Biol
(2012) ; Bioinformatics (2013)*

**InterEvScore : Discriminate
co-evolved interfaces**

**Contacts multi-
domaines**



**Patches
apolaires**

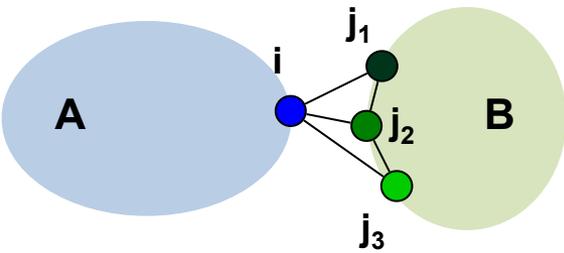


**Information
d'évolution**

Development of InterEvScore, a docking score taking evolution into account

Score inter-molecular interface contacts

For every possible couple or triad



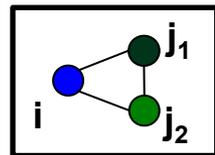
Find the most favorable environment ...

> 10 sequences

Contact propensities derived from InterEvol
(statistics on 1,289 interfaces)

 $i \text{---} j_1$ 0.53		 $i \text{---} j_1, j_2$ 0.04	
 $i \text{---} j_2$ -0.05		 $i \text{---} j_2, j_3$ -0.26	
 $i \text{---} j_3$ -0.11			

...taking into account evolution



Alignment for A

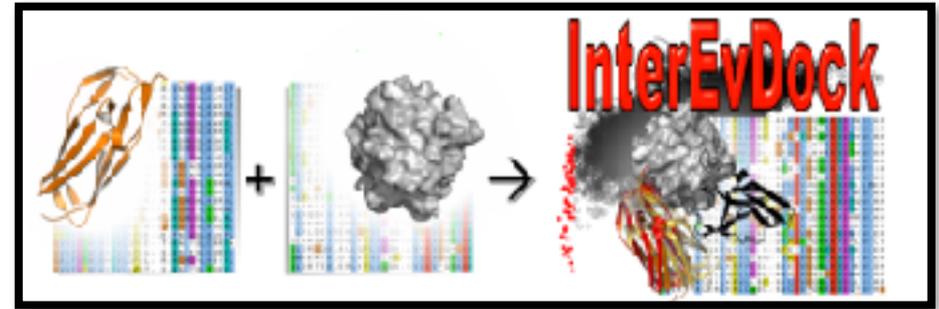
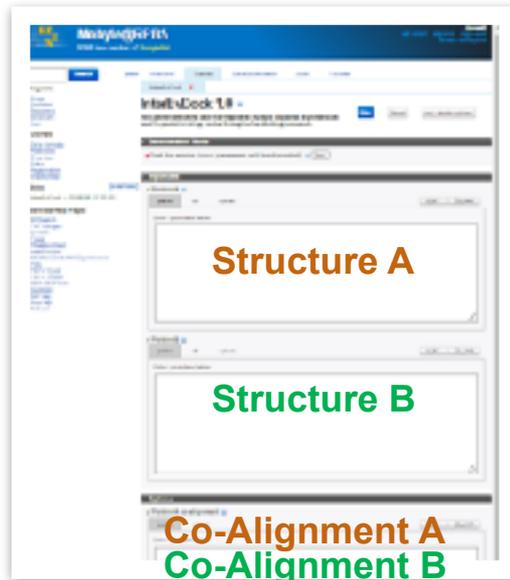
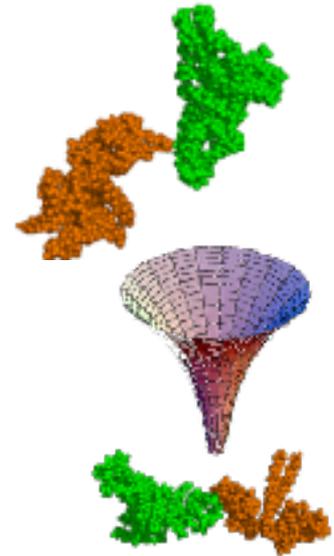
Alignment for B

