**Master 1 MABS**

**UE Bioinformatique des séquences**

**novembre - décembre 2013**

# APPLICATION DES ALIGNEMENTS MULTIPLES : L'IDENTIFICATION DE MOTIFS OU DE PATTERNS

# Un peu de vocabulaire

Familles, domaines, motifs, pattern, etc...

- domaine protéique: unité structurale (et fonctionnelle) indépendante, évolutivement conservée (doigt de zinc, boucle,...)

- motifs protéiques: plus courts
  - site de modification post-traductionnelle
  - site de liaison (ADN, métal,...)
  - site actif d'enzyme

- famille protéique: ensemble de protéines évolutivement reliées par un ou plusieurs domaines protéiques communs
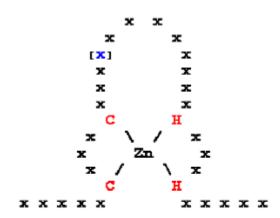
# Domaine conservé

- Exemple doigt de zinc

```
TYY1_HUMAN/383-407    YVCPF-DGCN---KKFAQSTNLKSHILT---H
YKQ8_CAEEL/78-102     YKCT---VCR---KDISSSESLRTHMFKQ-HH
BASO_HUMAN/719-742    FQCD---ICK---KTFKNACSVKIHHKN--MH
ZG29_XENLA/62-84      FVCT---VCG---KTYKYKHGLNTHLHS---H
P43_XENBO/106-130     LKCSV-PGCK---RSFRKKRALRIHVSE---H
IKAR_MOUSE/488-512    FECN---MCG---YHSQDRYEFSSHITRG-EH
Q92610/1043-1069      YTCG---YCTEDSPSFPRPSLLESHISL--MH
TRA1_CAEEL/306-331    YKCEF-ADCE---KAFSNASDRAKHQNR--TH
ZN10_HUMAN/383-405    YKCN---QCG---IIFSQNSPFIVHQIA---H
GLI1_XENLA/283-310    FVCHW-QDCSRELRPFKAQYMLVVHMRR---H
XFIN_XENLA/276-298    FRCS---ECS---RSFTHNSDLTAHMRK---H
TF3A_BUFAM/72-97      CKCET-ENCN---LAFTTASNMRLHFKR--AH
ZG58_XENLA/174-196    FVCT---ECN---LSFAGLANLRSHQHL---H
P43_XENBO/163-187     YRCSY-EDCQ---TVSPTWTALQTHLKK---H
TSH_DROME/354-378     FRCV---WCK---QSFPTLEALTTHMKDS-KH
ZN76_HUMAN/165-189    FRCGY-KGCG---RLYTTAHHLKVHERA---H
TF3A_BUFAM/219-244    YRCPR-ENCD---RTYTTKFNLKSHILT--FH
SUHW_DROAN/349-373    YACK---ICG---KDFTRSYHLKRHQKYS-SC
ZN76_HUMAN/285-309    YTCPE-PHCG---RGFTSATNYKNHVRI---H
SRYC_DROME/469-492    FKCN---YCP---RDFTNFPNWLKHTRR--RH
EVI1_HUMAN/761-784    YRCK---YCD---RSFSISSNLQRHVRN--IH
...
```

*Extrait de Pfam, entrée zf-C2H2*

# Domaine conservé

- Exemple doigt de zinc



**Motif Prosite:**

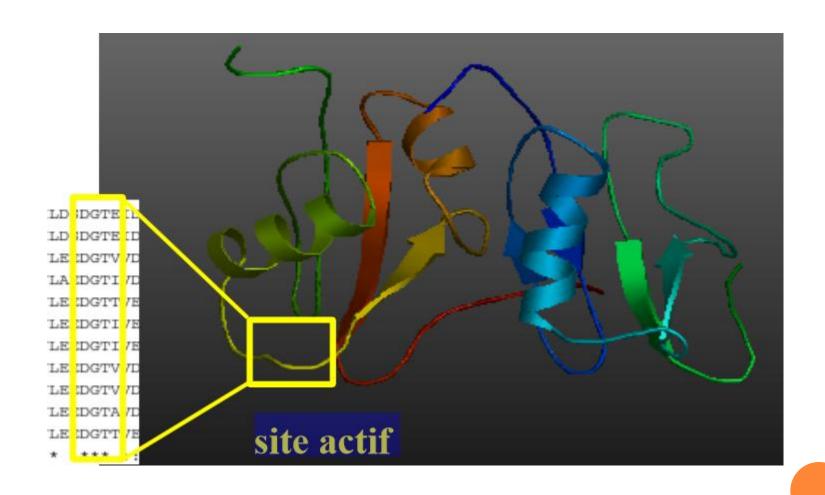C-x (2,4)-C-x (3)-[LIVMFYWC]-x (8)-H-x (3,5)-H

