



# Bioinformatique structurale : Un enjeu de l'après génome



IN'Tech Lyon le 23/10/2003

**Gilbert DELEAGE**

**PBIL-IBCP-CNRS UMR 5086**

Pôle Bioinformatique Lyonnais Lyon-Gerland

7, passage du Vercors

69367 Lyon cedex 07

Tél: +33 (0)4 -72-72-26-55

fax: +33 (0)4 -72-72-26 -01

mel: g.deleage@ibcp. fr



<http://www.ibcp.fr>

[g.deleage@ibcp.fr](mailto:g.deleage@ibcp.fr)



## Equipe Bioinformatique et RMN structurales

### Permanents

DELEAGE G. PR1 UCBL, **Bioinformatique**  
PENIN F. IR2-CNRS, NMR, Biophysique, Biochimie  
BETTLER E. MC2, **Bioinformatique**  
BLANCHET C. IR2-CNRS, **Bioinformatique**  
BOCKMANN A. CR2-CNRS, NMR  
COMBET C. CR2-CNRS **Bioinformatique**  
GEOURJON C. IR2-CNRS, **Bioinformatique**  
HUET E. CR2-CNRS, Biologie-Biochimie  
LAVERGNE J-P CR1-CNRS, Biochimie  
MONTSERRET R. IE2-CNRS, NMR, Biophysique

### Post-docs and CDD

DORKELD F. (EC), **Bioinformatique**  
LECLUSE A. (French ministry) **Bioinformatique**  
LACORNE N. (ACI GRID) **Bioinformatique**  
MISSEREY S. (Genopôle Rhône-Alpes) **Bioinformatique**

### Etudiants

BOULANT S. Biochimie  
GIRAUD N. NMR  
RATINIER M. Biochimie  
SAPAY N. **Bioinformatique**  
VERNOIS A. GRID computing, **Informatique**



**PBIL-IBCP**



Avant la **Bio-informatique**

Activité  
Biologique  
connue

Etude  
Biochimique  
Structure 3D

Séquence  
Protéine

Gène  
Mutagenèse

Relations  
structure-  
activité

**BIO-INFORMATIQUE**

Bases de données  
Prédiction des gènes

Identification de protéines  
Prédiction sites/signatures  
Prédiction de structure  
Modélisation moléculaire

Stockage  
Classification  
Intégration  
Criblage

Protéomique

Séquences  
génomiques

Séquences  
Protéiques

Prédiction  
Activités  
biologiques

Etudes  
Biochimiques  
Structures 3D

Génomique  
structurale

Aujourd'hui (depuis les programmes de  
séquençages massifs et la **Bio-informatique**)





# Réalisations bioinformatiques du groupe PBIL-IBCP



<http://pbil.ibcp.fr>

## o Database :

- o **HCVDB** : Base de données de séquences annotées d'HCV; <http://hepatitis.ibcp.fr>)

## o Méthodologies :

- o **MLRC** : Prédiction de la structure secondaire des protéines (Guermeur *et al.*, **Bioinformatics**, 1999) Coll. LIP6
- o **ProScan** : Scan PROSITE avec recherche pondérée par un système original
- o **PattInProt** : Algorithme de recherche original de signatures dans les banques. Outil de PROTEOMIQUE
- o **PROCSS** : Compatibilité des structures protéiques grâce aux structures secondaires (Geourjon *et al.*, **Protein Sci.**, 2001; Errami *et al.*, **Bioinformatics**, 2003)
- o **SOPM** : Méthode auto-optimisée de prédiction de structure secondaire (Geourjon *et al.*, **Protein Eng.**, 1994)
- o **SOPMA** : Méthode auto-optimisée de prédiction de structure secondaire avec alignements (Geourjon *et al.*, **Comput Appl Biosci.**, 1995)
- o **SuMo** : Détection de sites 3D dans les protéines (Jambon *et al.*, **Proteins**, 2003, Brevet international)

## o Logiciels :

- o **AnTheProt** : Logiciel intégré d'analyse de séquences client/serveur (Deléage *et al.*, **Comput Appl Biosci.**, 1988; Deléage *et al.*, **Comput Biol Med.**, 2001; <http://antheprot-pbil.ibcp.fr>)
- o **AnTheNuc** : Analyse de sites de restriction
- o **Bioread** : Interface graphique pour **extractblast** and **extractfasta** programmes "parser".
- o **DicroProt** : Analyse des spectres de dichroïsme circulaire des protéines (Deléage G et Geourjon C. **Comp. Appl. Biosci.**, 1993)
- o **MPSA** : integrated protein sequence analysis with client/server capabilities (Blanchet *et al.*, **Bioinformatics**, 2000; <http://mpsa-pbil.ibcp.fr>)
- o **SecTrace** : secondary structure plot

## o Services Web:

- o **NPS@** : Analyse intégrée de séquences sur le Web (Combet *et al.*, **Trends Biochem Sci**, 2000; Perrière, Combet *et al.*, **Nucleic Acids Res.**, 2003; <http://npsa-pbil.ibcp.fr>)
- o **Geno3d** : Modélisation moléculaire de protéines à grande échelle (Combet *et al.*, **Bioinformatics**, 2002; <http://geno3d-pbil.ibcp.fr>)
- o **SuMo** : protein common 3D sites detection (Jambon *et al.*, **Proteins**, 2003; <http://sumo-pbil.ibcp.fr>)

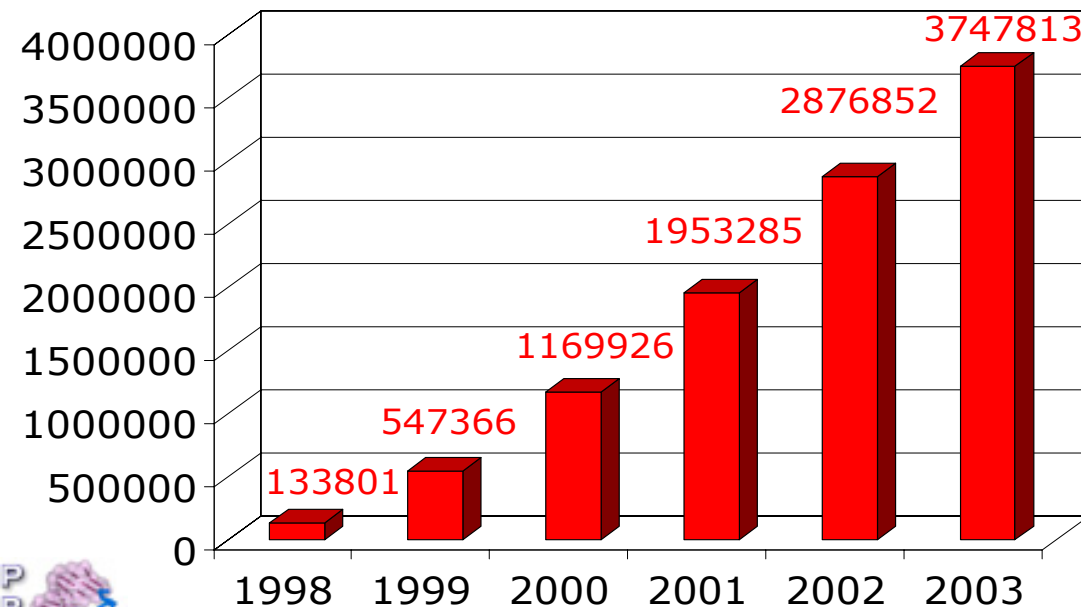
## o Applications biologiques:

- o Très nombreuses en collaboration avec des biologistes

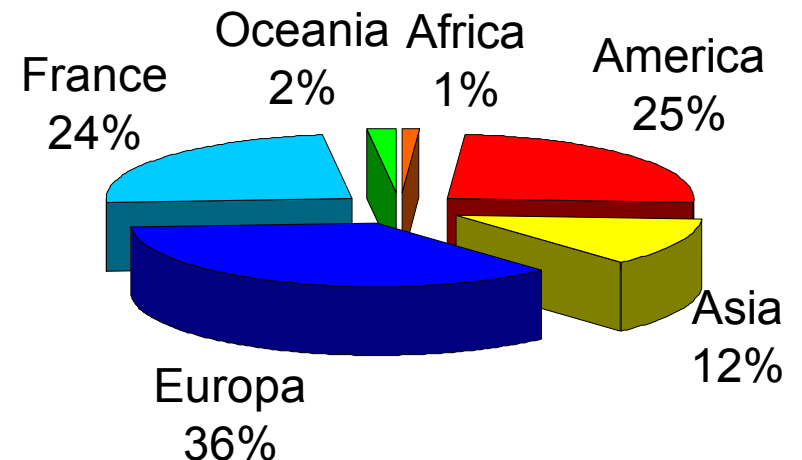




- **NPS@** <http://npsa-pbil.ibcp.fr>
- **Partie «protéine et structure»** du Pôle BioInformatique Lyonnais
- Interconnexion de **46 méthodes** d'analyse de séquences de protéine
- **Récupération automatique des données** dans des logiciels clients/serveurs d'analyse biologique MP5A, AnTheProt, Clustal X, RasMol, ...
- Liens hypertextes sur les données de **17 banques de données internationales** (SWISS-PROT, PROSITE, PDB, SCOP,...).
- **Mise à jour automatique** des bases de données disponibles sur le serveur (SP, SP-TrEMBL, NRL3D, Nr, PDB, PROSITE, etc.).
- **Références internationales:** Expasy, University of California, InfoBioGen, RSCB(PDB),....



