

**République Algérienne Démocratique et populaire**

**Université Batna 2**

**Cours de Bioinformatique  
(M1 Biotechnologie végétale +  
M1 EDP)**

**Cours 1 :**

**Les Banques de séquences biologiques**

**Réaliser par:**

**Mr. Ghedadba N**

# Contenu de la matière

**Chapitre 1- Les banques de séquences biologiques**

**Chapitre 2- Alignement de séquences biologiques**

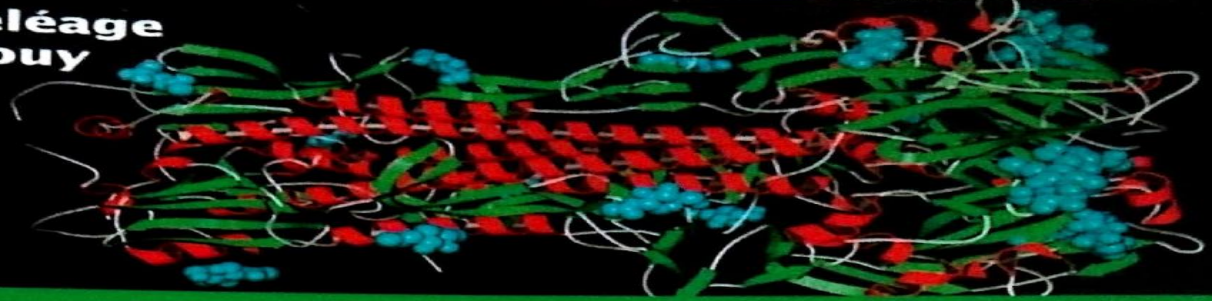
**Chapitre 3- La phylogénie moléculaire et l'annotation des génomes**

# Objectif

- *Initier les étudiants aux problématiques bio-informatiques liées à l'émergence des nouvelles biotechnologies.*
- Donner aux étudiants la connaissance et les moyens pour utiliser les logiciels existants sur le Web qui permettent déjà de traiter de manière puissante les données biologiques générées par les nouvelles biotechnologies (bases de données, logiciels de traitement de séquence, logiciels statistiques).

# Références

Gilbert Deléage  
Manolo Gouy



## Bioinformatique

Cours et cas pratique

Avec ce livre, des  
bonus sur le web



Licence 3  
Master  
Écoles d'ingénieurs

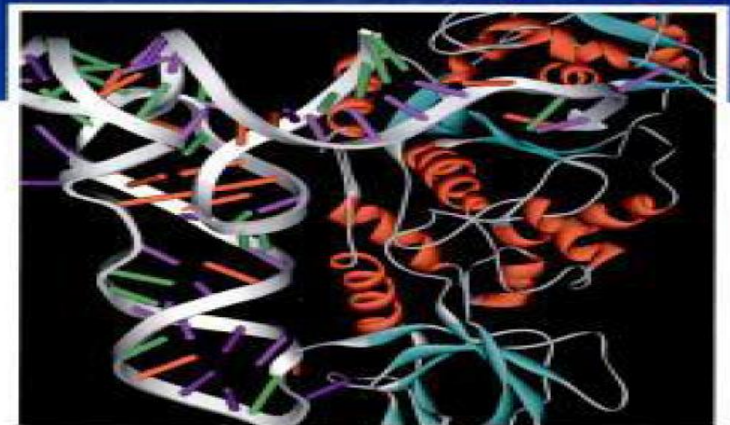
DUNOD

# BIOSCIENCES ET TECHNIQUES

Collection dirigée par J. Figarella et A. Calas

3<sup>e</sup> édition

G. Coutouly, E. Klein, E. Barbieri,  
M. Kriat



*BIBLIOTHEQUE SCIENTIFIQUE FB*  
**CONFIDENTIEL**

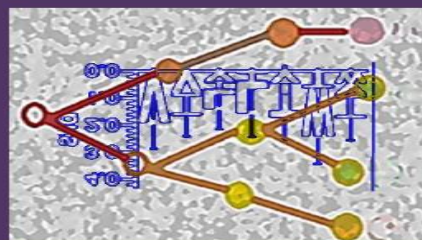
**Travaux dirigés**  
de biochimie,  
biologie moléculaire  
et bioinformatique



# Cours de phylogénie moléculaire

## Distances et constructions phylogénétiques

Support pédagogique de phylogénie moléculaire destiné aux étudiants du système LMD de Master (M1 et M2) et doctorants de Biotechnologie végétale, Biochimie et Microbiologie.