- H. sapiens*
- M. musculus*
- D. rerio*
- ⋮
- S. cerevisiae*

	i	j1
<i>H. sapiens</i>	W F H I S L E -	I L L H
<i>M. musculus</i>	M F H I S I F C	I I I H
<i>D. rerio</i>	W F F H I V L E R	M Q L H
<i>S. cerevisiae</i>	W F F H I V L E R	M Q L H
	W F H I A I I	I I I H
	F Q I S L D P E I	L H R R
	W F F H I S L E -	I L L H
	M F H I S I I	I I I H
	W F F H I S L E -	I L L H
	W F F H Y T L E -	I L L H

	j2	j1
<i>H. sapiens</i>	L L N T V	Q S L F T E Y
<i>M. musculus</i>	I I F V	Q S L F T E Y
<i>D. rerio</i>	L F E N L	S S L M F Q Y
<i>S. cerevisiae</i>	L R D K L	S S L I K Q Y
	I M E V	Q S I F I I Y
	E T V F Q	L Y K E V E Y
	L L Q T V	Q S L F T E Y
	I I I V	Q S I F I I Y
	L I E T V	Q S L Y T E Y
	L L D T V	R S L S E S E Y

INTEREVDOCK SERVER TO ACCOUNT FOR CO-EVOLUTION INFORMATION IN DOCKING



Returns 10 models (~1h) selected from a Consensus among the top solutions of :

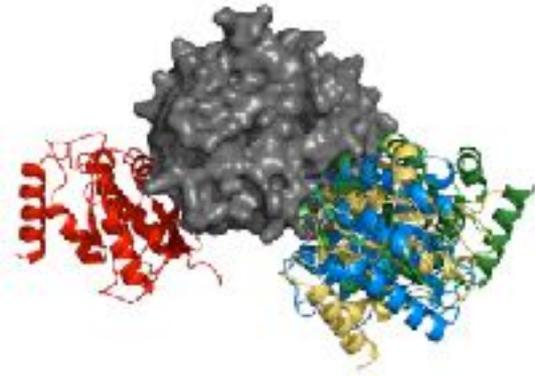
- **InterEvScore** (residue based / coevolution),
- **SOAP-PP** (statistical atomic based) (*A. Sali's lab*),
- **FRODOCK** (rigid-body + physics based) (*P. Chacon's lab*)

J. Yu et al, Nucleic Acids Research (2016)

Running at RPBS server (*coll. P. Tufféry*)

<http://bioserv.rpbs.univ-paris-diderot.fr/services/InterEvDock/>

InterEvScore success rates for individual test cases

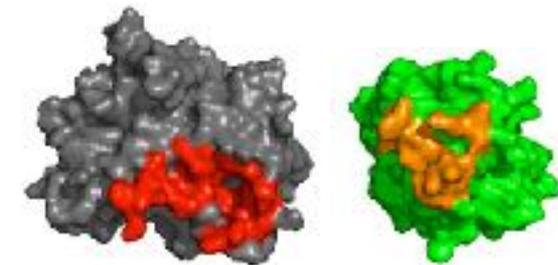


- Can we predict the **binding mode**?

With 10 models per test case

46 % success

(25 successful cases out of 54 test cases)



- Can we predict **binding sites**?

With 10 models per test case:

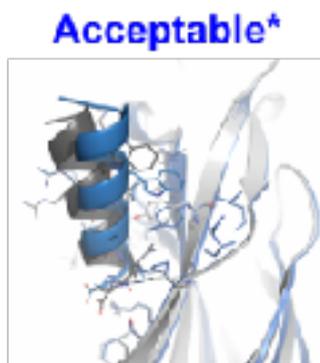
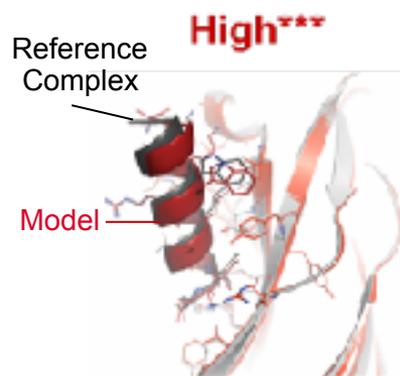
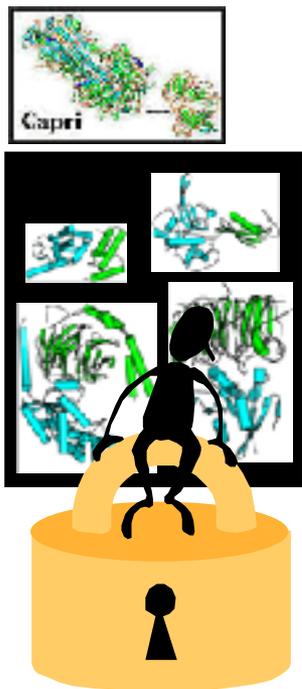
65 % success (35 cases) for both partners

91 % success (49 cases) for at least one partner
out of 54 test cases

- **Blind test for the different predictive strategies**

About every six months structures of individual binding partners are released to the community.