**2530 analyses / jour en 2002**  
**3584 analyses / jour depuis janvier 2003**





## o GRID computing :

- o **GPS@** : Portail d'analyse de séquences de protéines sur la Grille (Déploiement de NPS@ sur la grille, EU FP5 DATAGRID)
- o **GriPPS** : Services de recherche de signatures sur la grille (Ministère recherche ACI-GRID, CNRS-IN2P3, ENSL-LIP)
- o **RUGBI** : Prédiction de structures secondaires sur la grille GRID (Ministère recherche et CNRS-IN2P3 (V. BRETON), Biopôle Clermont-Limagne, CS, ECP)
- o **e-Toile** : Projet de grille expérimentale (Ministère recherche RNTL)
- o **GiGn** : Réseau de recherche français de Grille pour la GeNomique (Ministère recherche Action IMPG, with V. BRETON)
- o **HealthGrid** : Organisation de conférences (HealthGrid Conference, Lyon FR, January 2003)

## o Bioinformatique structurale et clinique :

- o **CPS@** : Version pour cluster de NPS@ (PBS, PVM and MPI algorithms)
- o **euHCVdb** : Base européenne de séquences HCV (EU FP5 HepCVax QLK2-CT-2002-01329 ; EU FP6 viRgil)
- o Collaboration avec D. KAHN (**PRODOM**)
  - o Amélioration des domaines PRODOM grâce aux prédictions de structures
  - o Connexion avec Geno3D
- o **StrucAnnot** Annotation structurale de genome *Arabidopsis thaliana* : (consortium France, Belgium (P. ROUZE, Gent) and SIB (A. BAIROCH, Geneva))
- o **MSFold** : Modélisation de structure 3D à partir des contraintes de MS. Protéomique structurale, (coll. E. Forest, IBS, Grenoble)
- o **MADPROTS** : Modélisation, analyses et serveur Web pour amarrage moléculaire
- o **SYNAPSE** : Système expert d'annotation de protéines

# Organisation hiérarchique des protéines

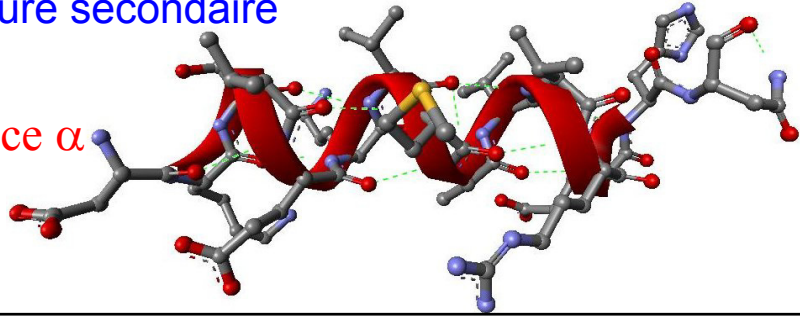
Structure primaire = séquence = mot écrit avec un alphabet de 20 lettres

MKLD**E**IARLAGVSRRTTASYVINGKAKQYRVSDKTVEKVMMAVVREHNYHPNAVAAGLRAGR

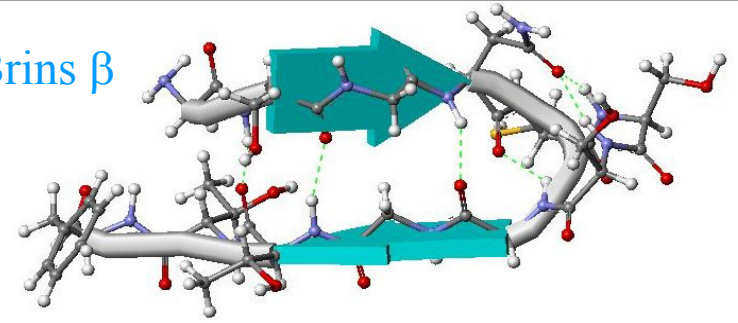


Structure secondaire

Hélice  $\alpha$

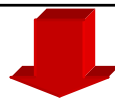


Brins  $\beta$

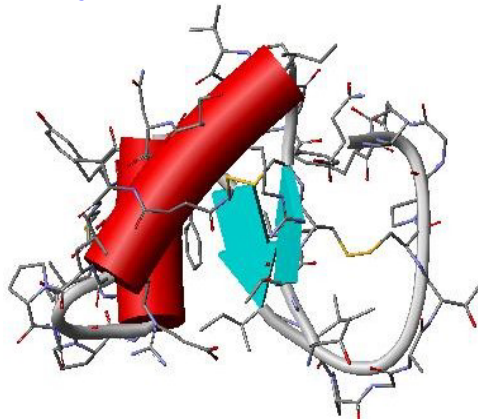


Structure secondaire = mot alphabet de 3 à 10 lettres

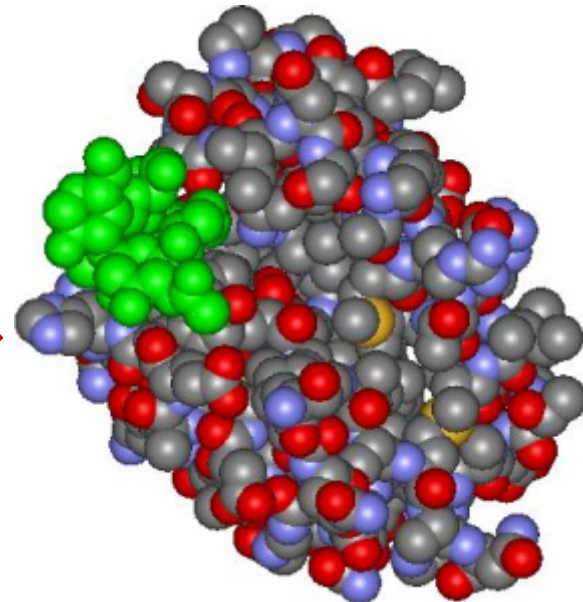
CCHHHHHHHHHHCCCEEEETTTEEEEECCCCHHHHHHHHHHHCCHHHHHHHHHGCCCC



Structure tertiaire = objet 3D



Fonction





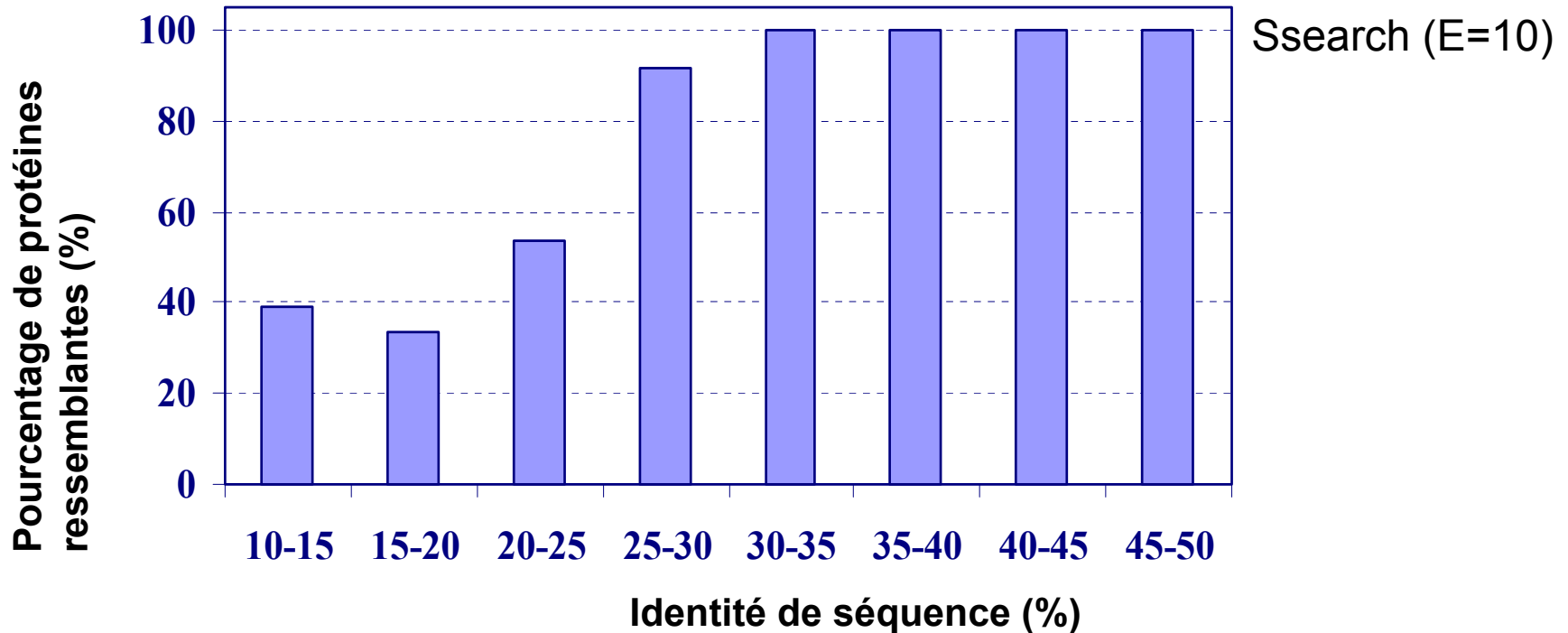
# Relation identité de séquence ressemblance structurale



1AJJ : LDL receptor

1CR8 : Low Density Lipoprotein Receptor Related Protein

		10	20	30	40																																									
		.....	.....	.....	.....	..																																								
1.pdb1ajj.ent	P	-	C	S	A	F	E	F	H	C	-	L	S	G	E	C	I	H	S	W	R	C	D	G	G	P	D	C	K	D	K	S	D	E	E	N	C	A	-	-	37					
2.pdb1cr8.ent	P	G	G	C	H	T	D	E	F	Q	C	R	L	D	G	L	C	I	P	L	R	W	R	C	D	G	D	T	C	M	D	S	D	E	K	S	C	E	G	V	42					
Identity	*	*	:	**	:	*	*	.	*	**																																				