# MOTIF (PATTERN)

```
GUN1_TRIRE/427-455    HWGQCGGI---GYSGC--K-TCTSGTTCQYSNDYYSQCL
GUX1_TRIRE/481-509    HYGQCGGI---GYSGP--T-VCASGTTCQVLNPYYSQCL
GUX1_PHACH/484-512    QWGQCGGI---GYTGS--T-TCASPYTCHVLNPYYSQCY
GUX2_TRIRE/30-58      VWGQCGGQ---NWSGP--T-CCASGSTCVYSNDYYSQCL
GUN5_TRIRE/209-237    LYGQCGGA---GWTGP--T-TCQAPGTCKVQNQWYSQCL
GUNF_FUSOX/21-49      IWGQCGGN---GWTGA--T-TCASGLKCEKINDWYYQCV
GUX3_AGABI/24-52      VWGQCGGN---GWTGP--T-TCASGSTCVKQNDFYSQCL
Q01763/473-500        --SQCGGL---GYAGP--TgVCPSPYTCQALNIYYSQCI
GUX1_PENJA/505-533    DWAQCGGN---GWTGP--T-TCVSPYTCTKQNDWYSQCL
GUXC_FUSOX/482-510    QWGQCGGQ---NYSGP--T-TCKSPFTCKKINDFYSQCQ
GUX1_HUMGR/493-521    RWQQCGGI---GFTGP--T-QCEEPYICTKLNDWYSQCL
GUX1_NEUCR/484-512    HWAQCGGI---GFSGP--T-TCPEPYTCAKDHDIYSQCV
Q9Y894/23-53          PWGQCGGP---GWTGPttT-CCVTGCTCPVTND-YSQCL
PSBP_PORPU/26-54      LYEQCGGI---GFDGV--T-CCSEGLMCMKMGPYYSQCR
GUNB_FUSOX/29-57      VWAQCGGQ---NWSGT--P-CCTSGNKCVKLNDFYSQCQ
PSBP_PORPU/69-97      PYGQCGGM---NYSGK--T-MCSPGFKCVELNEFFSQCD
GUNK_FUSOX/339-370    AYYQCGGSKSAYPNGN--L-ACATGSKCVKQNEYYSQCV
PSBP_PORPU/128-156    EYAACGGE---MFMGA -K-CCKFGLVCYETSGKWSQCR
```

*Extrait de Prosite, entrée PS00562*

C-G-G-x(4,7)-G-x(3)-C-x(5)-C-x(3,5)-[NHG]-x-[FYWM]- x(2)-Q-C

```
                +----------------+
                |                |
                |       +-----|---------+
                |       |     |         |
                |       |     |         |
xxxxxxxCxxxxxxxxxxCxxxxxCxxxxxxxxxxCx          ← pattern
      *****************************
```

