

## PART II

### PROTEIN SEQUENCE ANALYSIS



Seq1:

meepqsdlsi elplsgetfs pksakralpt ntssspppkk meepqsdlsi elplsgetfs  
pksakralpt ntssspppkk meepqsdlsi elplsgetfs pksakralpt pksakralpt

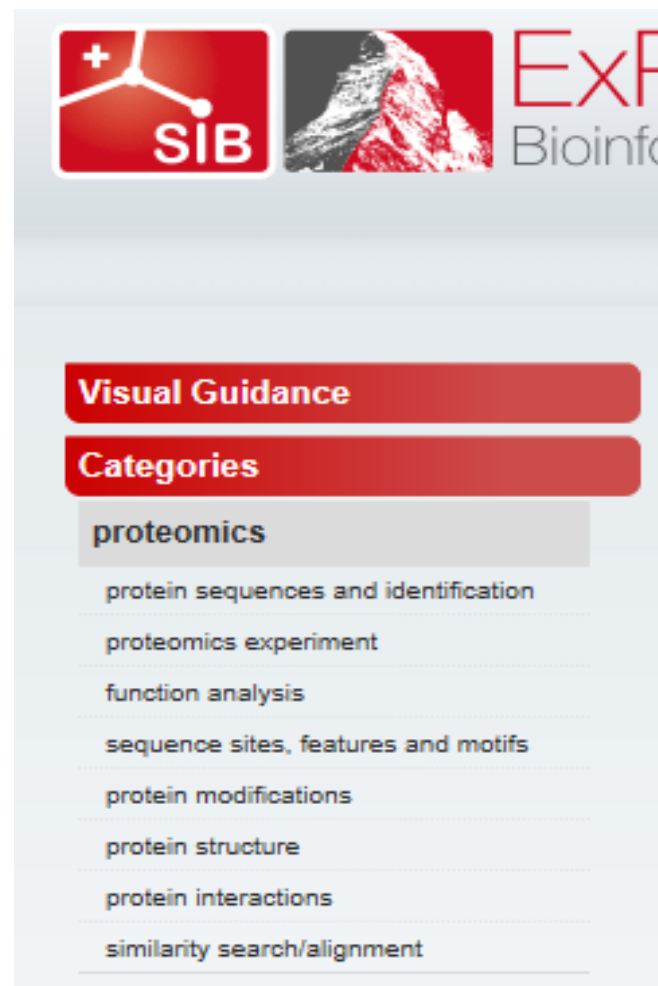
想…想象…

你想知道啥？



# SEQUENCE ANALYSIS OF PROTEINS

- Back-translation
- Molecular weights, pIs
- Amino acid composition
- Hydropathy profile



# 1. BACK—TRANSLATION

- Protein → DNA
- Use for cloning protein of interest where it may be present in low amount.
- Beware of codon bias and degeneracy of codons.



UUU-Phe  
UUC-Phe  
UUA-Leu  
UUG-Leu

UCU-Ser  
UCC-Ser  
UCA-Ser  
UCG-Ser

UAU-Tyr  
UAC-Tyr  
UAA-Stop  
UAG-Stop

UGU-Cys  
UGC-Cys  
UGA-Stop  
UGG-Trp

CUU-Leu  
CUC-Leu  
CUA-Leu  
CUG-Leu

CCU-Pro  
CCC-Pro  
CCA-Pro  
CCG-Pro

CAU-His  
CAC-His  
CAA-Gln  
CAG-Gln

CGU-Arg  
CGC-Arg  
CGA-Arg  
CGG-Arg

AUU-Ile  
AUC-Ile  
AUA-Ile  
AUG-Met

ACU-Thr  
ACC-Thr  
ACA-Thr  
ACG-Thr

AAU-Asn  
AAC-Asn  
AAA-Lys  
AAG-Lys

AGU-Ser  
AGC-Ser  
AGA-Arg  
AGG-Arg

GUU-Val  
GUC-Val  
GUA-Val  
GUG-Val

GCU-Ala  
GCC-Ala  
GCA-Ala  
GCG-Ala

GAU-Asp  
GAC-Asp  
GAA-Glu  
GAG-Glu

GGU-Gly  
GGC-Gly  
GGA-Gly  
GGG-Gly



# Biased codon usage

Amino acid	Codon	Bacteria	Yeast	Fruit Fly	Human
Leu	UUA				
	UUG		Preferred		
	CUU				
	CUC				
	CUA				
	CUG	Preferred		Preferred	Preferred
Val	GUU	Preferred	Preferred		
	GUC				
	GUA				
	GUG			Preferred	Preferred



# Back-translation Tool

## Sequence Manipulation Suite:

### Reverse Translate

Reverse Translate accepts a protein sequence as input and uses a codon usage table to generate a DNA sequence representing the most likely non-degenerate coding sequence. A consensus sequence derived from all the possible codons for each amino acid is also returned. Use Reverse Translate when designing PCR primers to anneal to an unsequenced coding sequence from a related species.

Paste the raw sequence or one or more FASTA sequences into the text area below. Input limit is 20,000,000 characters.

```
>sample sequence
ACDEFGHIKLMNPQRSTVWY*
```

Enter the codon table you wish to use (in GCG format). The default codon usage table was generated using all the E. coli coding sequences in GenBank. It was obtained from the [Codon Usage Database](#).

AmAcid	Codon	Number	/1000	Fraction	..
Gly	GGG	50527.00	11.12	0.15	
Gly	GGA	39036.00	8.59	0.12	
Gly	GGT	114185.00	25.14	0.34	
Gly	GGC	130043.00	28.63	0.39	

\*This page requires JavaScript. See [browser compatibility](#).

\*You can [mirror this page](#) or [use it off-line](#).

[new window](#) | [home](#) | [citation](#)



## 2. MOLECULAR WEIGHTS, PIs

- Aid in designing of purification experiments e.g. SDS-PAGE, IEF, 2-Dimensional Gel, Column chromatography etc.



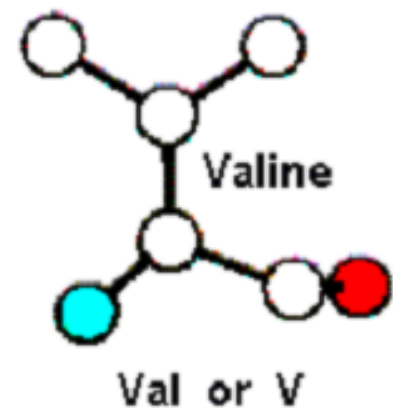
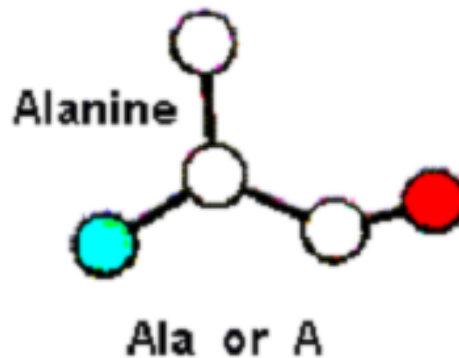
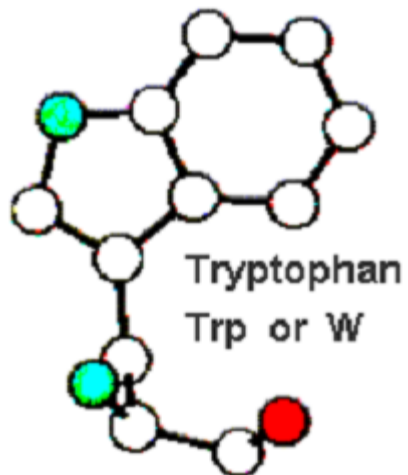
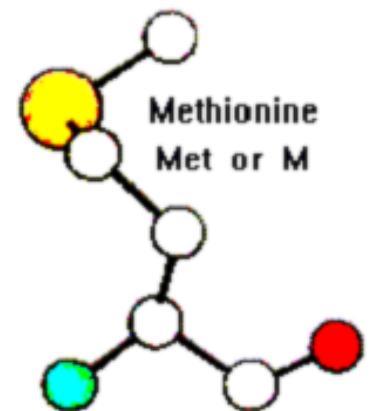
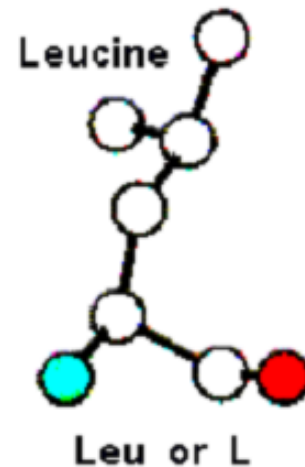
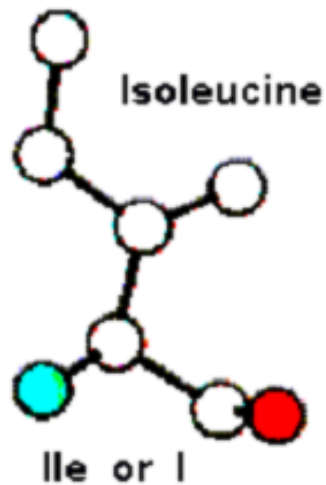
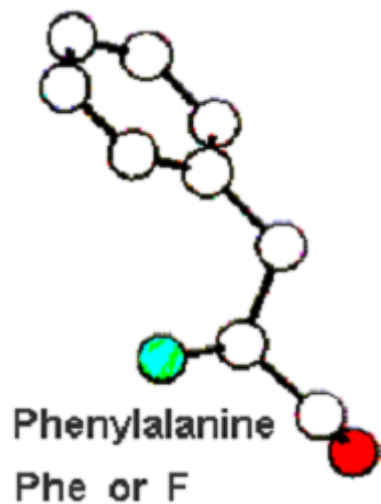