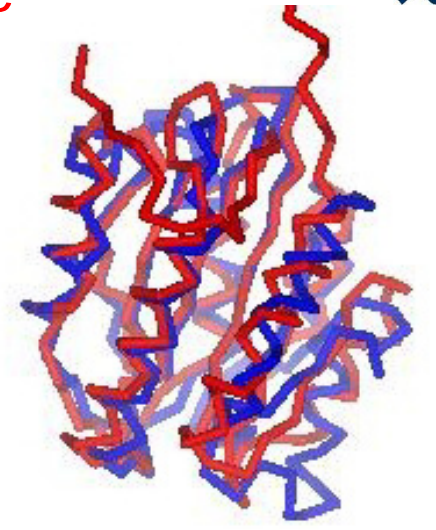
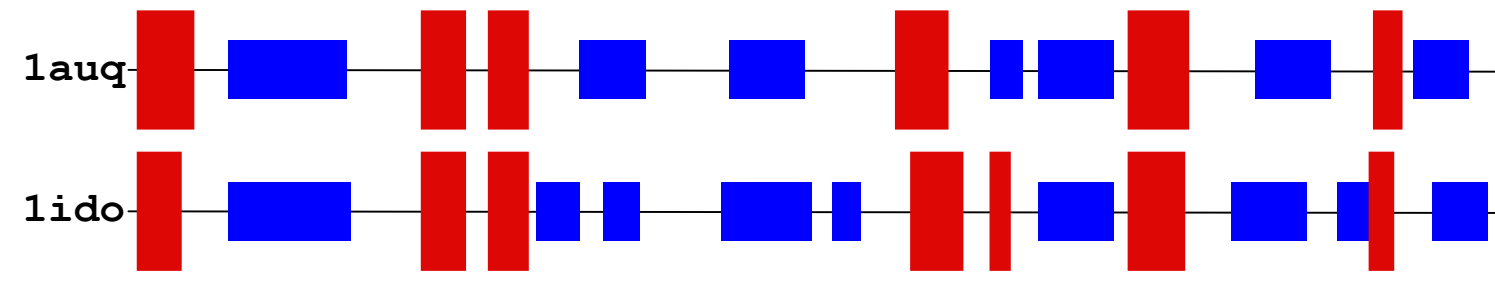


Nécessité de discriminer dans l'intervalle 10-30% d'identité



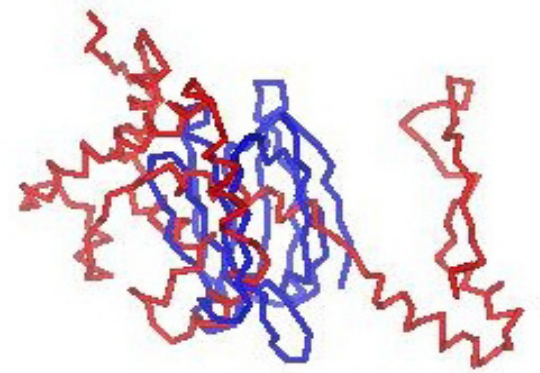
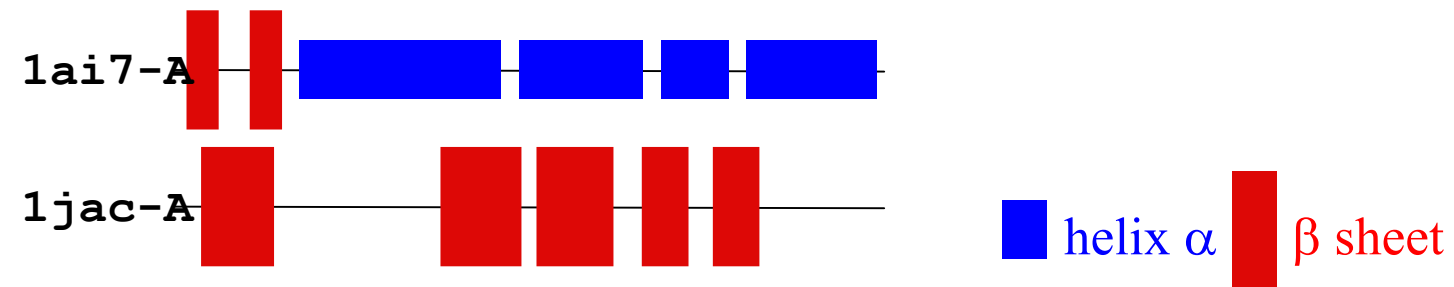
**Sov = 81**

15,9% identité ; Protéines apparentées (FSSP : RMSD = 2,3Å ; Z-score = 19,9)



**Sov = 9**

16% identité ; Protéines non apparentées



Les structures secondaires permettent de discriminer les protéines ayant des structures 3D proches « zone floue » de 10-30%



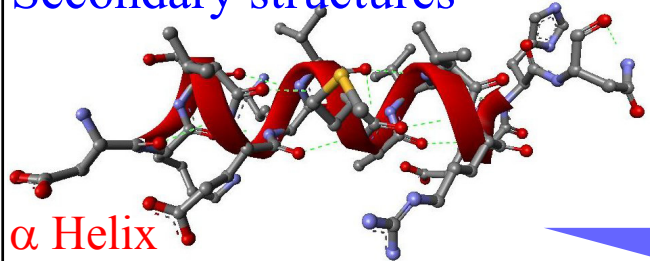
- Modélisation homologique « classique »
  - 2 protéines qui ont plus de 30% d'identité de séquences ont 80% de leurs C $\alpha$  superposables avec un écart quadratique moyen de 1 Å (RMSD=1Å)
- Modélisation par analogie (à faible taux d'identité)
  - 2 protéines qui ont la même fonction et une topologie « probablement » identique en dépit de l'absence de similarité importante (arguments expérimentaux ou de structures secondaires).
- Threading (en cours de développement)
  - Une séquence est testée sur une librairie de repliements pour déterminer sa compatibilité structure-séquence probable, méthode d'alignement séquence-structures tridimensionnelles.
- *Ab initio* (en progression dans CASP5, folding@home Stanford)
  - Structure directement déduite de la séquence à partir de règles empiriques.



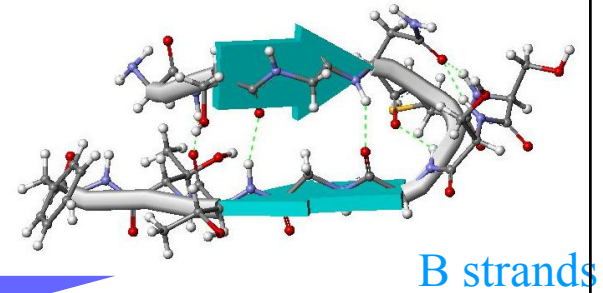
## Protein sequence

MKLDEIARLAGVSRRTASYVINGKAKQYRVSDKTVEKVMVREHNYHPNAVAAGLRAGR

## Secondary structures



# Geno3 Tools



Adresse [http://geno3d-pbil.ibcp.fr/cgi-bin/geno3d\\_automat.pl?page=/GENO3D/geno3d\\_home.html](http://geno3d-pbil.ibcp.fr/cgi-bin/geno3d_automat.pl?page=/GENO3D/geno3d_home.html)

**Pôle Bio-Informatique Lyonnais**  
**Geno3D**  
Geno3D is the IBCP contribution to PBIL in Lyon, France

[\[HOME\]](#) [\[GENO3D\]](#) [\[HELP\]](#) [\[REFERENCES\]](#) [\[NEWS\]](#)

Monday, March 5th 2001 : **GENO3D** is now available ([see news](#))

GENO3D : AUTOMATIC MODELING OF PROTEINS THREE-DIMENSIONAL STRUCTURE

[Abstract] [\[GENO3D help\]](#) [Original server]

Database :

Sequence name (optional) :

Paste a protein sequence below : [help](#)

```

DPDDVQFYTIENSVPVHLLRTGDEFATGTFFFDCKPCRLTHTWQTNRALG
LPPFLNSLPQSEGATNFGDIGVQDKRRGVTQMGNTNYITEATIHHPAEV
GYSAPYYSEASTQGFKTPIAAGRGAQTDENQAADGNPRYAFGRQHGQ
KTTTGTPEPFTYIAHQDTGRYPEGDIQINFNLPVINDVLLPTDPI
GGKTGINYINIFNTYGLTALNNVPPVYPNGQIUWKEFDLTKPRLHVNA
PFVCQNNCPGOLFVRVAPNLTNEYDPDASANMSRIVTYSDFWVKGLVFK
AKLRASHTUNPIQMSINVDNQFNYPVPSNIGGMKIVYEKSQLAPRKLY