# LE PROFIL DE CONSERVATION ISSU DE L'ALIGNEMENT MULTIPLE



CLUSTAL FORMAT for T-Coffee version_3.13 (http://www.tcoffee.org), CPU=1.78 sec, SCORE=41, Nseq=11, Len=382

```
tr|O61464|O61464_DROME        METAANSG--------------------D----------SKKPFKVKDVTRNIKKAVCA
tr|Q28ZV7|Q28ZV7_DROPS        MPNAMETTT-------------------S----------SKKPFKVKDVTRNIKKAVCA
tr|Q66K97|Q66K97_XENTR        MQGALDYANALSPKSLIRSVTNVGTSLTRRVLFPPLPE-PPQRPFRVSNSDRSSKKGIVA
unk|VIRT1655|Blast_submission  MEVTGDAG--------------VPESGEIR---------TLKPCLLRRNYSREQHGVAA
Q96AQ7|CIDEC_HUMAN            MEYAMKSLSLLYPKSLSRHVSVRTSVVTQQLLSEPSPKAPRARPCRVSTADRSVRKGIMA
sp|P56198|CIDEC_MOUSE         MDYAMKSLSLLYPRSLSRHVAVSTAVVTQQLVSKPSRETPRARPCRVSTADRKVRKGIMA
tr|Q5XI33|Q5XI33_RAT          MDYAMKSLSLLYPRSLSRHVAVSTAVVTQQLVSKPSRETPRARPCRVSTADRKVRKGIMA
sp|O60543|CIDEA_HUMAN         MEAARDYAG-----ALIRPLTFMGSQTKRVLFTP---LMHPARPFRVSNHDRSSRRGVMA
tr|A4FUX1|A4FUX1_BOVIN        METARDCAG-----ALLRPLTFMGSQTKKVLFTP---FMHPARPFRVSNHDRSSRRGVMA
sp|O70303|CIDEB_MOUSE         ----MEYLSAFNPNGLLRSVSTVSSELSRRVWNS---APPPQRPFRVCDHKRTVRKGLTA
sp|Q3T191|CIDEB_BOVIN         ----MEYLSNLDPSSLLRSVSNMSADLGRKVWTS---APPRQRPFRVCDNKRTTRKGLTA
                               .                           :* :      ::.: *
```

```
tr|O61464|O61464_DROME        SSLEEIRSKVAEKFEKCDH--PTIHLD DGT IDDEEYFRTLDENTELVAVF GEHWID
tr|Q28ZV7|Q28ZV7_DROPS        ASLEEIRDKVAEKFGKCDH--PTIHLD DGT IDDEEYFRTLDENTELVAVF GEHWID
tr|Q66K97|Q66K97_XENTR        GTLKELIEKASETLFIHSD--VTLVLE DGT VDTEDFFQSLEDNTQFLLLE QKWTQ
unk|VIRT1655|Blast_submission  SCLEDLRSKACDILAIDKSLT VTLVLA DGT VDDDDYFLCLPSNTKFVALAS NEKWAY
Q96AQ7|CIDEC_HUMAN            YSLEDLLLKVRDTLMLADK--FFLVLE DGT VETEEYFQALAGDTVFMVLQ QKWQP
sp|P56198|CIDEC_MOUSE         HSLEDLLNKVQDILKLKDK--FSLVLE DGT VETEEYFQALAKDTMFMVLL QKWKP
tr|Q5XI33|Q5XI33_RAT          HSLEDLLGKVQDILKLKDK--FSLVLE DGT VETEEYFQALPRDTVFMVLQ QKWKS
sp|O60543|CIDEA_HUMAN         SSLQELISKTLDALVIATG--VTLVLE DGT VDTEEFFQTLGDNTHFMILE QKWMP
tr|A4FUX1|A4FUX1_BOVIN        SSLQELLSKTLDALVVASQ--VTLVLE DGT VDTEEFFQTLGDNTHLMVLE QKWTP
sp|O70303|CIDEB_MOUSE         ASLQELLDKVLETLLLRG---LTLVLE DGT VDSEDFFQLLEDDTCLMVLE GQSWSP
sp|Q3T191|CIDEB_BOVIN         ATRQELLDKALEALVLSG---LTLVLE DGT VESEEFFQLLEDDTCLMVLE GQSWSP
                               :::  *. : :         : * *** :: :::*  *  :* :: :   : *
```

* = résidu parfaitement conservé
: = substitution conservative
. = substitution semi-conservative

```
tr|O61464|O61464_DROME        PTHYVTITTPHGNEAGTGNGELNGGGEG----------DTTDANNSES-ARIRQLVGQLQ
tr|Q28ZV7|Q28ZV7_DROPS        PTHYVTITTPHGSETVTGNGDISSGGVGGGVGGSCDGGDTTDANHSESAARIRQLVGQLQ
tr|Q66K97|Q66K97_XENTR        ERNSKRAV--------------------------------------------------Q
unk|VIRT1655|Blast_submission  NNSDGGTA-------WISQESF---------------DVDETDSGAG-LKWKNVARQLK
Q96AQ7|CIDEC_HUMAN            PSEQGTRH--------------------PLSL-----------------SH--------K
sp|P56198|CIDEC_MOUSE         PSEQRKKR--------------------AQLAL----------------SQ--------K
tr|Q5XI33|Q5XI33_RAT          PSEQRKKK--------------------AQLSL----------------SQ--------K
sp|O60543|CIDEA_HUMAN         GSQHVP---------------------------------------------TC------S
tr|A4FUX1|A4FUX1_BOVIN        AGHQTP------------------------------------------------AR-----R
sp|O70303|CIDEB_MOUSE         KS-G---M-------------------LSYGLG----------------RE--------K
sp|Q3T191|CIDEB_BOVIN         RRSG---V-------------------LSYGLG----------------QE--------K
```