→ X-ray structure of every query complex is kept secret.



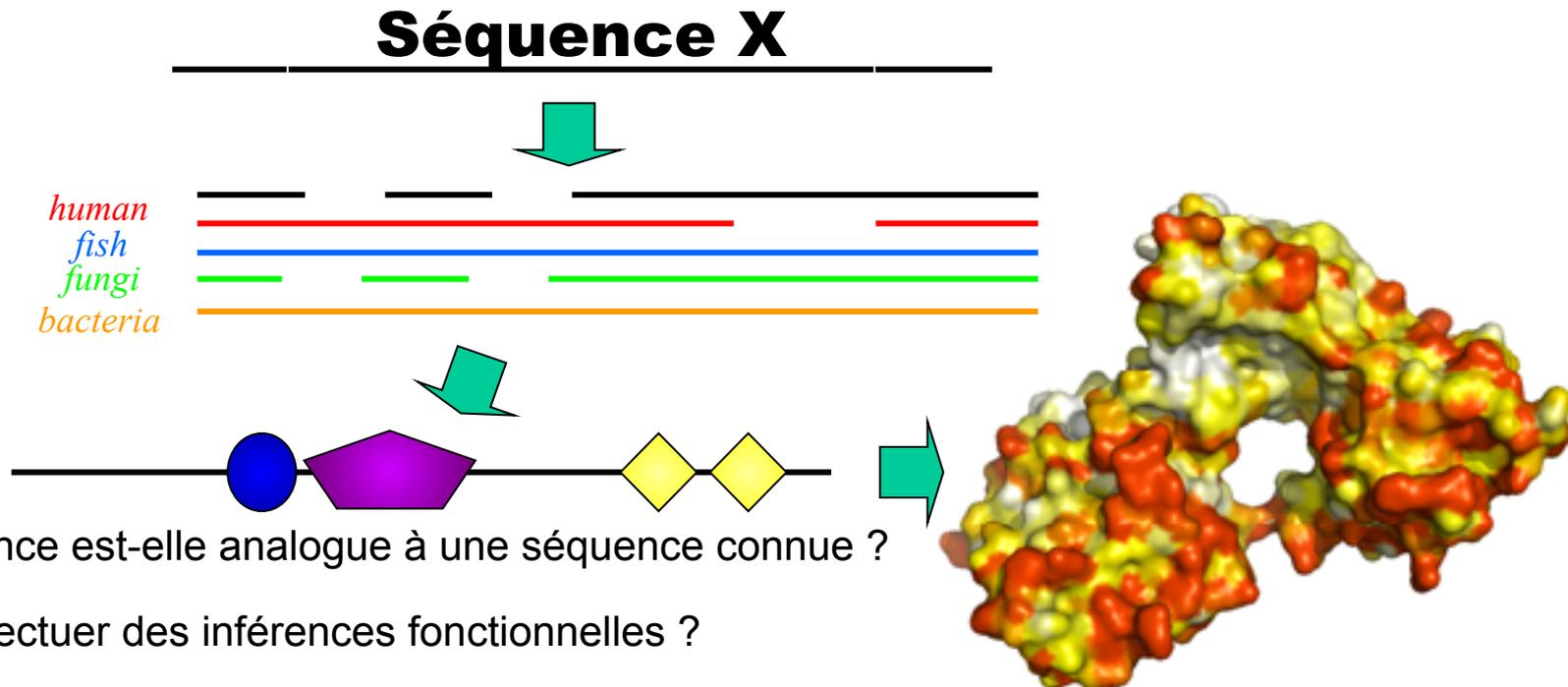
THE THREE MOST RECENT CAPRI EVALUATION MEETINGS.