UNIVERSITE  
CONSTANTINE 1

Faculté des Sciences de  
la Nature et de la Vie

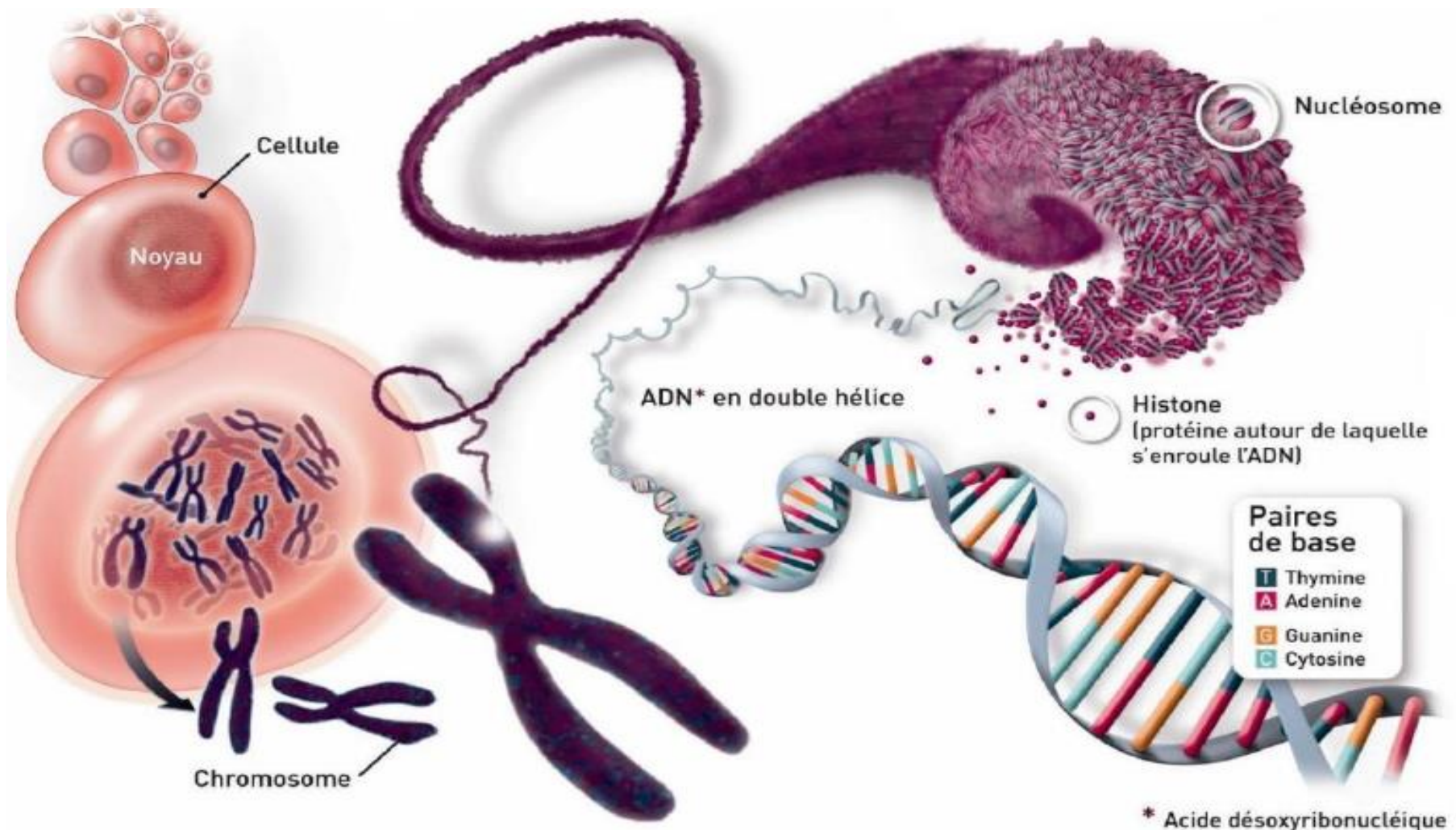
Pr DJEKOUN A.

Pr HAMIDECHI M. A.