**Qu'est ce qu'il y a de si spécial ici ???**

```
tr|O61464|O61464_DROME        NNLCNVSVMNDADLDSLSNMDPNSLVD------ITGKEFMEQLKDAGRPLCAKRNAEDRL
tr|Q28ZV7|Q28ZV7_DROPS        NNLCNVSVMNDADLDSLSNMDPNSLVD------ITGKEFMEQLKDAGRPLCAKRNAEDRL
tr|Q66K97|Q66K97_XENTR        HEKKTGIANLTFDLYKLNP----------------------------------------
```

site actif

```
GHEGVGKVVKLGAGA
GHEKKGYFEDRGPSA
GHEGYGGRSRGGGYS
GHEFEGPKGCGALYI
GHELRGTTFMPALEC
```

```
GHE--G----------    Consensus 100%
GHE--G-----G---     Consensus 60 %
```

```
GHE-x(2)-G-x(5)-[GA]
```

*pattern*
signature

profil

précision
sensibilité

`<A-x-[ST](2)-x(0,1)-[APTL]-x(4,10)-C-{V}`

**<A** en N terminal

**x** = n'importe quel AA

**ST](2)** = Ser ou Thr 2 fois

**x(0,1)** 1 aa ou aucun

**x(4,10)** entre 4 et 10 aa quelconques

**{V}** tout sauf une Val

Code IUPAC pour les nucléotides

| Code | Description |
|------|-------------|
| A | Adénine |
| C | Cytosine |
| G | Guanine |
| T | Thymine |
| U | Uracile |
| R | Purine (A ou G) |
| Y | Pyrimidine (C, T, ou U) |
| M | C ou A |
| K | T, U, ou G |
| W | T, U, ou A |
| S | C ou G |
| B | C, T, U, ou G (pas A) |
| D | A, T, U, ou G (pas C) |
| H | A, T, U, or C (pas G) |
| V | A, C, or G (pas T, pas U) |
| N | Toutes les bases (A, C, G, T, ou U) |

| ExPASy Home page | Site Map | Search ExPASy | Contact us | PROSITE |
|---|---|---|---|---|

Search PROSITE for P00174    Go    Clear

# NiceSite View of PROSITE: PS00191

| General information about the entry | |
|---|---|
| Entry name | CYTOCHROME_B5_1 |
| Accession number | PS00191 |
| Entry type | PATTERN |
| Date | APR-1990 (CREATED); DEC-2004 (DATA UPDATE); SEP-2005 (INFO UPDATE). |
| PROSITE documentation | PDOC00170 |

Pattern

| Name and characterization of the entry | |
|---|---|
| Description | Cytochrome b5 family, heme-binding domain signature. |
| Pattern | [FY]-[LIVMK]-{I}-{Q}-H-P-[GA]-G. |

[FY]-[LIVMK]-{I}-{Q}-H-P-[GA]-G

**Numerical results**

- UniProtKB/Swiss-Prot release number: **48.1**, total number of sequence entries in that release: **195058**.
- Total number of hits in UniProtKB/Swiss-Prot: **86 hits in 86 different sequences**
- Number of hits on proteins that are known to belong to the set under consideration: **80 hits in 80 different sequences**
- Number of hits on proteins that could potentially belong to the set under consideration: **0 hits in 0 different sequences**
- Number of false hits (on unrelated proteins): **6 hits in 6 different sequences**
- Number of known missed hits: **4**
- Number of partial sequences which belong to the set under consideration, but which are not hit by the pattern or profile because they are partial (fragment) sequences: **2**
- Precision (true hits / (true hits + false positives)): **93.02 %**
- Recall (true hits / (true hits + false negatives)): **95.24 %**