CAPRI evaluation meeting, year and targets						
	2009, 12 targets		2013, 14 targets		2016, 20 targets	
Rank	Group	Success	Group	Success	Group	Success
1	Vajda/Kozakov	6/4 ^{***} /2 ^{**}	Bonvin	9/1 ^{***} /3 ^{**}	Guerois	10/1^{***}/8^{**}
2	Zacharias	6/4 ^{***} /1 ^{**}	Bates	8/2 ^{**}	Zacharias	10/3 ^{***} /2 ^{**}
3	Zou	6/3 ^{***} /2 ^{**}	Vakser	7/1 ^{***}	ClusPro	9/3 ^{**}
4	Eisenstein	6/3 ^{***} /1 ^{**}	Kozakov/ Vajda	6/2 ^{***} /3 ^{**}	Kozakov/ Vajda	8/3 ^{***} /2 ^{**}
5	Wolfson	6/3 ^{***} /1 ^{**}	Shen	6/1 ^{***} /3 ^{**}	Seok	8/3 ^{***} /2 ^{**}
6	Weng	6/2 ^{***} /2 ^{**}	Fernandez-Recio	6/1 ^{***} /3 ^{**}	Fernandez-Recio	7/1 ^{***} /3 ^{**}
7	Zhou	6/2 ^{***} /2 ^{**}	ClusPro	6/4 ^{**}	Zou	7/1 ^{***} /2 ^{**}
8	Bonvin	6/1 ^{***} /4 ^{**}	Zou	6/1 ^{***} /2 ^{**}	Weng	6/1 ^{***} /4 ^{**}
9	ClusPro	5/1 ^{***} /3 ^{**}	Zacharias	6/1 ^{***}	Vakser	6/2 ^{***} /2 ^{**}
10	Fernandez-Recio	5/2 ^{**}	Eisenstein	5/1 ^{***} /2 ^{**}	Bates	6/3 ^{**}

*Kozakov D et al (2017)
Nat. Protocols*

Et en absence de structure ?

Quelles informations peut-on extraire d'une analyse de séquence ?



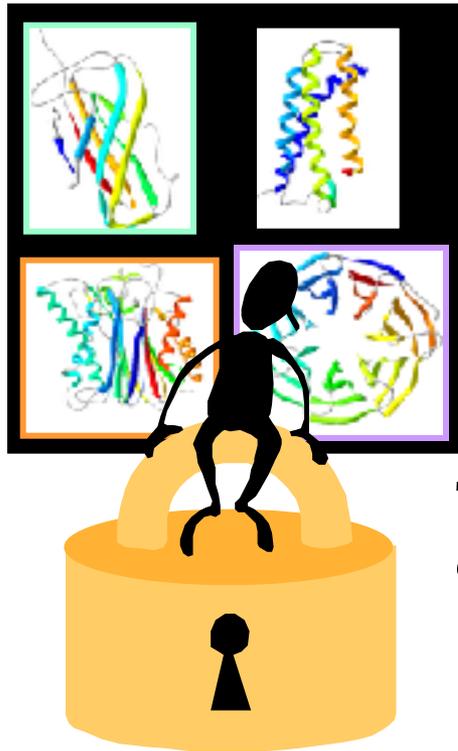
- 1 - Ma séquence est-elle analogue à une séquence connue ?
- 2 - Puis-je effectuer des inférences fonctionnelles ?
- 3 - Comment découper ma séquence en modules fonctionnels repliés?
 - Quelles sont les régions fonctionnelles et structurales ?
- 4 - Peut-on prédire son organisation structurale ?
- 5 - Quelles sont les régions soumises à des pressions de sélection particulières ?

Comment évaluer les meilleures méthodes : CASPn

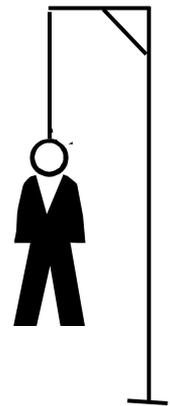
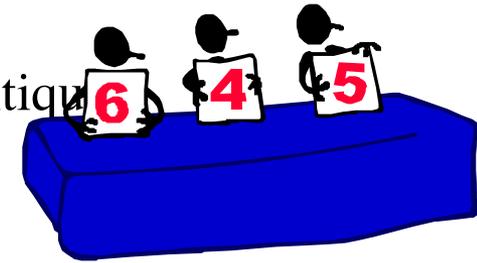
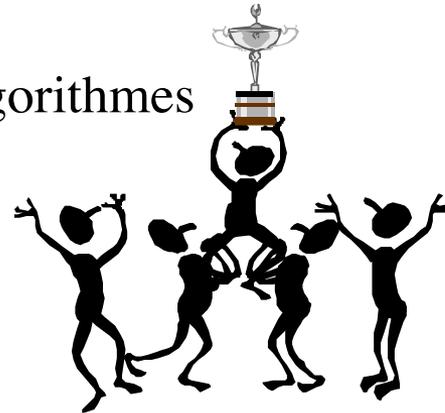
<http://predictioncenter.org/>

Tous les deux ans, un ensemble de nouvelles structures sont gardées secrètes durant la prédiction.