# Dogme central de la biologie moléculaire



**Le génome:** c'est-à-dire l'ensemble des gènes d'un organisme, est structurellement défini par les chromosomes au sein de chaque cellule, selon la théorie chromosomique de l'hérédité. Chaque chromosome est constitué d'une molécule d'ADN (Figure 1), support physique des gènes, et de protéines associées





# L'ADN

- L'ADN est le **support de l'information génétique.**
- L'ADN est une **longue molécule, faite de deux brins s'enroulant en une double hélice.**
- Les deux brins de la double hélice suggèrent un mécanisme de réplication de l'ADN
- Chaque brin est le support d'une **succession de nucléotides**
- **Quatre types de nucléotides : (Adénine, Cytosine, Guanine, Thymine).**

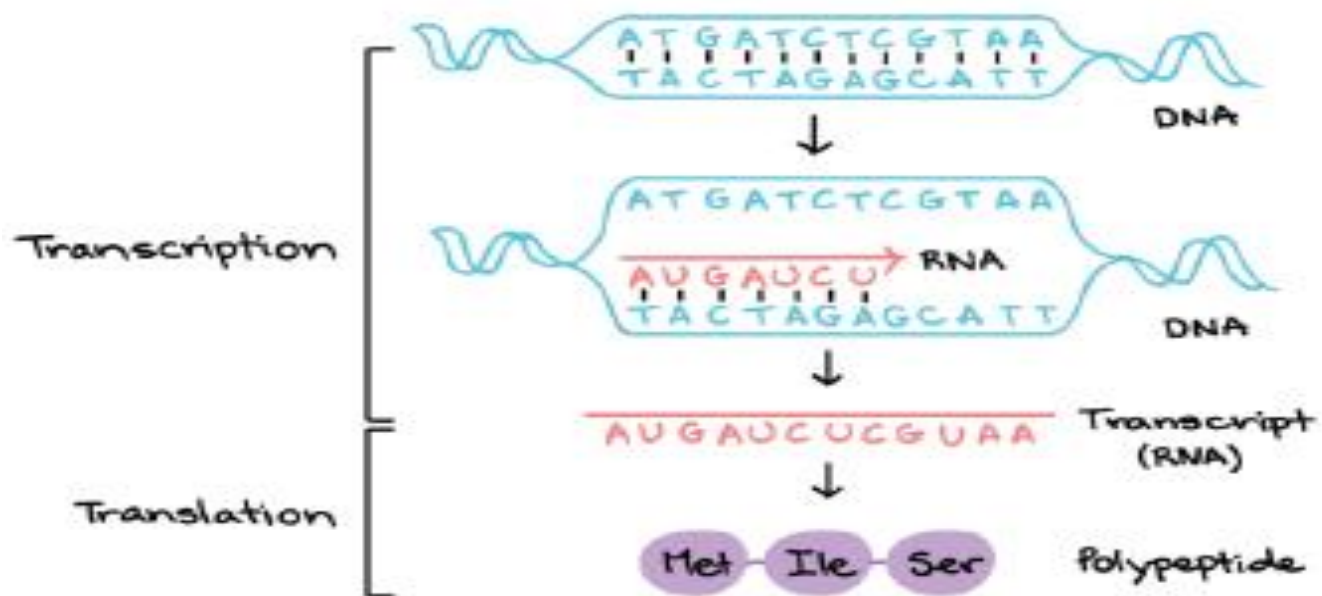
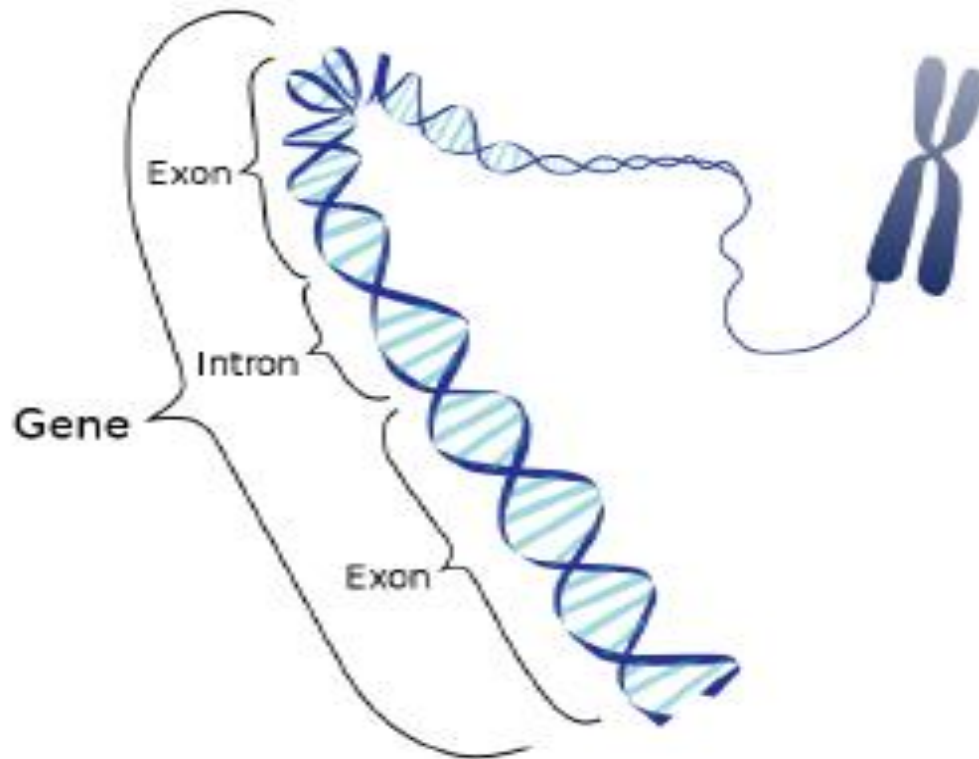
**La structure de l'ADN représenté par la figure 2, les deux brins sont orientés, on parle de l'orientation 5' → 3' par des considérations biochimiques.**