```

Number of PSI-BLAST run

<http://geno3d-pbil.ibcp.fr>



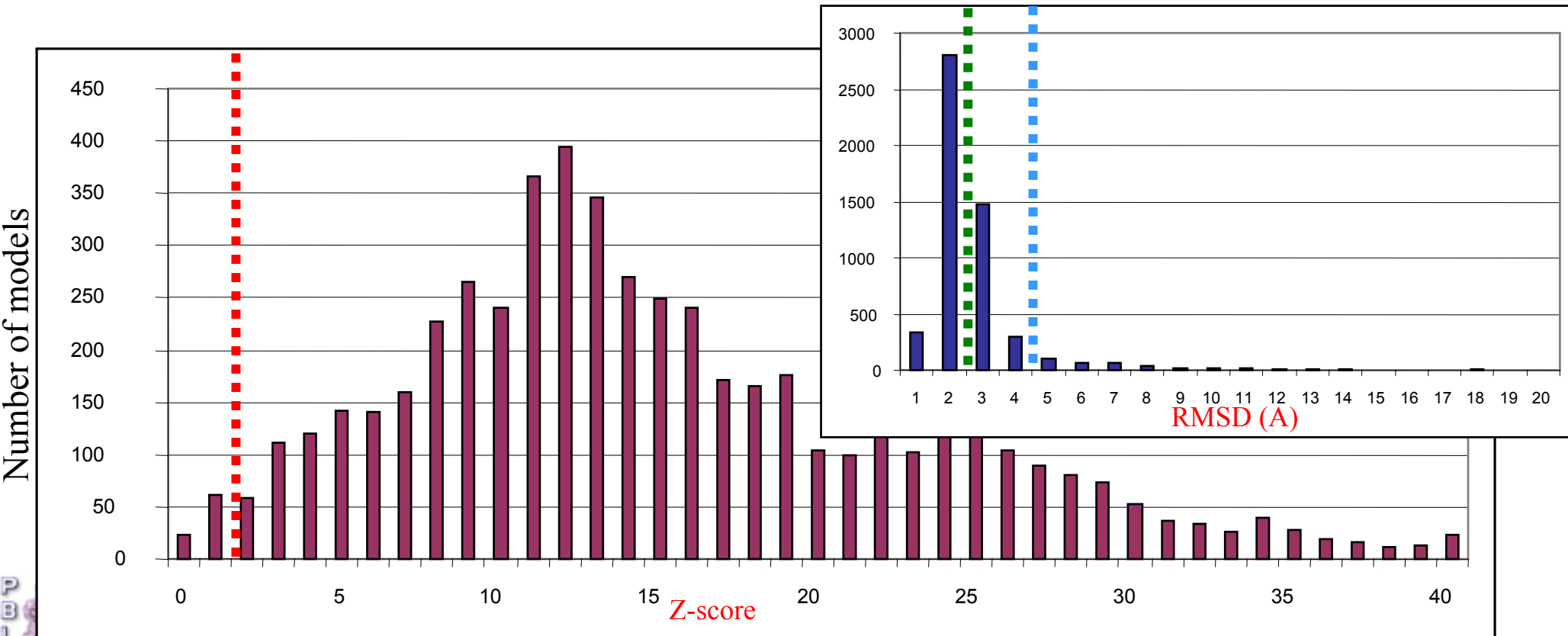


# 3D-Crunch : Modélisation à grande échelle



## Stratégie

- ✓ Modélisation à faible taux d'identité (entre 10-35%)
- ✓ 1,315 protéines de structures 3D connues (pdb 25%)
- ✓ PSI-BLAST sur chaque entrée (max 5 run).
- ✓ Les 5,390 protéines ayant un pourcentage d'identité entre 10 et 35% ont été modélisées
- ✓ Les calculs ont été réalisés au CC-IN<sub>2</sub>P<sub>3</sub> Villeurbanne (France) 7,881 heures de CPU



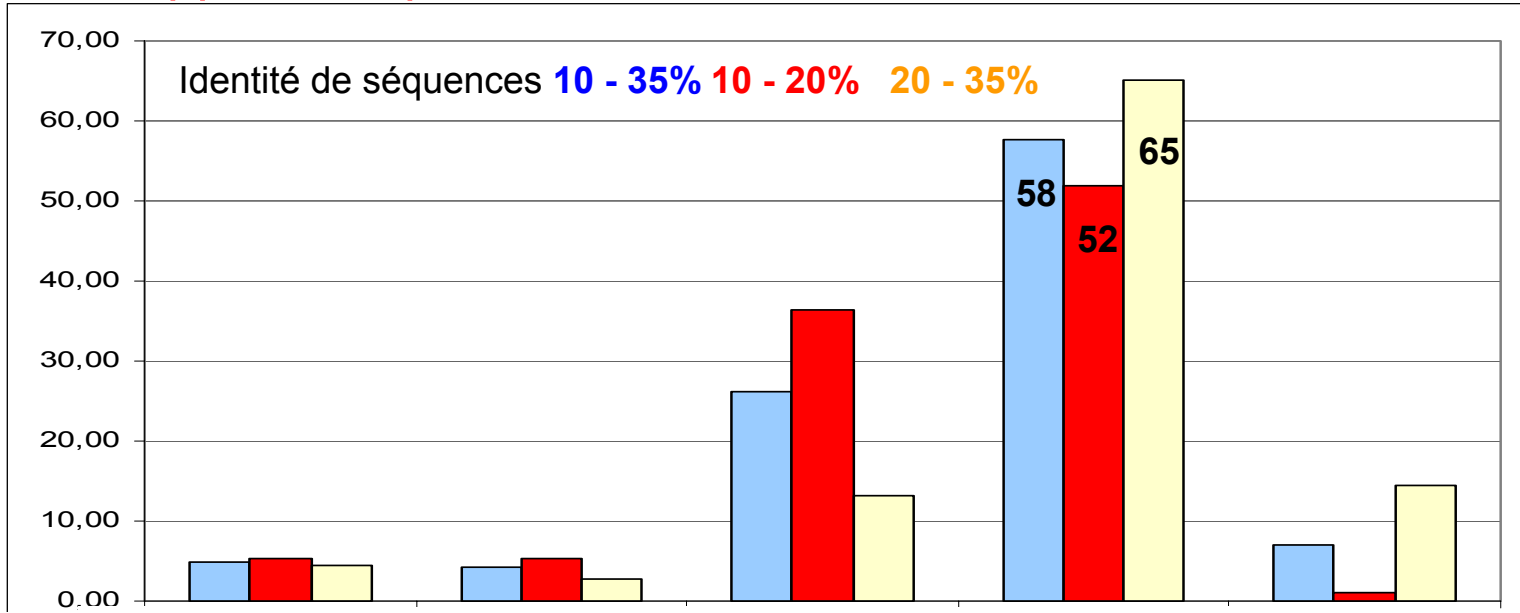




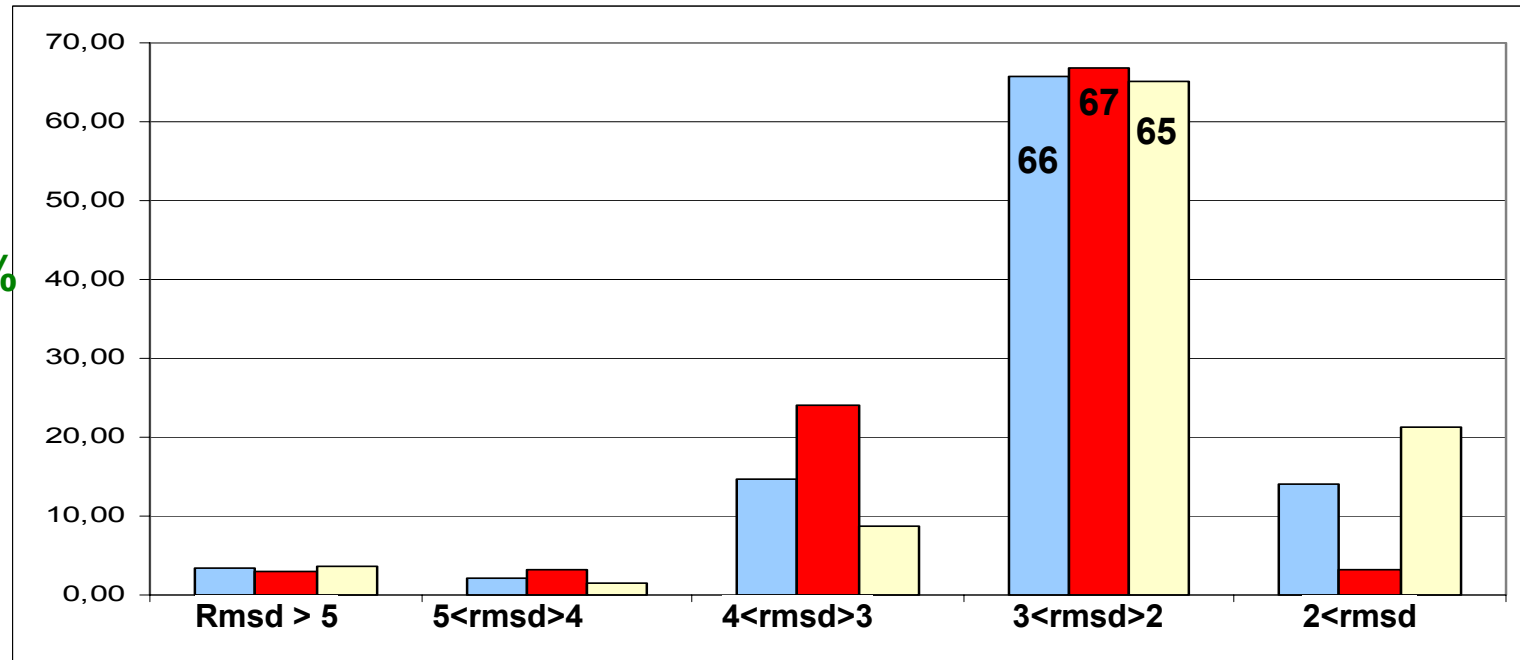
# Apport des prédictions de structures secondaires