### 3. AMINO ACID COMPOSITION

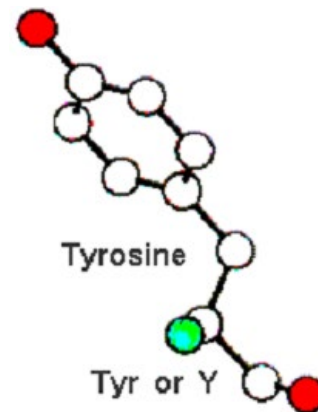
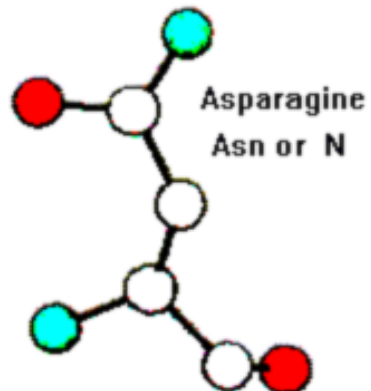
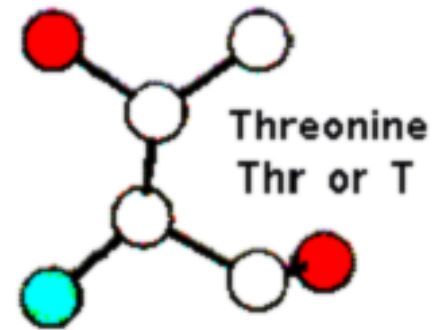
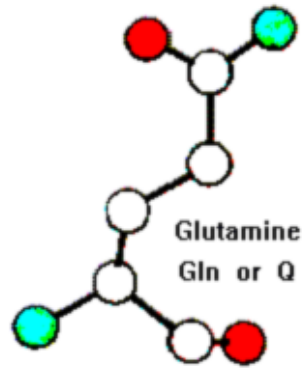
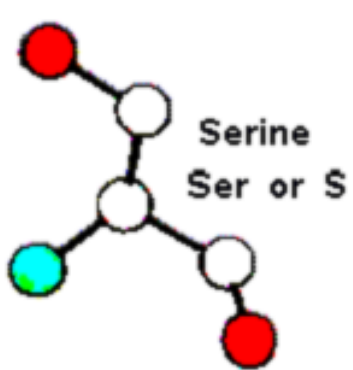
- Determine the percentages of amino acid residues present in a protein molecule.
- Uses:
  - determine the lifestyles of organisms: high percentages of Glu (−) and both Lys and Arginine (+) in hyperthermophiles vs. mesophiles → absent (*Tekalia et al.*, 2002).
  - predict structural class (Luo *et al.*, 2002).



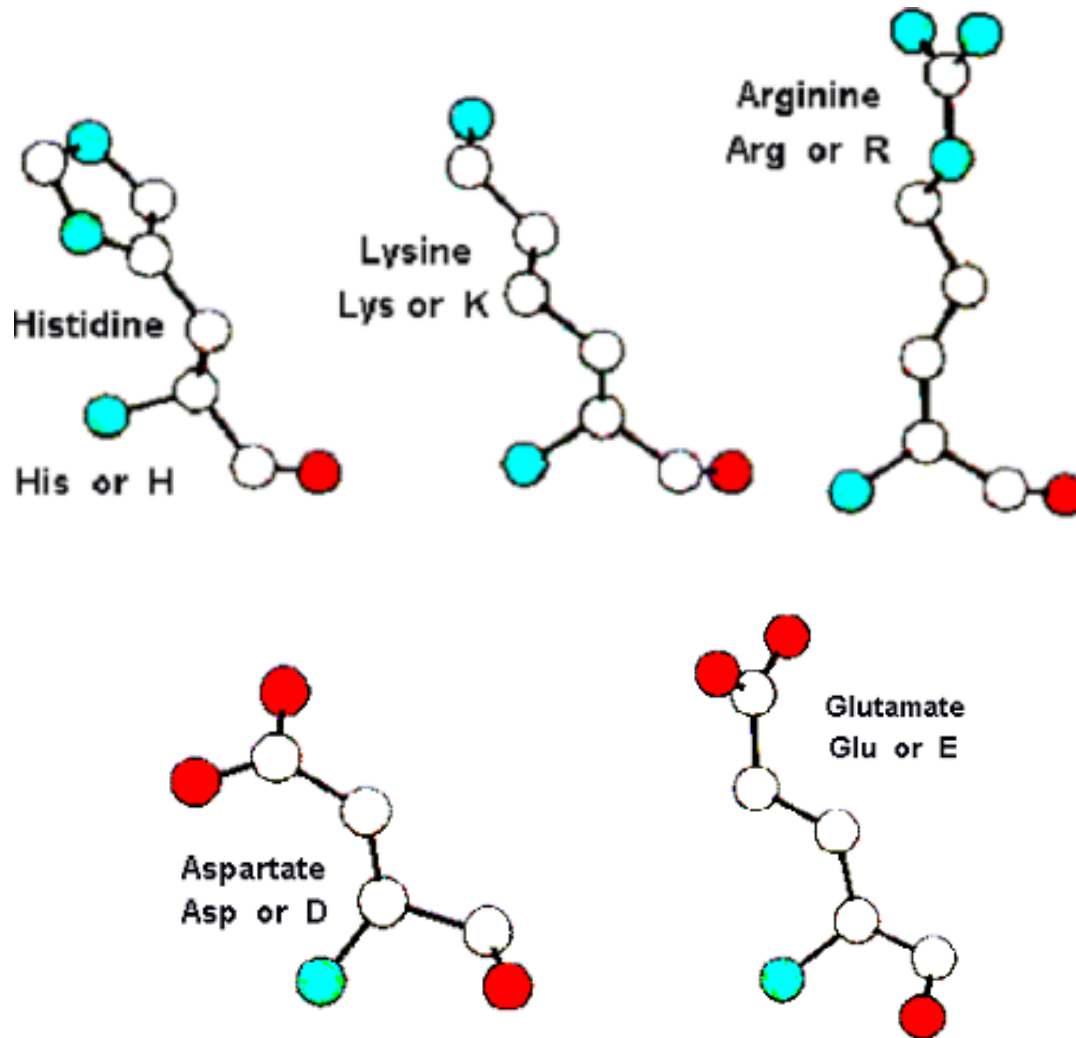
# NONPOLAR AMINO ACIDS (FILMWAV)



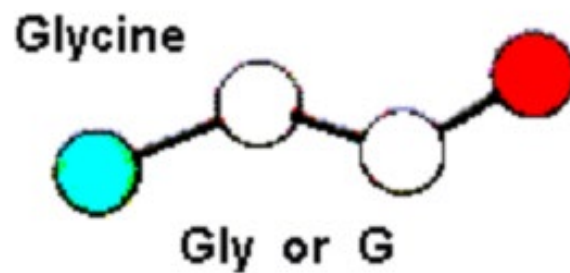
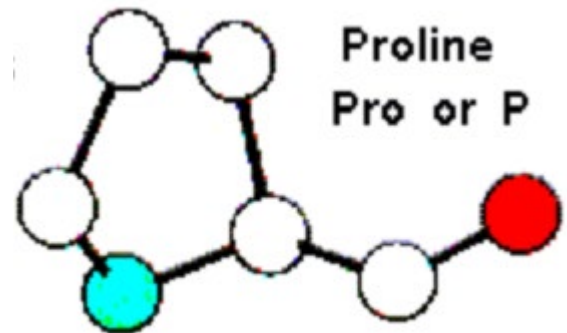
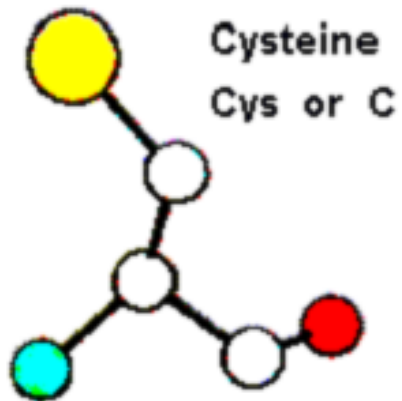
## POLAR UNCHARGED ( $S^-Q^+T^-N^+Y^-$ )