**Comments**

- Taxonomic range: **Eukaryotes, Prokaryotes (Bacteria), Eukaryotic viruses**
- Maximum known number of repetitions of the pattern in a single protein: **1**
- 'Interesting' site in the pattern: **5,heme_iron**
- VERSION: **1**

**Cross-references**

| | True positive hits: |
|---|---|
| | ACO1_AJECA  (Q12618),  ACO1_YEAST  (P21147),  CYB2_HANAN  (P09437), |
| | CYB2_YEAST  (P00175), CYB51_ARATH (Q42342), CYB52_ARATH (O48845), |
| | CYB52_SCHPO (Q9USM6), CYB5L_MIMIV (Q5UR80), CYB5L_NEUCR (Q8X0J4), |
| | CYB5M_HUMAN (O43169), CYB5M_MOUSE (Q9CQX2), CYB5M_PONPY (Q5RDJ5), |
| | CYB5M_RAT   (P04166)  CYB5P_DROME (P19967)  CYB5P_DROVI (P50266) |

## Alignement

**AATTGA**

**AGGTCC**

**AGGATG**

**AGGCGT**

## Matrice de position

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| A | 4 | 1 | 0 | 1 | 0 | 1 |
| C | 0 | 0 | 0 | 1 | 1 | 1 |
| G | 0 | 3 | 3 | 0 | 2 | 1 |
| T | 0 | 0 | 1 | 2 | 1 | 1 |

## Matrice de fréquences

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| **A** | 1 | 0.25 | 0 | 0.25 | 0 | 0.25 |
| **C** | 0 | 0 | 0 | 0.25 | 0.25 | 0.25 |
| **G** | 0 | 0.75 | 0.75 | 0 | 0.50 | 0.25 |
| **T** | 0 | 0 | 0.25 | 0.50 | 0.25 | 0.25 |

Matrice de
fréquences

$$\log \left[ \frac{f_b}{p_b} \right]$$

Matrice de poids
de position
(Position Weight Matrix ou
Position Specific Scoring Matrix
PSSM)

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| **A** | 1.2 | 0 | -1.6 | 0 | -1.6 | 0 |
| **C** | -1.6 | -1.6 | -1.6 | 0 | 0 | 0 |
| **G** | -1.6 | 0.96 | 0.96 | -1.6 | 0.59 | 0 |
| **T** | -1.6 | -1.6 | 0 | 0.59 | 0 | 0 |

C  G  T  A  T  G  T  A  A  G  G  T  G  T  A  C  G  T  A  G

Pour trouver si la séquence contient un motif, les bases de données appliquent sur les séquences soumises les matrices qu'elles ont correspondant toutes à des motifs décrits, et notamment la matrice ci-dessous.

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| **A** | 1.2 | 0 | -1.6 | 0 | -1.6 | 0 |
| **C** | -1.6 | -1.6 | -1.6 | 0 | 0 | 0 |
| **G** | -1.6 | 0.96 | 0.96 | -1.6 | 0.59 | 0 |
| **T** | -1.6 | -1.6 | 0 | 0.59 | 0 | 0 |

Calcul de score par fenêtre glissante

C  G  T  A  T  G  T  A  A  G  G  T  G  T  A  C  G  T  A  G

C G T A T G [T A A G G T] G T A C G T A G

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| A | 1.2 | [0] | [-1.6] | 0 | -1.6 | 0 |
| C | -1.6 | -1.6 | -1.6 | 0 | 0 | 0 |
| G | -1.6 | 0.96 | 0.96 | [-1.6] | [0.59] | 0 |
| T | [-1.6] | -1.6 | 0 | 0.59 | 0 | [0] |

Score = -4.21

C G T A T G T [A A G G T G] T A C G T A G

|   | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| A | [1.2] | [0] | -1.6 | 0 | -1.6 | 0 |
| C | -1.6 | -1.6 | -1.6 | 0 | 0 | 0 |
| G | -1.6 | 0.96 | [0.96] | [-1.6] | 0.59 | [0] |
| T | -1.6 | -1.6 | 0 | 0.59 | [0] | 0 |