Classement, analyse critique des points forts et faibles de chaque méthode.

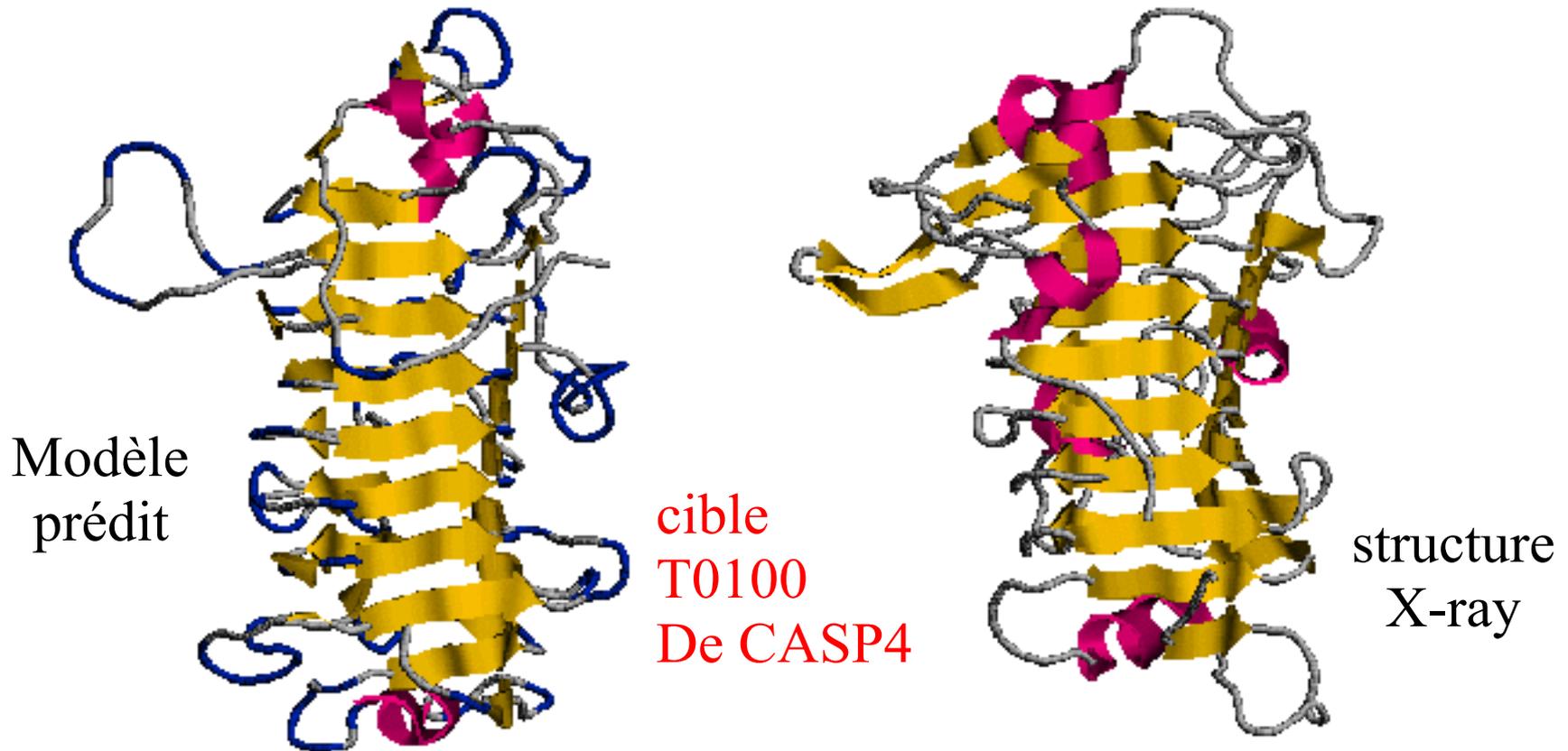


Test en aveugle des différents algorithmes



Exemple de prédiction

Pas d'homologie de séquence évidente...