## POLAR CHARGED (KHERD)



# UNIQUE PROPERTIES



# Protein functions from specific residues

- C Disulphide-rich, zinc fingers
- DE Acidic proteins (unknown)
- G Collagens
- H Histidine-rich glycoprotein
- KR Nuclear proteins, nuclear localisation
- P Collagen, filaments
- ST Mucins (high molecular weight glycosylated proteins that form a major part of a protective biofilm on the surface of epithelial cells)

# Protein functions from specific residues

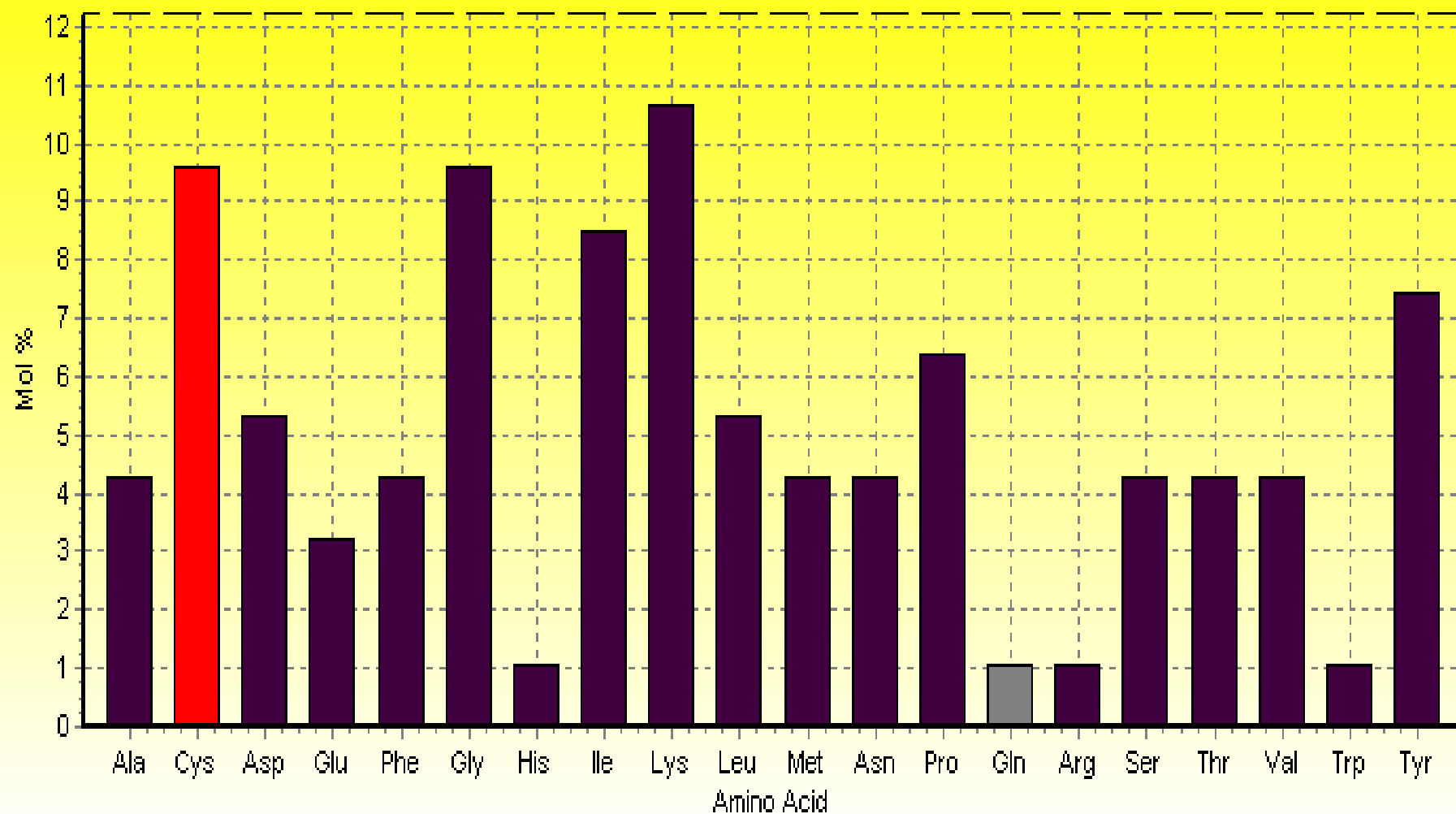
- Polar (C, D, E, H, K, N, Q, R, S, T) – active sites
- Aromatic (F, H, W, Y) – protein ligand-binding sites
- $\text{Zn}^{2+}$ -coordinates (C, D, E, H, N, Q) – active site, zinc finger
- $\text{Ca}^{2+}$ -coordinates (D, E, N, Q) – ligand-binding site
- Mg/Mn-coordinates (D, E, N, S, R, T) –  $\text{Mg}^{2+}$  or  $\text{Mn}^{2+}$  catalysis, ligand binding
- Phosphate-binding (H, K, R, S, T) – phosphate and sulfate binding

Protein: *Hottentotta judaica*

Length: 94 amino acids

Molecular weight: 10510.93 Daltons

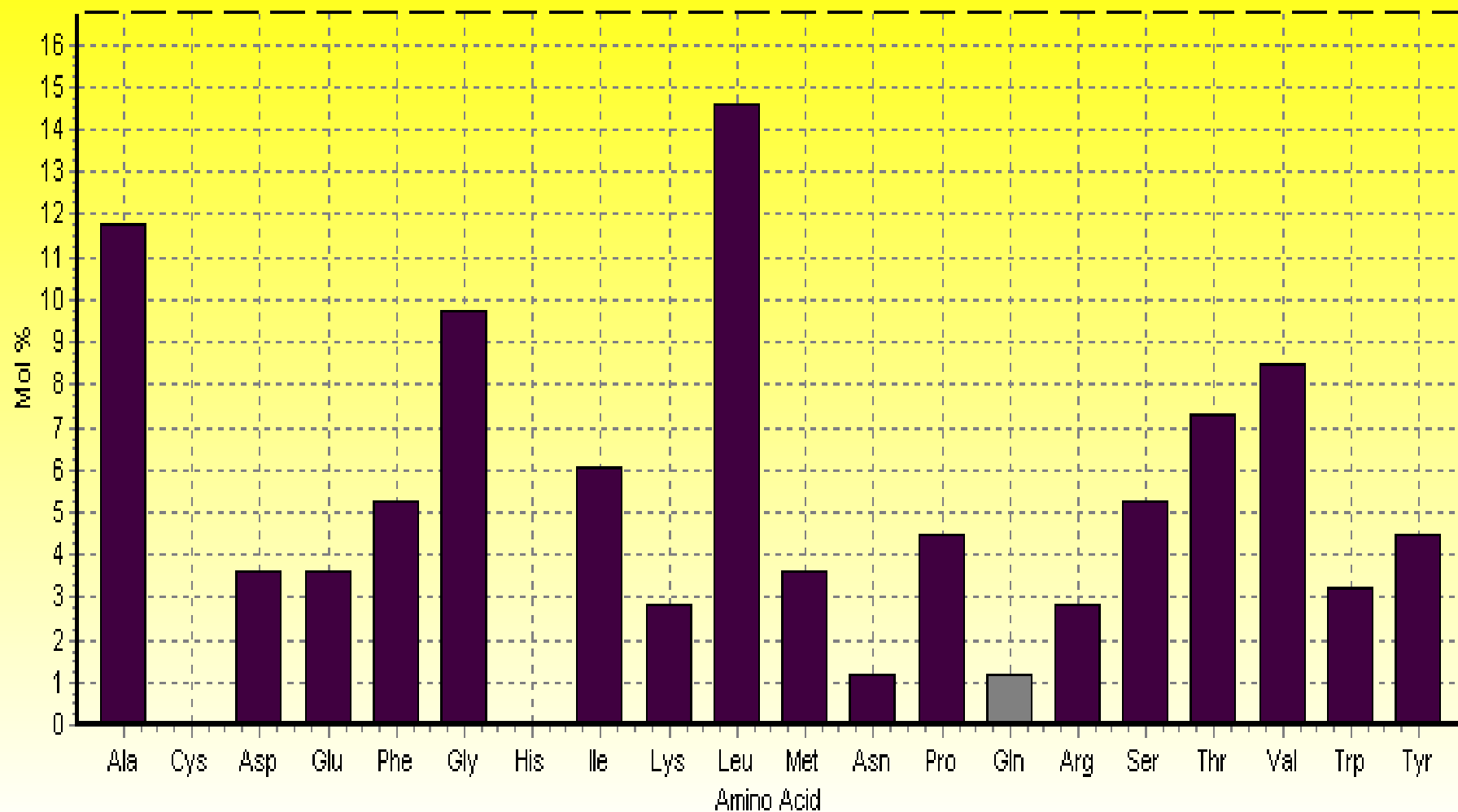
Amino Acid Composition  
*Hottentotta judaica*