Sans Sov



Avec Sov >60%

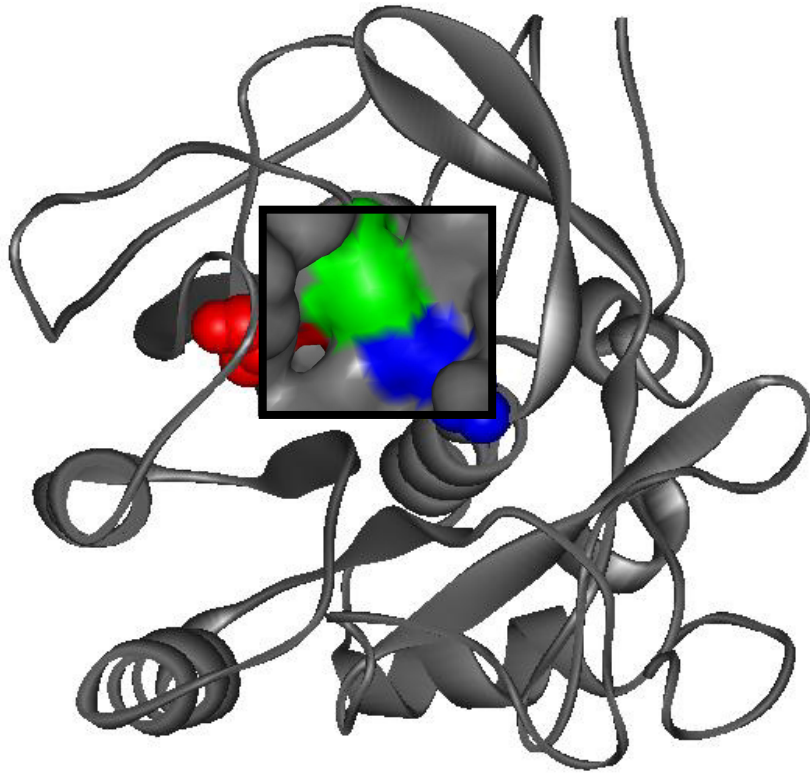




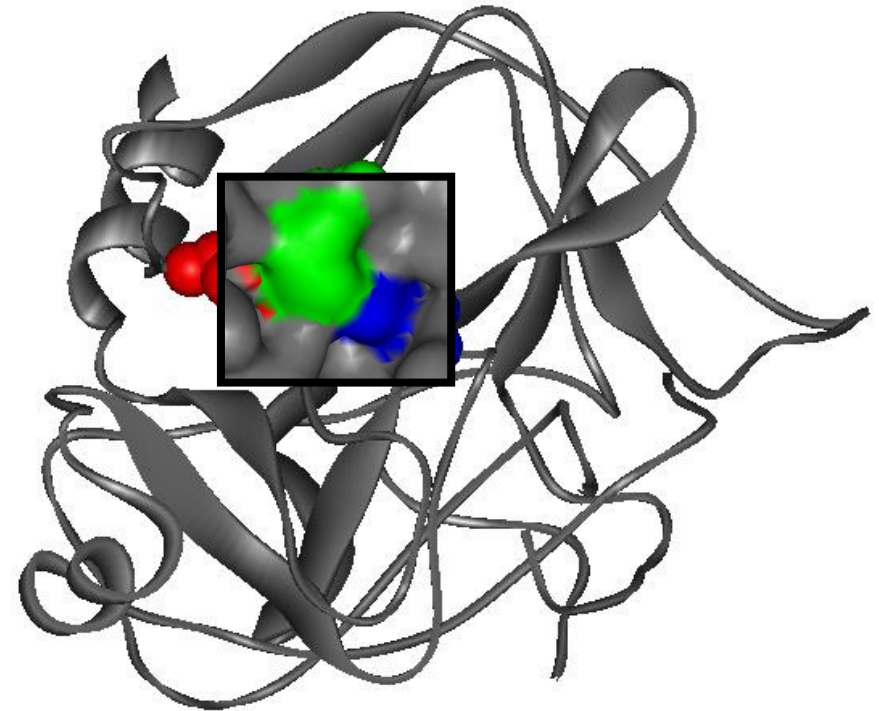
- Bonne fiabilité y compris à bas taux d'identité
- Ne pas introduire d'*a priori* au niveau des «gaps» et insertions dans l'alignement.
- Possibilité d'avoir une estimation de la qualité du modèle obtenu.
- Possibilité de combiner des données provenant de protéines proches (séquences ou fonctions) mais aussi des données expérimentales.
- Possibilité de modéliser des dimères ou trimères
- Possibilité de modéliser des protéines par domaines
- Possibilité d'inclure les ligands (géométriquement puis minimisation)
- Disponible également sur serveur sécurisé Contrat de collaboration avec des industriels

Utilisation dans le cadre de modélisation moléculaire à grande échelle  
Protéome complet *Arabidopsis thaliana*. Programme GENOPLANTE

Problème 1 : séquences non alignables, repliement différent,  
même fonction, même site actif



**Subtilisine**



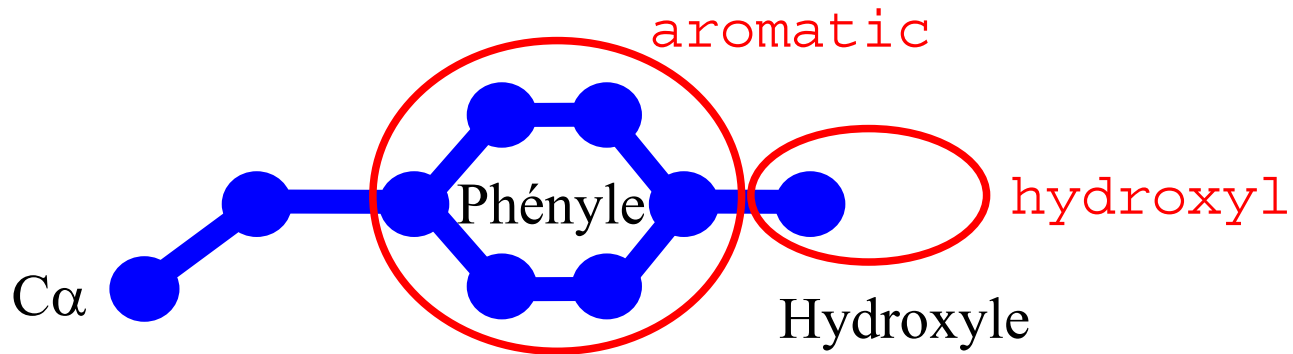
**Chymotrypsine**

*Protéases à sérine*





*Etape 1 : découpage en groupements chimiques*

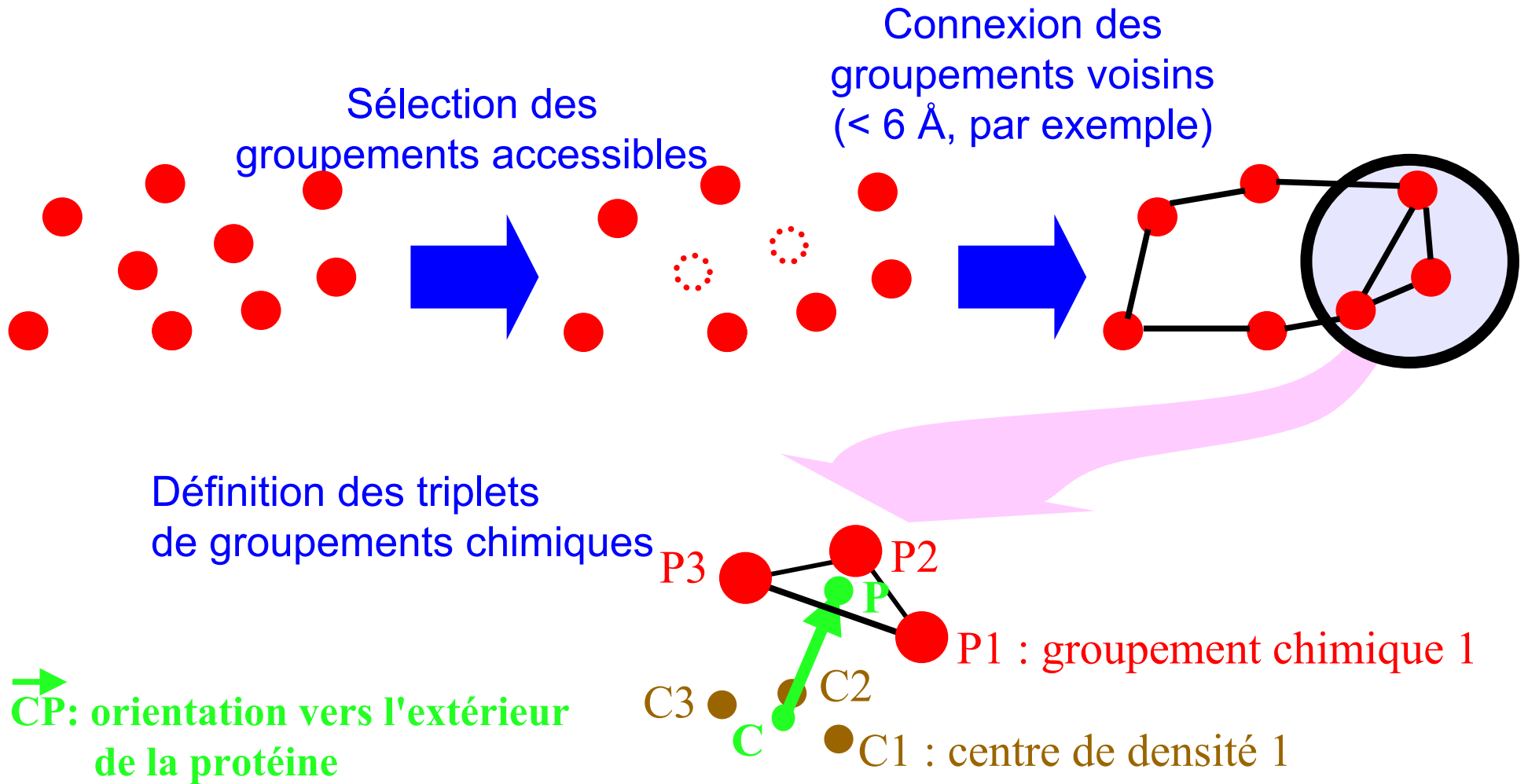


Exemple : tyrosine = aromatic,  
hydroxyl

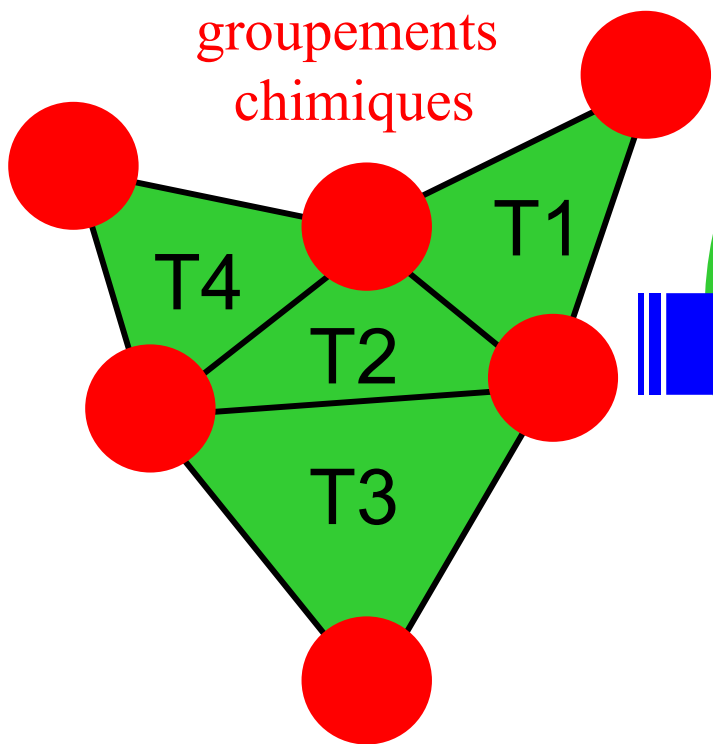
(Jambon *et al.*, **Proteins**, 2003, Brevet international)