Score = 0.56

C G T A T G T A **A G G T G T** A C G T A G

| | 1 | 2 | 3 | 4 | 5 | 6 |
|---|---|---|---|---|---|---|
| A | 1.2 | 0 | -1.6 | 0 | -1.6 | 0 |
| C | -1.6 | -1.6 | -1.6 | 0 | 0 | 0 |
| G | -1.6 | 0.96 | 0.96 | -1.6 | 0.59 | 0 |
| T | -1.6 | -1.6 | 0 | 0.59 | 0 | 0 |

Score = 4.3

Si dans la base de données interrogée, le seuil minimal de détection de motif est par exemple S = 4 alors seule la séquence AGGTGT sera considérée comme pattern

# EXEMPLE : RECHERCHE DE MOTIF SUR UNE SÉQUENCE VEM-1 DE *CAENORHABDITIS ELEGANS*

## Vema (Mammalian ventral midline antigen) related protein 1, isoform a

UniProtKB/Swiss-Prot: Q9TY05

GenPept    Graphics

```
>gi|75024827|sp|Q9TY05|Q9TY05_CAEEL Vema (Mammalian ventral midline antigen)
related protein 1, isoform a
MYTVSVTFLHKSFFTMDLSSWFEFTMYDAVFLVVVLGFFFYWLTRSEQPLPAPPKELAPLPMSDMTVEEL
RKYDGVKNEHILFGLNGTIYDVTRGKGFYGPGKAYGTLAGHDATRALGTMDQNAVSSEWDDHTGISADEQ
ETANEWETQFKFKYLTVGRLVKNSSEKADYGNRKSFVRGAESLDSIINGGDEGTKKDN
```

HOME | SEARCH | BROWSE | FTP | HELP | ABOUT

**Pfam 27.0 (March 2013, 14831 families)**

The Pfam database is a large collection of protein families, each represented by **multiple sequence alignments** and **hidden Markov models (HMMs)**. **More...**

**QUICK LINKS**

**SEQUENCE SEARCH**

**VIEW A PFAM FAMILY**

**VIEW A CLAN**

**VIEW A SEQUENCE**

**VIEW A STRUCTURE**

**KEYWORD SEARCH**

**JUMP TO**

**ANALYZE YOUR PROTEIN SEQUENCE FOR PFAM MATCHES**

Paste your protein sequence here to find matching Pfam families.

Go    Example

This search will use and an E-value of 1.0. You can set your own search parameters and perform a range of other searches here.

# Exemple : recherche de motif sur une séquence vem-1 de *Caenorhabditis elegans*

# EXEMPLE: RECHERCHE DE MOTIF SUR UNE SÉQUENCE VEM-1 DE *CAENORHABDITIS ELEGANS*

# Family: *Cyt-b5* (PF00173)

159 architectures   5879 sequences   3 interactions   836 species   98 structures

**Summary** | Domain organisation | Clan | Alignments | HMM logo | Trees | Curation & m... | Species | Interactions | Structures



Jump to... ⓘ

[enter ID/acc] [Go]

## Summary: Cytochrome b5-like Heme/Steroid binding domain

Pfam includes annotations and additional family information from a range of different sources. These sources can be accessed via the tabs below.

**Wikipedia: Cytochrome b5** | Pfam | InterPro

This is the Wikipedia entry entitled "Cytochrome b5🗗". **More...**

### Cytochrome b5   [Edit Wikipedia article]

**Cytochromes b$_5$** are ubiquitous electron transport hemoproteins found in animals, plants, fungi and purple phototrophic bacteria. The microsomal and mitochondrial variants are membrane-bound, while bacterial and those from erythrocytes and other animal tissues are water-soluble. The family of cytochrome b$_5$-like proteins includes (besides cytochrome b$_5$ itself) hemoprotein domains covalently associated with other redox domains in flavocytochrome cytochrome b$_2$ (L-lactate dehydrogenase; EC 1.1.2.3🗗), sulfite oxidase (EC 1.8.3.1🗗), plant and fungal nitrate reductases (EC 1.7.1.1🗗, EC 1.7.1.2🗗, EC 1.7.1.3🗗), and plant and fungal cytochrome b$_5$/acyl lipid desaturase fusion proteins.

| **Contents** [hide] |
| --- |
| 1 Structure |
| 2 Cytochrome b$_5$ in some biochemical reactions |
| 3 See also |
| 4 References |
| 5 External links |