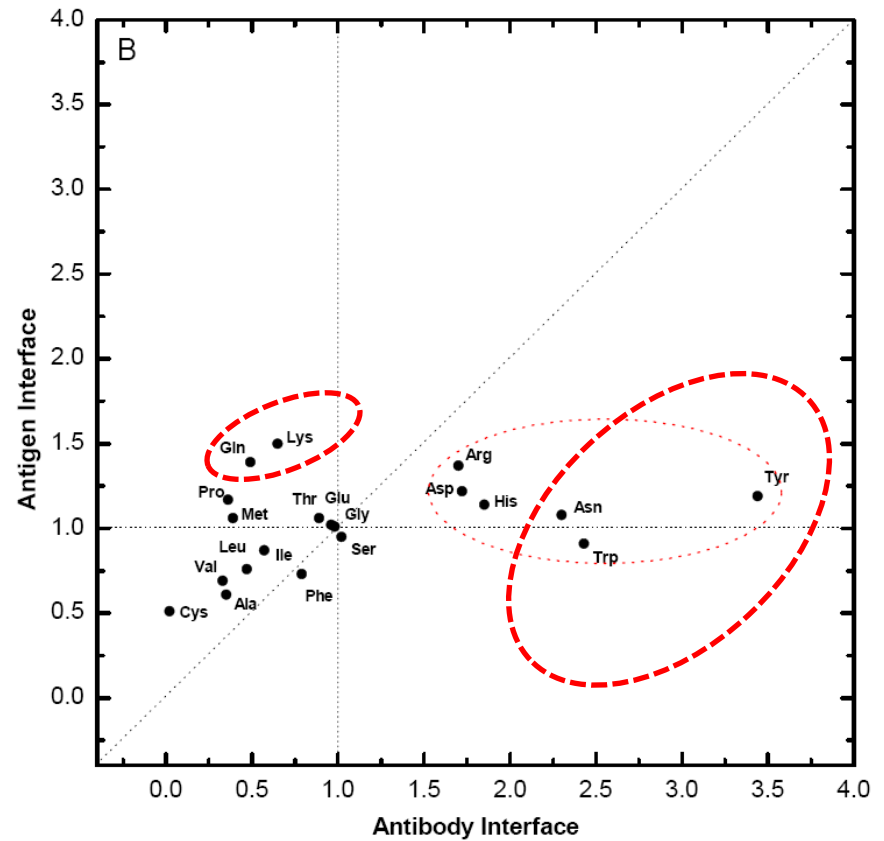
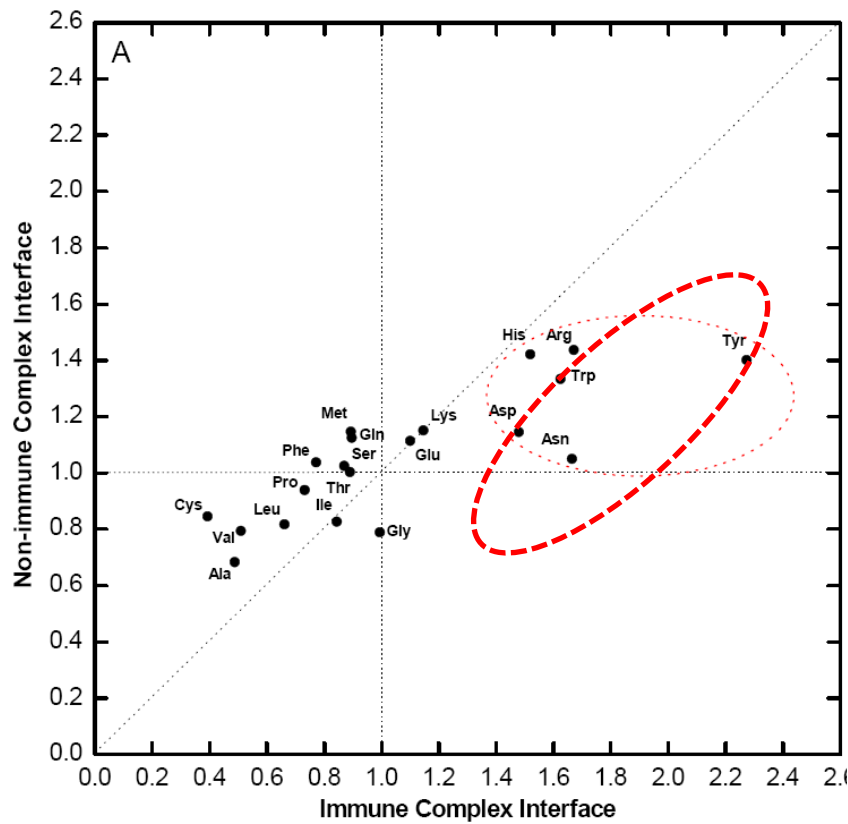


Protein: Halobacterium halobium  
 Length: 247 amino acids  
 Molecular weight: 26723.79 Daltons