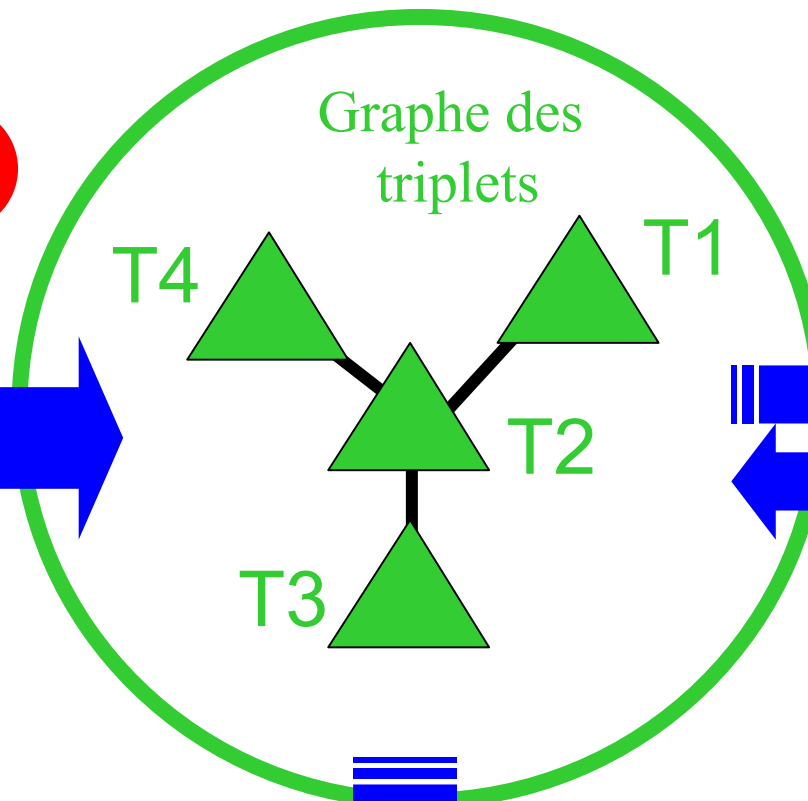
# Etape 2 : génération d'un graphe de triplets de groupements chimiques



Graphe des groupements chimiques

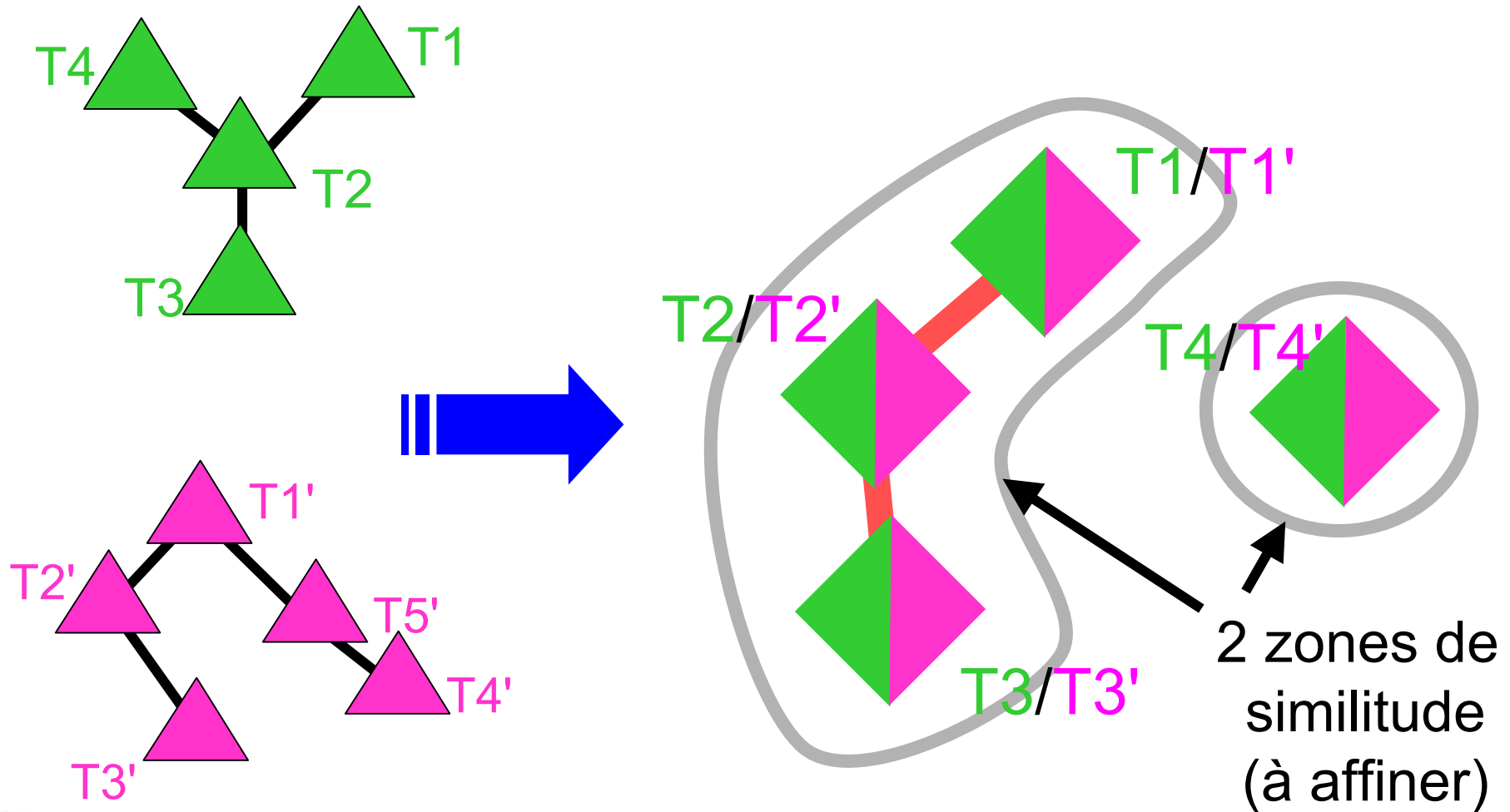


Graphe des triplets



Base de données

Utilisation  
(comparaisons)







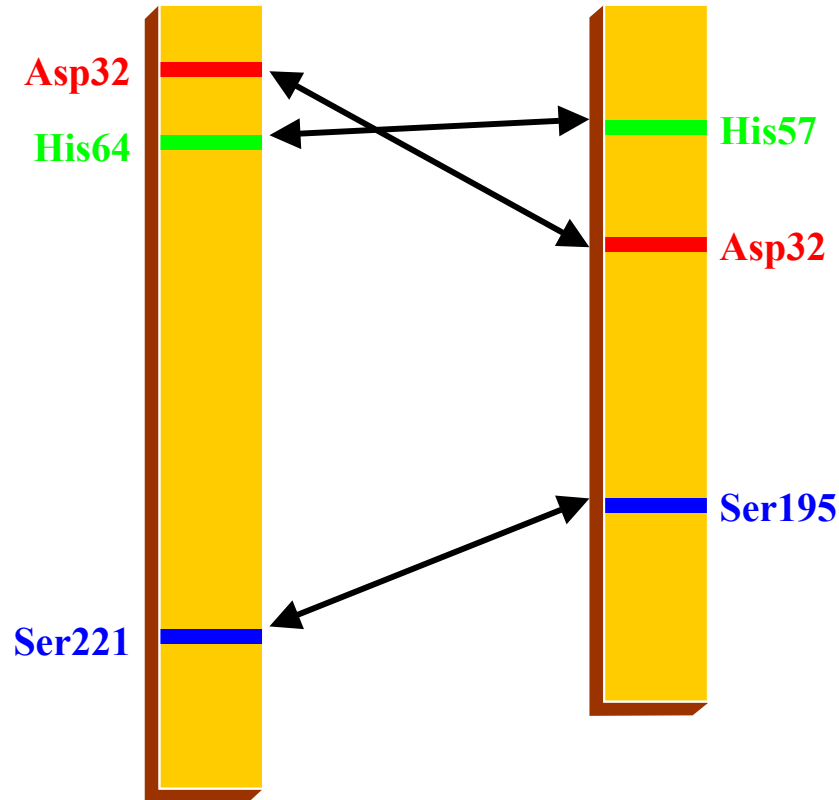
# Exemple 1 : protéases à sérine



## Comparaison subtilisine 1SBC / chymotrypsine 1AFQ



Résultat :