### Structure

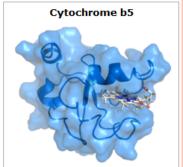3-D structures of a number of cytochrome b$_5$ and yeast flavocytochrome b$_2$ are known. The fold belongs to the α+β class, with two hydrophobic cores on each side of a β-sheet. The larger hydrophobic core constitutes the heme-binding pocket, closed off on each side by a pair of helices connected by a turn. The smaller hydrophobic core may have only a structural role and is formed by spatially close N-terminal and C-terminal segments. The two histidine residues provide the fifth and sixth heme ligands, and the propionate edge of the heme group lies at the opening of the heme crevice. Two isomers of cytochrome b$_5$, referred to as the A (major) and B (minor) forms, differ by a 180° rotation of the heme about an axis defined by the α- and γ-meso carbons.

### Cytochrome b$_5$ in some biochemical reactions

EC 1.6.2.2🗗 cytochrome-b$_5$ reductase

NADH + H$^+$ + 2 ferricytochrome b$_5$ → NAD$^+$ + 2 ferrocytochrome b$_5$

EC 1.10.2.1🗗 L-ascorbate—cytochrome-b$_5$ reductase

**Cytochrome b5**

Rat cytochrome b5 bound to heme

| Identifiers | |
| --- | --- |
| **Symbol** | CYB5A |
| **Alt. symbols** | CYB5 |
| **Entrez** | 1528🗗 |
| **HUGO** | 2570🗗 |
| **OMIM** | 250790🗗 |
| **PDB** | 1JEX🗗 |
| **RefSeq** | NM_001914🗗 |
| **UniProt** | P00167🗗 |
| **Other data** | |
| **Locus** | Chr. 18 q23🗗 |

**Cytochrome b5**

| Identifiers | |
| --- | --- |
| **Symbol** | Cyt_B5 |
| **Pfam** | PF00173🗗 |
| **InterPro** | IPR001199🗗 |
| **PROSITE** | PDOC00170🗗 |

# Family: *Cyt-b5* (PF00173)

Loading page components (1 remaining)...

60 architectures  1547 sequences  2 interactions  316 species  63 structures

Summary

Domain organisation

Alignments

HMM logo

Trees

Curation & models

Species

Interactions

Structures

Jump to...

enter ID/acc  Go

## HMM logo

HMM logos are one way of visualising profile HMMs. They provide a quick overview of the properties of an HMM in a graphical form. You can see a more detailed description of HMM logos and find out how you can interpret them here. **More...**
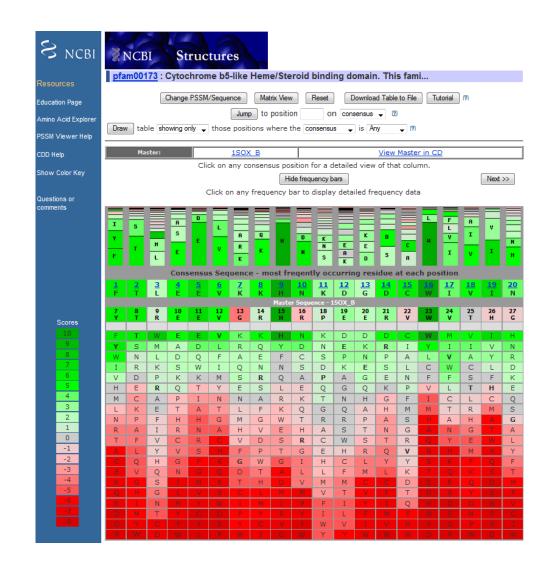
# WEB LOGO ENTIER DE CYT-B5