# Amino Acid Composition Halobacterium halobium

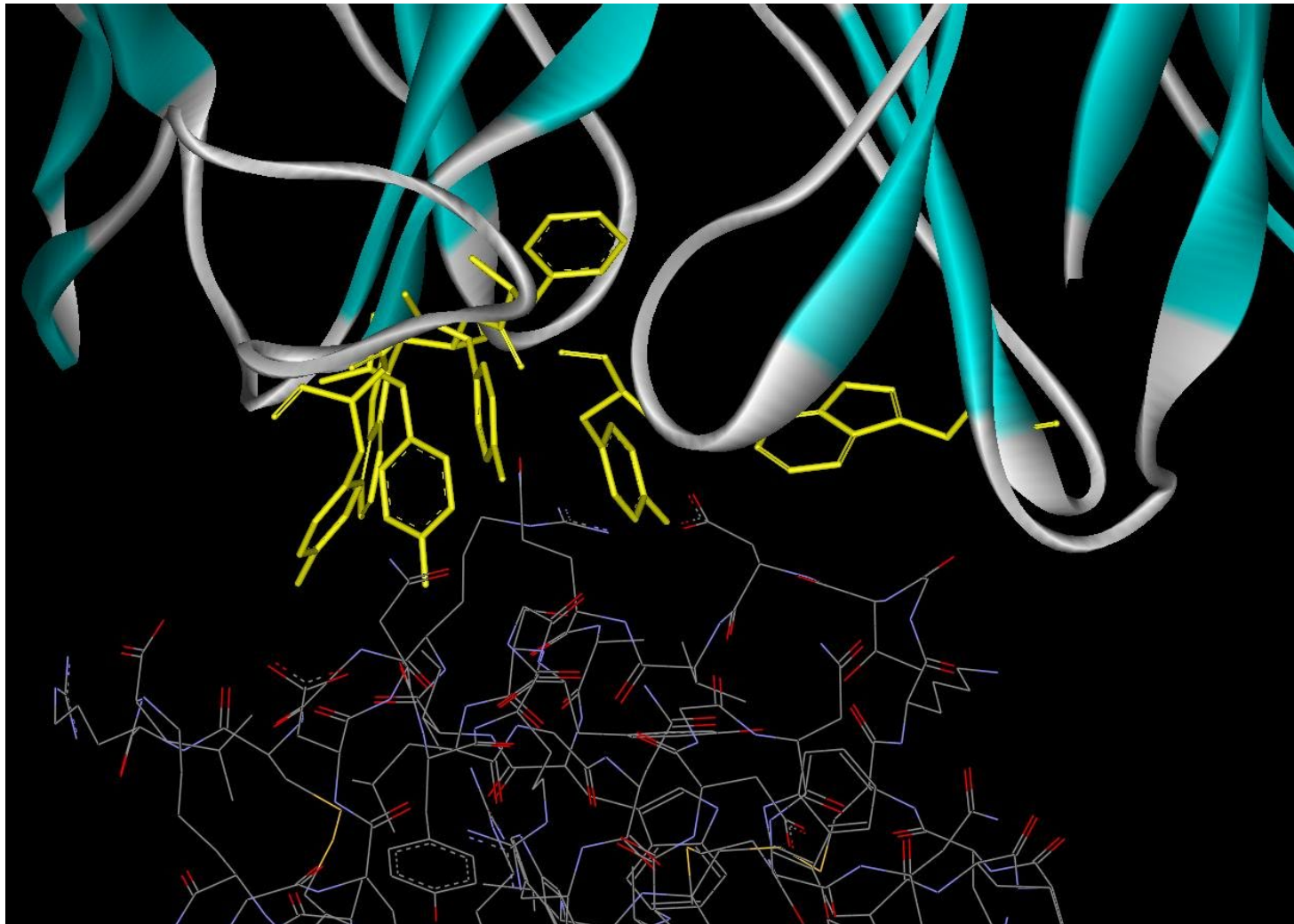


# 作用界面**抗体侧**氨基酸组成规律

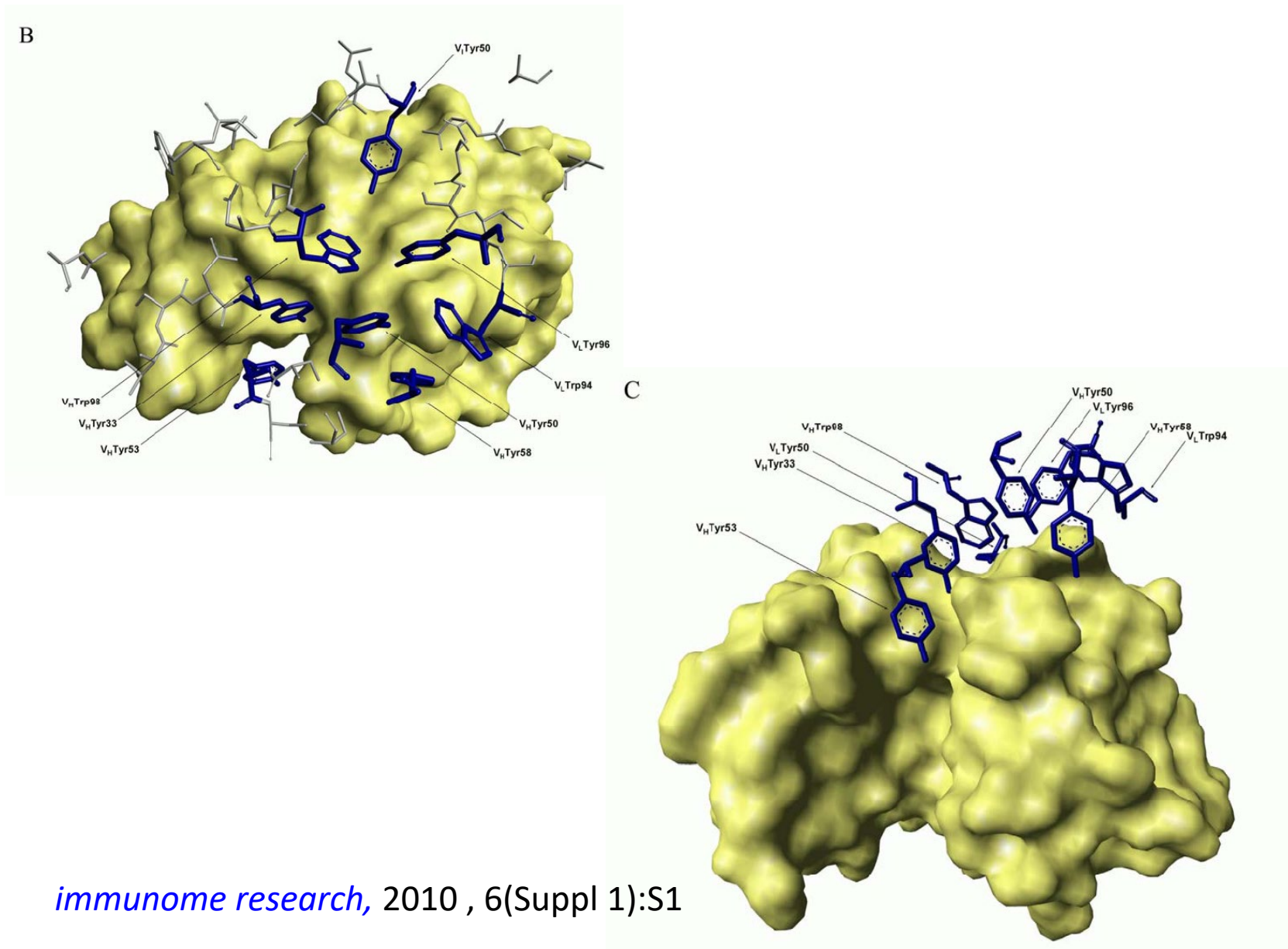




1BVK



# Ab-Ag界面上富集的芳香族氨基酸空间分布特征



## AACompldent tool

**AACompldent** is a tool which allows the identification of a protein from its amino acid composition [\[references\]](#) It searches the [Swiss-Prot](#) and / or [TrEMBL](#) databases for proteins, whose amino acid compositions are closest to the amino acid composition given.

[Documentation](#) is available.

Few amino acid analysis techniques produce composition results for all amino acids. We currently have indexed Swiss-Prot and TrEMBL for the following constellations. Please choose one of them:

1. **Constellation 0: ALL amino acids:** Ala, Ile, Pro, Val, Arg, Leu, Ser, Thr, Gly, Met, His, Phe, Tyr, Lys, Asp, Asn, Gln, Glu, Cys and Trp.
2. **Constellation 1:** Ala, Ile, Pro, Val, Arg, Leu, Ser, Asx, Thr, Glx, Gly, Met, His, Phe and Tyr.  
(Asp+Asn=Asx; Gln+Glu=Glx; Lys, Cys and Trp are not considered).
3. **Constellation 2:** Ala, Ile, Pro, Val, Arg, Leu, Ser, Asx, Lys, Thr, Glx, Gly, Met, His, Phe and Tyr.  
(Asp+Asn=Asx; Gln+Glu=Glx; Cys and Trp are not considered).
4. **Constellation 5:** Ala, Ile, Pro, Val, Arg, Leu, Ser, Asx, Lys, Thr, Glx, Gly, Met, His, Phe, Tyr and Cys.  
(Asp+Asn=Asx; Gln+Glu=Glx; Trp is not considered).

## HYDROPATHY PROFILES

- Hydrophathy – describe the hydrophobicity and hydrophilicity of a protein sequence.
- A graph in which hydrophathy values are calculated within a sliding window and plotted for each residue in a protein sequence.



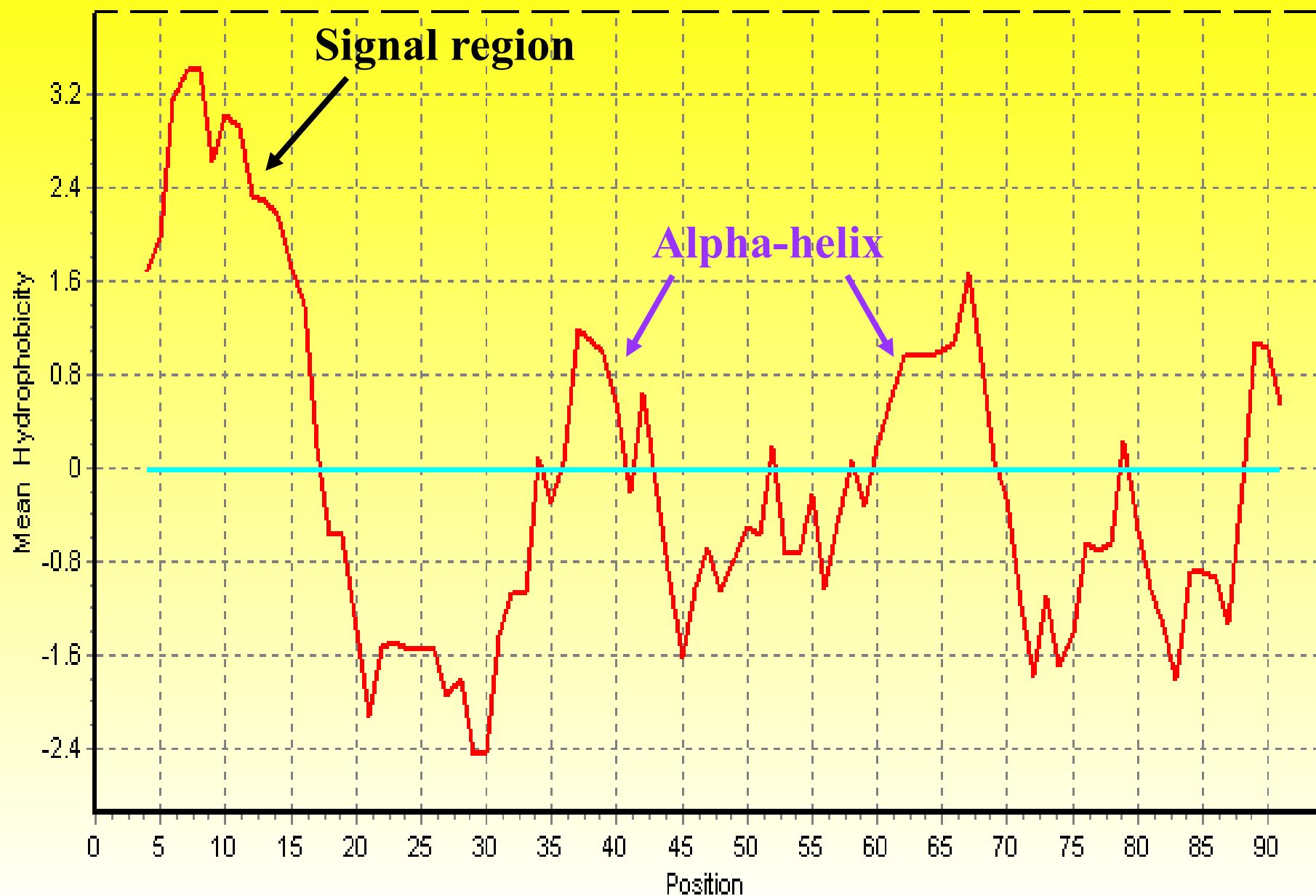
## A SLIDING WINDOW

M	K	F	F	L	M	C	L	I	I	F	P	I	M	G	V	L	G
---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---	---