Résultats CASP11 (2014):

<http://predictioncenter.org/casp11/results.cgi>

Best servers for template-based modelling

❖ {Zhang-Server, QUARK} (**I-TASSER**)

❖ **nns** (Lee group)

❖ **RosettaServer** – gets much better with best models

❖ **MyProteinMe** (Skwark group)

❖ **HHpred** ranked best in CASP9 (2010) and has been integrated into other servers/approaches

Controle continu avec CAMEO

CAMEO | Continuous Automated Model Evaluation

★ Try our BETA! Home 3D - Protein Structure QE - Model Quality Estimation CP - Contact Prediction More Login

CAMEO continuously evaluate the accuracy and reliability of predictions

- 3D - Protein Structure**
280 weeks, 4922 targets, 8 predictors.
- QE - Model Quality Estimation**
170 weeks, 23262 structural models, 18 predictors.
- CP - Contact Prediction**
12 weeks, 68 targets, 4 predictors.

★ Predictions in all categories are evaluated against reference structures released by the PDB on a weekly basis.

CAMEO is a community project

⇒ CAMEO continuously applies quality assessment criteria established by the protein structure prediction community. Since the accuracy requirements for different scientific applications vary, there is no "one fits all" score. CAMEO therefore offers a variety of scores - assessing different aspects of a prediction (coverage, local accuracy, completeness, etc.) to reflect these requirements.

⇒ CAMEO is a community project - please feel free to suggest additional/alternative ways how CAMEO can support users and developers of structure prediction.

Join CAMEO today...

We invite developers of prediction methods to participate in the continuous evaluation by registering their servers (**REGISTER**). We also invite developers of scoring and evaluation methods to suggest alternative scoring schemes. Please contact us **directly**.

Servers of the following groups are registered so far:

A. Salt ¹, L. McGuffin ², T. Schwede ³, J. Seeding ⁴, D. Baker ⁵, A. Fiser ⁶, M. Sternberg ⁷, Y. Zheng ⁸, G. Floudas ⁹, S. Tosatto ¹⁰, J. Xu ¹¹, Y. Zhou ¹², O. Brock ¹³, B. Walther ¹⁴, A. Elofsson ¹⁵, D. Labudde ¹⁶, C. Venclovas ¹⁷, J. Cheng ¹⁸, D. Taghan Bishop ¹⁹, Y. An-Suel ²⁰.

<http://cameo3d.org/>

Analyse des séquences protéiques en absence d'homologues détectables

- ❖ **Analyse des peptides signaux**

 - G. von Heijne SignalP <http://www.cbs.dtu.dk/services/SignalP>

- ❖ **Low complexity**

- ❖ **Domaines Coiled-coils**

 - A. Lupas COILS (http://www.ch.embnet.org/software/COILS_form.html)

 - B. Berger Multicoil (<http://groups.csail.mit.edu/cb/multicoil/cgi-bin/multicoil.cgi>)

 - M. Delorenzi Marcoil (<http://bcf.isb-sib.ch/webmarcoil/webmarcoilC1.html>)

Analyses combinées dans SMART : <http://smart.embl-heidelberg.de/>

Nombreux outils regroupés dans le Bioinformatics Toolkit

<http://toolkit.tuebingen.mpg.de/>

- ❖ **Segments transmembranaires**

 - cf diapo suivante

- ❖ **Structures Secondaires**

 - cf diapos suivantes

- ❖ **Reconnaissance de repliement Threading**

Sources d'inspiration pour le cours et pour aller plus loin

- **Sur les alignements et l'évolution structurale des protéines**
<http://www.people.virginia.edu/~wrp/papers/ismb2000.pdf>
- **Très bon chapitre d'ouvrage par Koonin accessible on-line**
<http://www.ncbi.nlm.nih.gov/books/bv.fcgi?rid=sef.chapter.166>

MERCI MERCI MERCI:

Jessica Andreani : jessica.andreani@cea.fr

Raphaël Guerois : guerois@cea.fr

CEA Saclay, Institut Joliot

91190 Gif sur Yvette Cedex

RéNaFoBiS

Grand tournoi
bb-foot.

Regardez cet air combatif !!!



Allez-vous le laisser gagner ?