Données biochimiques

```
[Comparison result of (1AFQ,1SBC):
----- Patch number 1 ----...
Number of groups: 4
Score = 1.116 [RMSD = 0.708] [penalty = 0.408]

Selected pairs of groups:
ammonium B HIS 57 | 0.55 | 30.245 2.96...
ammonium HIS 64 | 0.54 | 30.407 2.70...

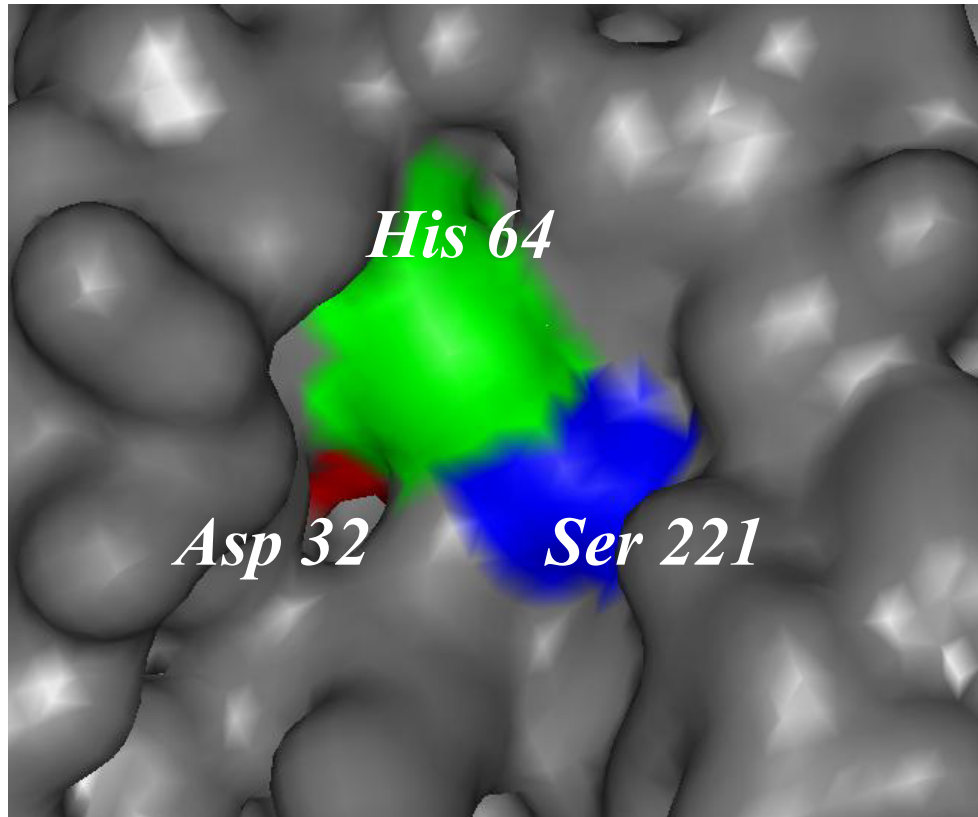
aromatic B HIS 57 | 0.56 | 31.077 3.61...
aromatic HIS 64 | 0.56 | 31.309 3.27...

acyl B ASP 102 | 0.67 | 32.459 7.49...
acyl ASP 32 | 0.66 | 31.730 7.17...

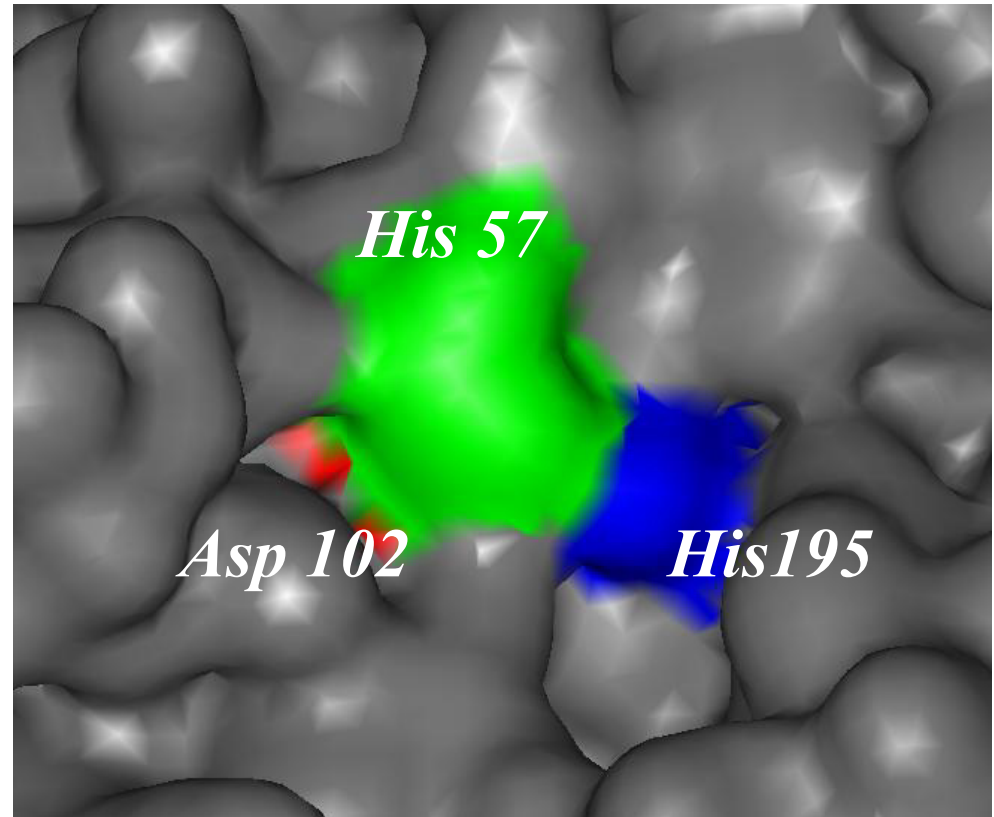
hydroxyl C SER 195 | 0.59 | 27.628 2.51...
hydroxyl SER 221 | 0.57 | 27.964 3.42...
```



Subtilisine (1SBC)



Chymotrypsine (1AFQ)





## Exemple 2 : lectines de légumineuses

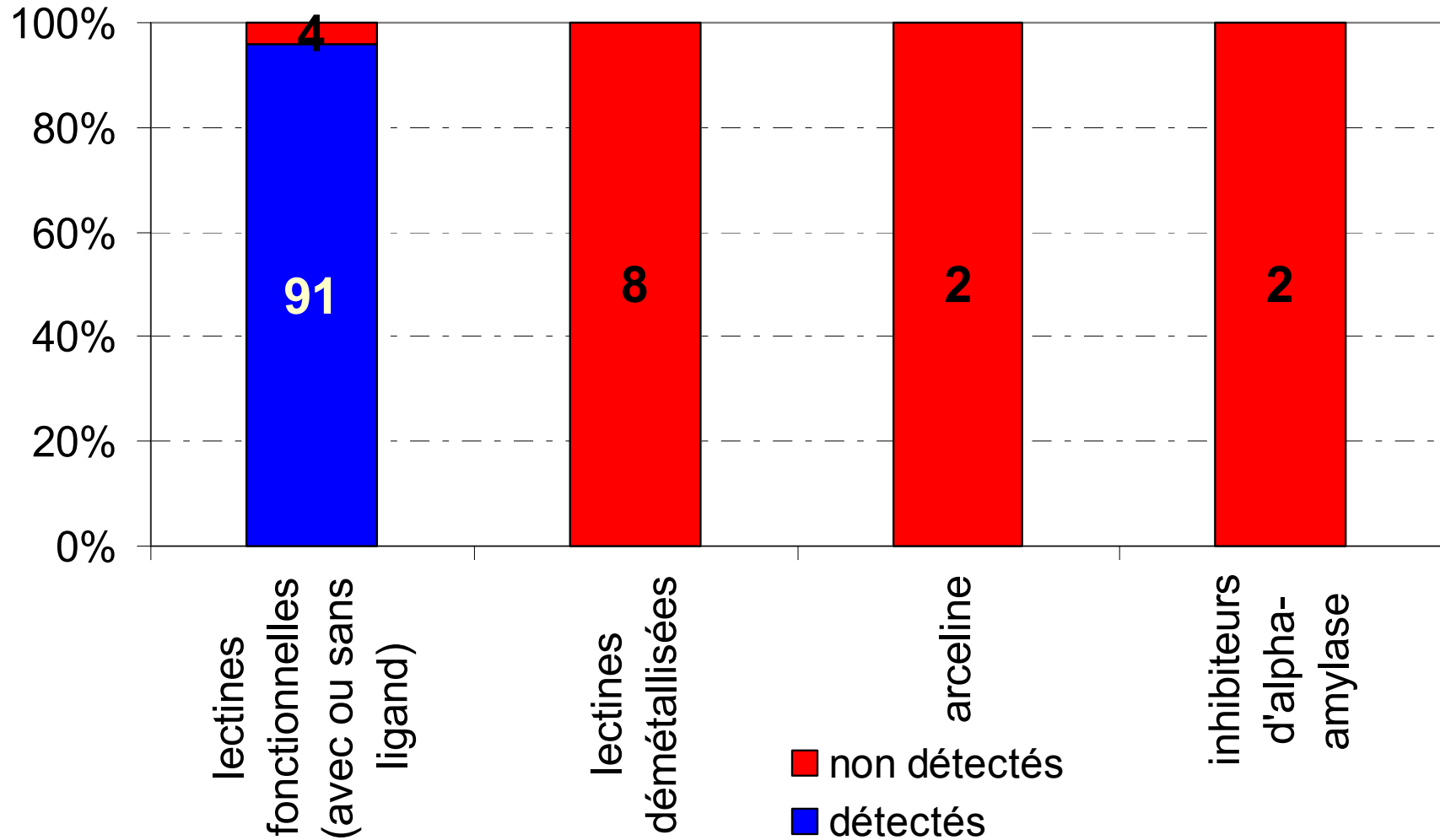


*Collaboration avec Anne Imberty, CERMAV (Grenoble)*

- Travail sur la famille des lectines de légumineuses
- Séquences voisines
- Même repliement