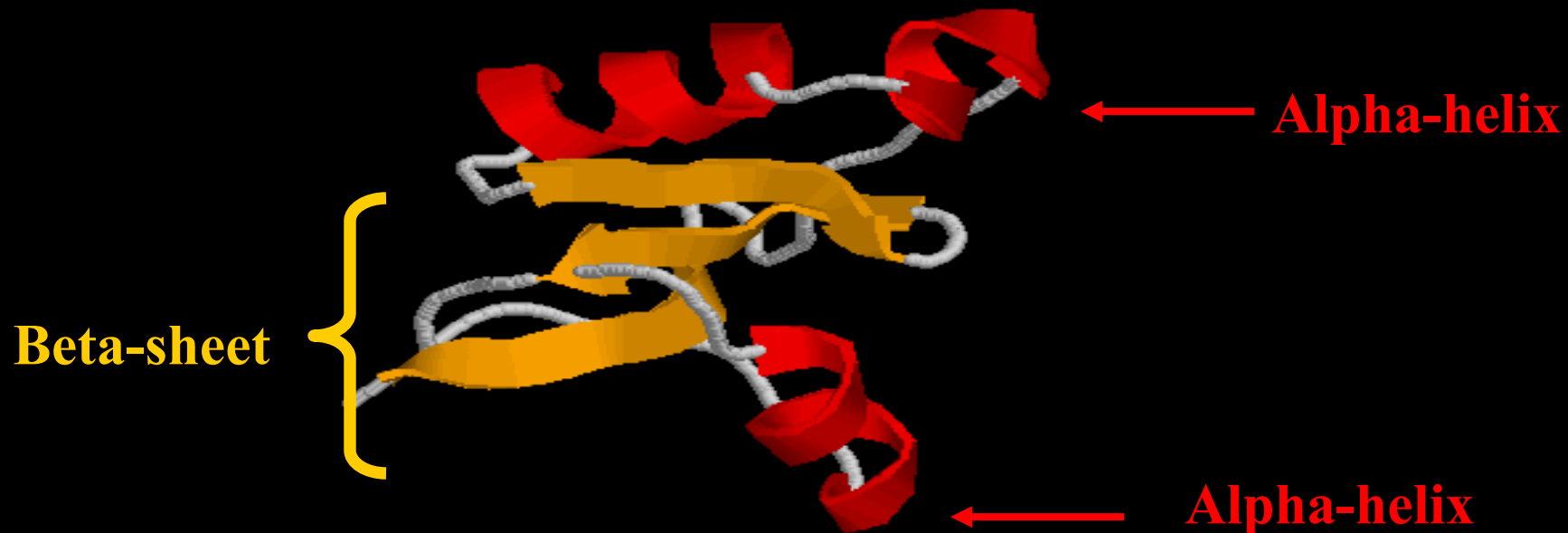




Kyte & Doolittle Scale Mean Hydrophobicity Profile  
Scan-window size = 7



MKFFLMCLIIFPIMGVLGKKNGYPLDRNGKTTECSGVNAIAPHYCNSECT  
KVYYAKSGYCCWGACYCFGLEDDKPIGPMKDITKKYCDVQIIPS  
(Signal region:1-18 Toxin region:19-94  
Disulphide      )



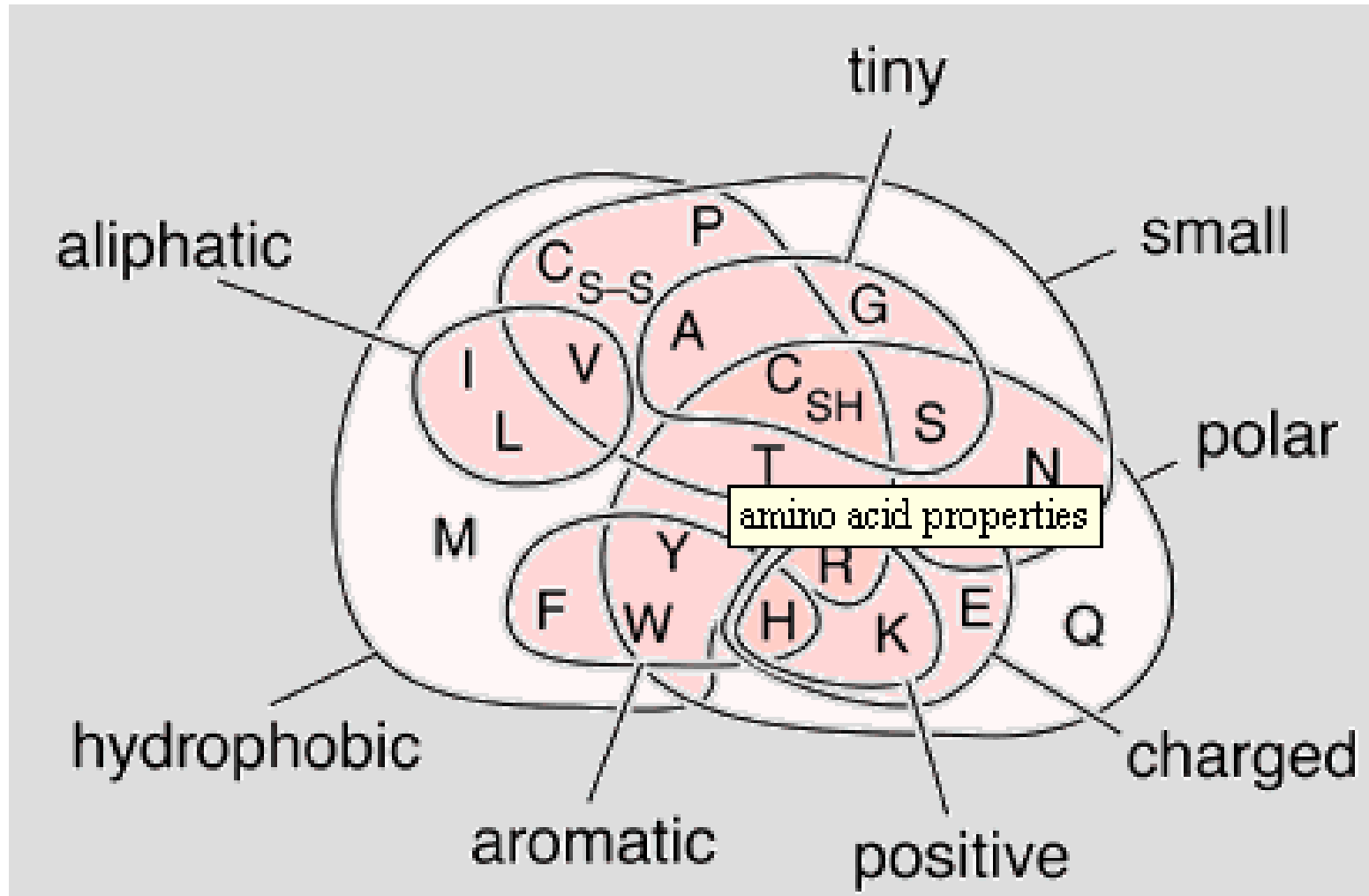
A SCHEMATIC REPRESENTATION OF A 3-D STRUCTURE OF A  
SCORPION TOXIN

# HYDROPATHY PROFILES

- Hydropathy scale – each amino acid is assigned a value reflecting its relative hydrophobicity and hydrophilicity.
- 2 broad classes of scales:
  - Environmental characteristics of protein residues.
  - Experimental measurements of amino acid physiochemical properties.