# DESCRIPTION DE CYT-B5 DANS LA CONSERVED DOMAINS DATABASE DU NCBI

# PSSM
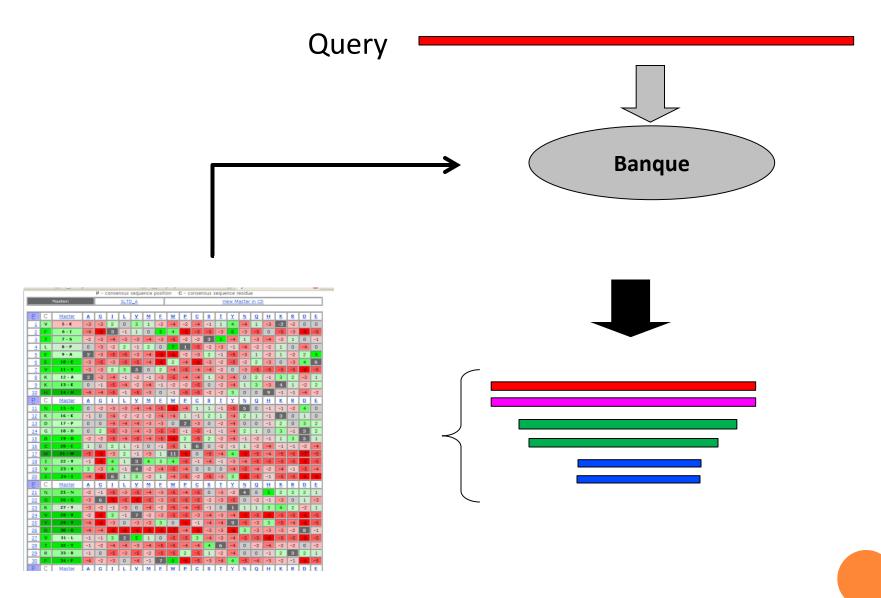
**PSI-BLAST (Position-Specific Iterative)**

- alignements multiples de hits ayant les meilleurs scores dans un blast classique

- génération d'un profil en calculant un score pour chacune des positions de l'alignement (PSSM)

- utilisation façon itérative de ce profil pour faire de nouvelles recherches et affinage à chaque itération

**PHI-BLAST (Pattern Hit Initiated BLAST )**

Pattern donné par l'utilisateur puis PSI-BLAST
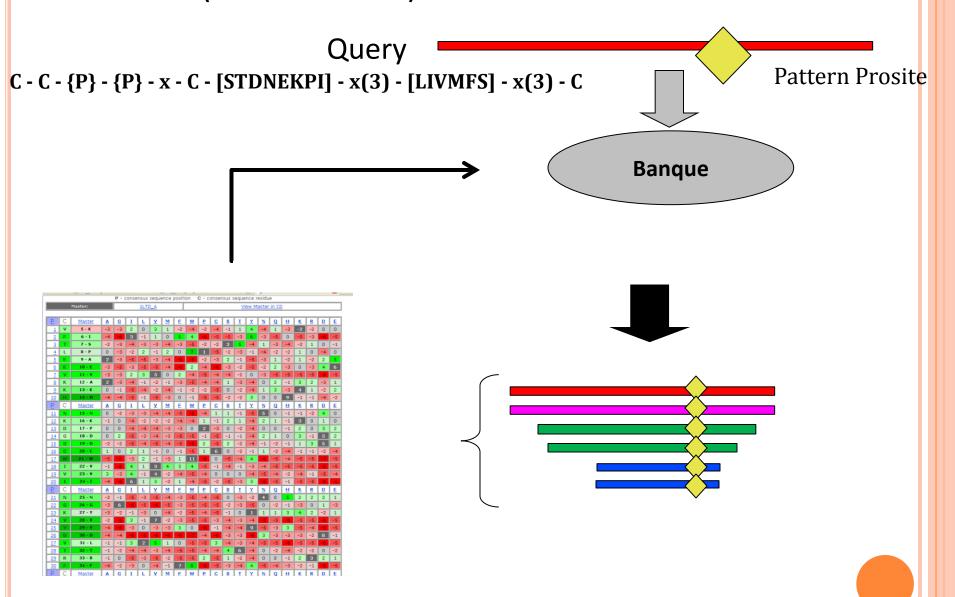
Intérêt : recherche de familles de protéines
          détecter des membres que BLAST ne trouve pas

# PSI-BLAST (Position-Specific Iterative)



PSSM *(position specific score matrix)*
(matrice de poids de position)

# PHI-BLAST (Pattern Hit Initiated)

Query

C - C - {P} - {P} - x - C - [STDNEKPI] - x(3) - [LIVMFS] - x(3) - C

Pattern Prosite

Banque



PSSM *(position specific score matrix)*
(matrice de poids de position)