# Recherche de site à sucre dans la famille des lectines de légumineuses



**Site recherché : lectine de cacahuète (2PEL)**

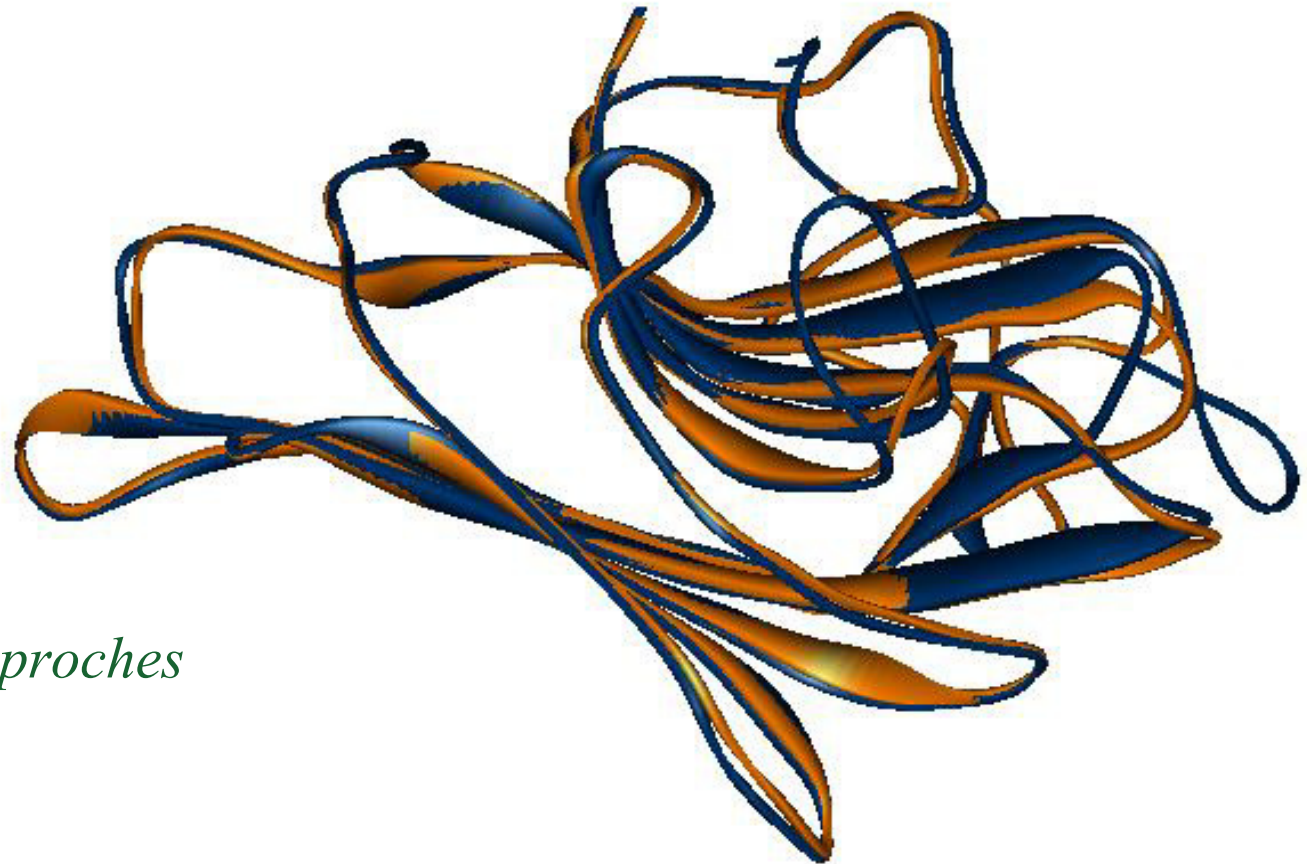
**Nombre de protéines analysées = 107**

**Pourcentage de réussite globale = 96 %**



Orange : concanavaline A fonctionnelle

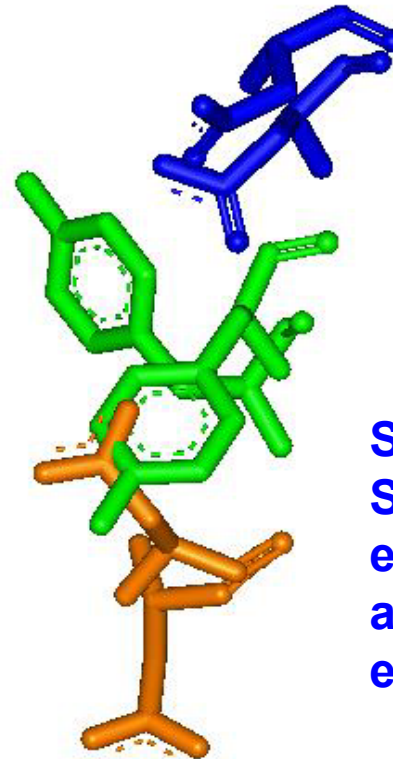
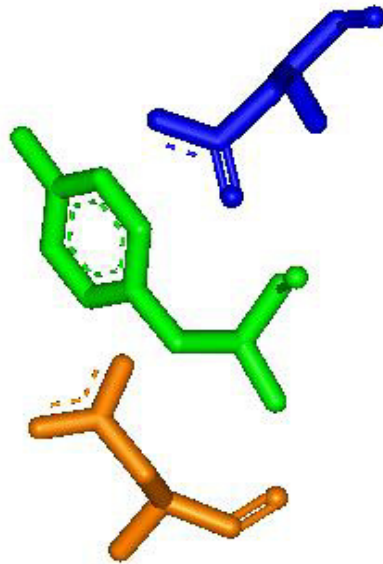
Bleu : concanavaline A non fonctionnelle (démétallisée)





- $RMSD = 0,9 \text{ \AA}$
- *séquences très proches*

# Vue des sites 1DQ1 et 1DQ2 : concanavaline A native et démétaillée

Site de fixation du sucre (mannose)



**Superposition :  
Site fonctionnel  
et site détruit en  
absence de calcium  
et de zinc**

**General Information about the Entry**

Entry name: SWISSPROT:YTFR\_ECOLI

Prim. accession #: P39326

Sec. accession #: P39327

Created: Release 31, 1-FEB-1995

Last sequence update: Release 40, 16-OCT-2001

Last annotation update: Release 41, 28-FEB-2003

**Description and Origin of the Protein**

Keywords: Hypothetical protein; ATP-binding; Transp

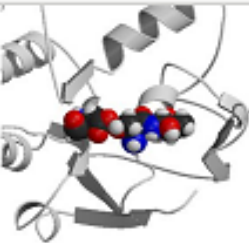

Description: Hypothetical abc transporter atp-binding p

	YTFR_ECOLI	1J10 - ATP	Deformation (+ coef.)	Deviation (Å)	Depth difference	Weight
delta_minus	SER 35	THR 45	1.48 (4.97)	0.175	0.014	0.6
delta_plus	backbone SER 35	backbone THR 45	1.94 (4.59)	0.163	0.003	0.6
hydroxyl	SER 35	THR 45	1.48 (5.37)	0.175	0.014	0.65
delta_plus	backbone THR 36	backbone THR 46	1.48 (4.56)	0.143	0.011	0.6
hydroxyl	THR 36	THR 46	1.18 (3.69)	0.148	0.014	0.65
delta_minus	d1 ASP 156	d1 ASP 163	2.58 (3.66)	0.300	0.020	0.6

[RasMol] | [PDB] | [PDB file] | [RasMol] | [PDB] | [PDBsum]

YTFR\_ECOLI | 1J10 - ATP

TRANSPORT PROTEIN - CRYSTAL STRUCTURE ANALYSIS OF THE ABC TRANSPORTER FROM THERMOTOGA MARITIMA

**Pôle Bio-Informatique L. Geno3D**

Geno3D is the IBCP contribution to PBIIL in Lyon

[HOME](#) [GENO3D](#) [HELP](#) [REFERENCES](#) [NEWS](#)

Tuesday, February 17th 2003 - GENO3D is now running on Linux th

**FIRST STEP :**  
Select template(s) to use for each chain in one or more pdb target :


**PSI-BLAST run 1 for 'YTFR\_ECOLI'**

TEMPLATE	E	FIRST	LAST	ID	ALIGNEMENT	COMMENT
<input type="checkbox"/> p0b121A-1	1e-06	230	477	24	<a href="#">see alignment</a>	TRANSPORT PROTEIN
<input checked="" type="checkbox"/> p0b1b0uA-0	3e-13	15	221	27	<a href="#">see alignment</a>	TRANSPORT PROTEIN
<input type="checkbox"/> p0b1b0uA-1	1e-07	268	462	22	<a href="#">see alignment</a>	TRANSPORT PROTEIN
<input type="checkbox"/> p0b1g7A-0	1e-11	20	222	28	<a href="#">see alignment</a>	PROTEIN TRANSPORT
<input type="checkbox"/> p0b1g7A-1	2e-08	283	475	28	<a href="#">see alignment</a>	PROTEIN TRANSPORT
<input type="checkbox"/> p0b1f7vC-0	4e-08	32	233	24	<a href="#">see alignment</a>	TRANSPORT PROTEIN/HYDROLASE
<input type="checkbox"/> p0b1f7vC-1	2e-05	403	492	30	<a href="#">see alignment</a>	TRANSPORT PROTEIN/HYDROLASE

- Sequence of this chain :  
LSKFFPGVKA LGGWDFSLRR GEIMALLGEM GAGSTLIKA LTGVYEADRG  
TMIKGGQIIS PMNTAHAGQL QIGTVYQGEVM LPHNHWAVM LPIQREPKF  
GLLESPHEEK PATELHASYG FSLQVREPLN RFSVAPQQIV AICRAIDLSA  
KVLILDEPTA SLDTQVELL FDLNRQLDR GVSLIFVTH LIQVYQSDR  
ITVLSNG

- This chain was modelled using 1 templates :

Template	Alignment	Secondary information (Sov)
p0b1b0uA_0	all_aligned	all_aligned 77,24





# Etude structure fonction des protéines. Génomique structurale