# VENN DIAGRAM OF THE 20 AMINO ACID PHYSIOCHEMICAL PROPERTIES



# HYDROPATHY PROFILES

- Basic ranking: internal {FILMV}, external {DEHKNQR}, ambivalent {ACGPSTWY}

Scale Name	Residue ranking
Kyte	IVLFCMAGTSWYPHDNEQKR
Eisenberg	IFVLWMAGCYPTSHENQDKR
Cornette	FILVMHYCWAQRTGESPNDK



## HYDROPATHY PROFILES

- Detect possible transmembrane domains (consecutive 20–25 runs of hydrophobic amino acids).
- Hydrophobic protein cores



protein sequence analysis tool

ExPASy: SIB Bioinformatics Resource Portal

+



← → ↺ 🏠

🔒 https://www.expasy.org

⋮ 🛡️ ☆

⬇️ 📄 📖 🗺️ 🗨️ ☰

⚙️ 最常访问 📁 火狐官方网站 🌐 新手上路 📁 常用网址 🇨🇳 JD 京东商城



ExPASy  
Bioinformatics Resource Portal

Home About Contact

Query all databases ▾

search help

Visual Guidance

Categories

proteomics

genomics

structure analysis

systems biology

evolutionary biology

population genetics

transcriptomics

biophysics

imaging

IT infrastructure

medicinal chemistry

glycomics

Resources A..Z


Help

ExPASy is the **SIB Bioinformatics Resource Portal** which provides access to scientific databases and software tools (i.e., *resources*) in different areas of life sciences including proteomics, genomics, phylogeny, systems biology, population genetics, transcriptomics etc. (see **Categories** in the left menu). On this portal you find resources from many different SIB groups as well as external institutions.

Featuring today


Compute pI/MW

Compute isoelectric point and molecular weight from protein sequence  
[\[details\]](#)




FindPept

Identify peptides resulting from unspecific protein cleavage from their experimental masses  
[\[details\]](#)





FindMod


Predict potential protein post-translational modifications and find potential single amino




Popular resources

 UniProtKB

 SWISS-MODEL

 STRING

 PROSITE

Latest News

Protein Spotlight: A way in -  
2020-03-26

As children in Scotland, back home from school and when the weather was dry, we would fling our schoolbags into the hall and grab a few golf clubs, a ball and a tee. There was no need for a change of clothes or shoes, ...[More](#).

New neXtProt application release  
2020-02-28 - 2020-02-28

Modifications to SPARQL queries For

## ProtScale






**ProtScale** [[Reference](#) / [Documentation](#)] allows you to compute and represent the profile produced by any amino acid scale on a selected protein.

An **amino acid scale** is defined by a numerical value assigned to each type of amino acid. The most frequently used scales are the hydrophobicity or hydrophilicity scales and the secondary structure conformational parameters scales, but many other scales exist which are based on different chemical and physical properties of the amino acids. This program provides 57 predefined scales entered from the litera

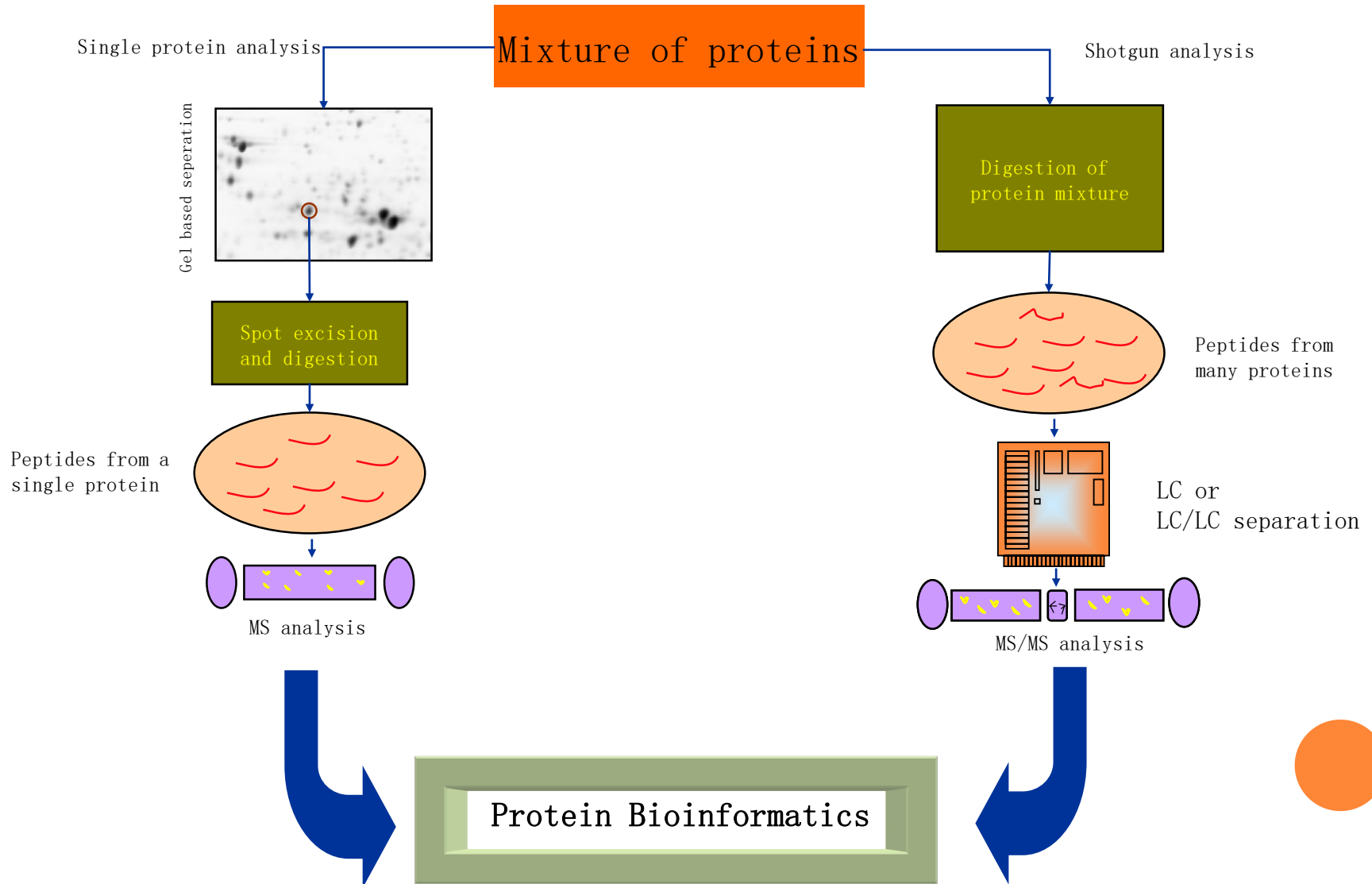
- |  |   |
|--|---|
| <input type="radio"/> Molecular weight                     | <input type="radio"/> Number of codon(s)          |
| <input type="radio"/> Bulkiness                            | <input type="radio"/> Polarity / Zimmerman        |
| <input type="radio"/> Polarity / Grantham                  | <input type="radio"/> Refractivity                |
| <input type="radio"/> Recognition factors                  | <input type="radio"/> Hphob. / Eisenberg et al.   |
| <input type="radio"/> Hphob. OMH / Sweet et al.            | <input type="radio"/> Hphob. / Hopp & Woods       |
| <input checked="" type="radio"/> Hphob. / Kyte & Doolittle | <input type="radio"/> Hphob. / Manavalan et al.   |
| <input type="radio"/> Hphob. / Abraham & Leo               | <input type="radio"/> Hphob. / Black              |
| <input type="radio"/> Hphob. / Bull & Breese               | <input type="radio"/> Hphob. / Fauchere et al.    |
| <input type="radio"/> Hphob. / Guy                         | <input type="radio"/> Hphob. / Janin              |
| <input type="radio"/> Hphob. / Miyazawa et al.             | <input type="radio"/> Hphob. / Rao & Argos        |
| <input type="radio"/> Hphob. / Roseman                     | <input type="radio"/> Hphob. / Tanford            |
| <input type="radio"/> Hphob. / Wolfenden et al.            | <input type="radio"/> Hphob. / Welling & al       |
| <input type="radio"/> Hphob. HPLC / Wilson & al            | <input type="radio"/> Hphob. HPLC / Parker & al   |
| <input type="radio"/> Hphob. HPLC pH3.4 / Cowan            | <input type="radio"/> Hphob. HPLC pH7.5 / Cowan   |
| <input type="radio"/> Hphob. / Rf mobility                 | <input type="radio"/> HPLC / HFBA retention       |
| <input type="radio"/> HPLC / TFA retention                 | <input type="radio"/> Transmembrane tendency      |
| <input type="radio"/> HPLC / retention pH 2.1              | <input type="radio"/> HPLC / retention pH 7.4     |
| <input type="radio"/> % buried residues                    | <input type="radio"/> % accessible residues       |
| <input type="radio"/> Hphob. / Chothia                     | <input type="radio"/> Hphob. / Rose & al          |
| <input type="radio"/> Ratio hetero end/side                | <input type="radio"/> Average area buried         |
| <input type="radio"/> Average flexibility                  | <input type="radio"/> alpha-helix / Chou & Fasman |
| <input type="radio"/> beta-sheet / Chou & Fasman           | <input type="radio"/> beta-turn / Chou & Fasman   |
| <input type="radio"/> alpha-helix / Deleage & Roux         | <input type="radio"/> beta-sheet / Deleage & Roux |
| <input type="radio"/> beta-turn / Deleage & Roux           | <input type="radio"/> Coil / Deleage & Roux       |
| <input type="radio"/> alpha-helix / Levitt                 | <input type="radio"/> beta-sheet / Levitt         |
| <input type="radio"/> beta-turn / Levitt                   | <input type="radio"/> Total beta-strand           |
| <input type="radio"/> Antiparallel beta-strand             | <input type="radio"/> Parallel beta-strand        |



## Primary structure analysis

- **ProtParam**  - Physico-chemical parameters of a protein sequence (amino-acid and atomic compositions, isoelectric point, extinction coefficient, etc.)
- **Compute pI/Mw**  - Compute the theoretical isoelectric point (*pI*) and molecular weight (*Mw*) from a UniProt Knowledgebase entry or for a user sequence
- **ScanSite pI/Mw** - Compute the theoretical *pI* and *Mw*, and multiple phosphorylation states
- **MW, pI, Titration curve** - Computes *pI*, composition and allows to see a titration curve
- **Scratch Protein Predictor**
- **HeliQuest** - A web server to screen sequences with specific alpha-helical properties
- **Radar** - De novo repeat detection in protein sequences
- **REP** - Searches a protein sequence for repeats
- **REPRO** - De novo repeat detection in protein sequences
- **T-REKS** - De novo detection and alignment of repeats in protein sequences
- **TRUST** - De novo repeat detection in protein sequences
- **XSTREAM** - De novo tandem repeat detection and architecture modeling in protein sequences
- **SAPS**  - Statistical analysis of protein sequences at EMBnet-CH [Also available at **EBI**]
- **Coils**  - Prediction of coiled coil regions in proteins (Lupas's method) at EMBnet-CH [Also available at PBIL]
- **Paircoil** - Prediction of coiled coil regions in proteins (Berger's method)
- **Paircoil2** - Prediction of the parallel coiled coil fold from sequence using pairwise residue probabilities with the Paircoil algorithm.
- **Multicoil** - Prediction of two- and three-stranded coiled coils
- **2ZIP** - Prediction of Leucine Zippers
- **ePESTfind** - Identification of PEST regions
- **HLA\_Bind** - Prediction of MHC type I (HLA) peptide binding
- **PEPVAC** - Prediction of supertypic MHC binders
- **RANKPEP** - Prediction of peptide MHC binding
- **SYFPEITHI** - Prediction of MHC type I and II peptide binding
- **ProtScale**  - Amino acid scale representation (Hydrophobicity, other conformational parameters, etc.)
- **Drawhca** - Draw an HCA (Hydrophobic Cluster Analysis) plot of a protein sequence

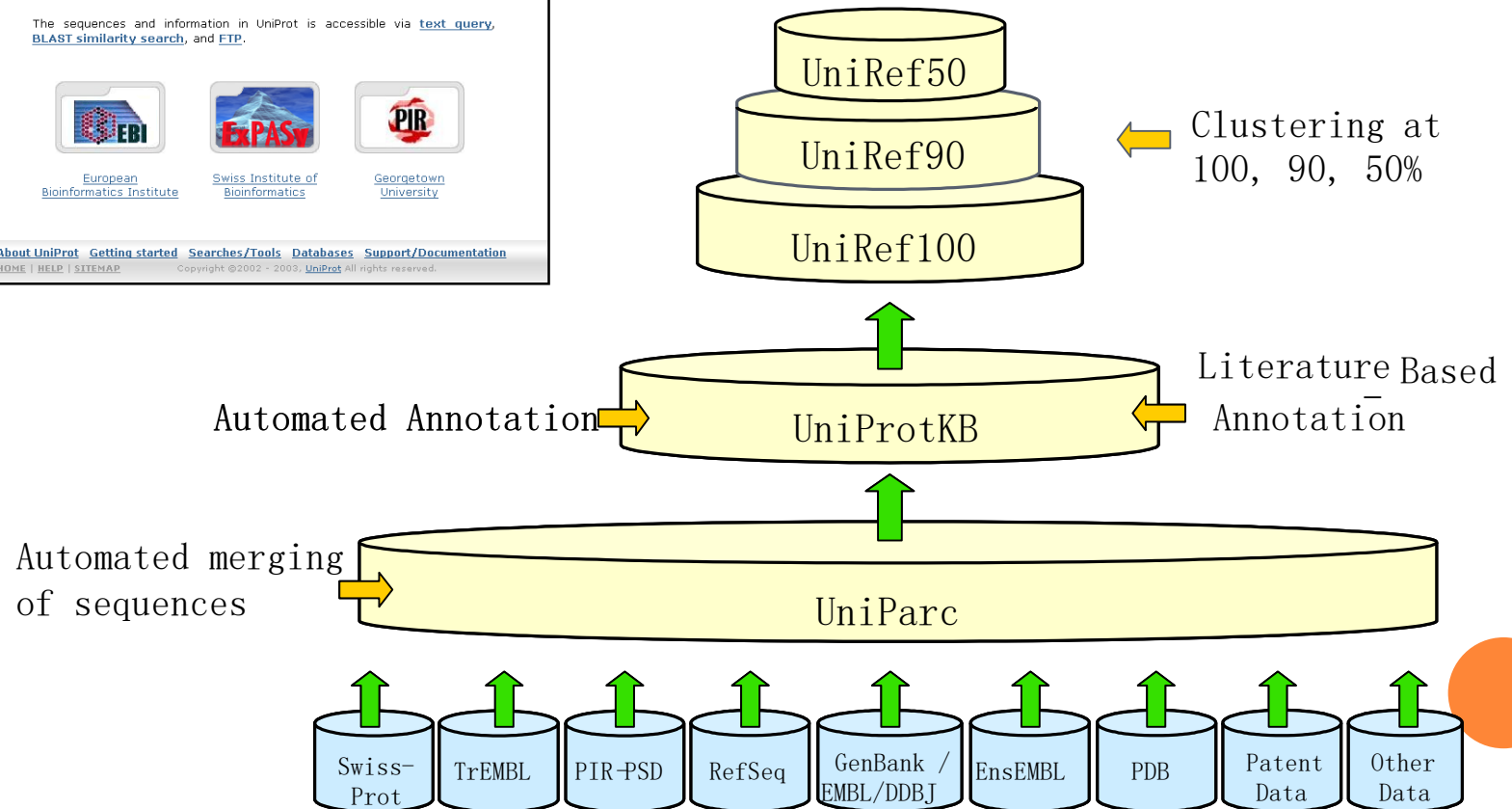
# Single protein and shotgun analysis



# UNIVERSAL PROTEIN RESOURCE



<http://www.uniprot.org/>



protein sequence analysis tool

ExPASy: SIB Bioinformatics Resource Portal

+

← → ↺ 🏠

🔒 https://www.expasy.org


⋮ 🛡️ ☆

⬇️ 📄 📖 📧 📧 📧 📧

⚙️ 最常访问 📁 火狐官方网站 🌈 新手上路 📁 常用网址 📺 JD 京东商城

+

SIB



ExPASy

Bioinformatics Resource Portal

Home About Contact

Query all databases

✖

search help

Visual Guidance

Categories

proteomics

genomics

structure analysis

systems biology

evolutionary biology

population genetics

transcriptomics

biophysics

imaging

IT infrastructure

medicinal chemistry

glycomics

Resources A..Z


Help

ExPASy is the **SIB Bioinformatics Resource Portal** which provides access to scientific databases and software tools (i.e., *resources*) in different areas of life sciences including proteomics, genomics, phylogeny, systems biology, population genetics, transcriptomics etc. (see **Categories** in the left menu). On this portal you find resources from many different SIB groups as well as external institutions.

Featuring today


Compute pI/MW

Compute isoelectric point and molecular weight from protein sequence  
[\[details\]](#)




FindPept

Identify peptides resulting from unspecific protein cleavage from their experimental masses  
[\[details\]](#)





FindMod


Predict potential protein post-translational modifications and find potential single amino




Popular resources

 UniProtKB

 SWISS-MODEL

 STRING

 PROSITE

Latest News

Protein Spotlight: A way in -  
2020-03-26

As children in Scotland, back home from school and when the weather was dry, we would fling our schoolbags into the hall and grab a few golf clubs, a ball and a tee. There was no need for a change of clothes or shoes, ...[More](#).

New neXtProt application release  
2020-02-28 - 2020-02-28

Modifications to SPARQL queries For

# TRY SWISS INSTITUTE OF BIOINFORMATICS

- <https://www.expasy.org/>
- <http://www.uniprot.org>
- How to calculate theoretical MW/PI from protein sequence?