- **L'extrémité 5' se termine par un groupement phosphate.**
- **L'extrémité 3' se termine par un groupement hydroxyle.**

**Il faut retenir qu'il y a un brin sens orienté 5' → 3' et un brin anti-sens orienté 3' → 5', et généralement lorsqu'on veut connaître la séquence des nucléotides le long d'un brin, on va la lire traditionnellement dans le sens 5' → 3'.**



**Figure 2.** Représentation schématique d'une succession de nucléotides d'un fragment d'ADN double brins.



- Le texte génomique est écrit dans un alphabet de 4 lettres : A, C, G, T

## Language protéique :

### ■ Acides aminés : codes à 1 et 3 lettres

- Acide aspartique (D, Asp)
- Acide glutamique (E, Glu)
- Alanine (A, Ala)
- Arginine (R, Arg)
- Asparagine (N, Asn)
- Cystéine (C, Cys)
- Glutamine (Q, Gln)
- Glycine (G, Gly)
- Histidine (H, His)
- Isoleucine (I, Ile)
- Leucine (L, Leu)
- Lysine (K, Lys)
- Méthionine (M, Met)
- Phénylalanine (F, Phe)
- Proline (P, Pro)
- Sérine (S, Ser)
- Thréonine (T, Thr)
- Tryptophane (W, Trp)
- Tyrosine (Y, Tyr)
- Valine (V, Val)

# **I. Introduction à la Bioinformatique**

# I. Définition de la Bio-informatique

- C'est un traitement de **l'information Biologique** stockée sous forme de données accessibles aisément et exploitable
- **La bioinformatique est définie comme l'utilisation de bases de données et d'algorithmes informatiques pour analyser, les gènes, les protéines, et la collection complète d'acide désoxyribonucléique (ADN) d'un organisme vivant (le génome)**

# I. Définition moderne de la Bio-informatique

- Est une multi discipline (discipline hybride) théorique de l'analyse *in Silico* de l'information biologique contenu dans les séquences nucléiques (ADN et ARN) et protéiques

# La Bio-information

## - La Bio-informatique s'intéresse aux données :

- **Le génome:** est la totalité de l'ADN d'un individu ou d'une espèce.
- **Le transcriptome:** est l'ensemble des ARN messagers transcrits à partir du génome.
- **Le protéome** (l'ensemble des protéines biosynthétisées).
- **le sécrétome** (l'ensemble des protéines sécrétées)
- **le métabolome:** l'ensemble des métabolites produits par une cellule (molécules organiques qui ne sont ni ADN, ni protéines)



# Domaines de la Bio-informatique (Applications)

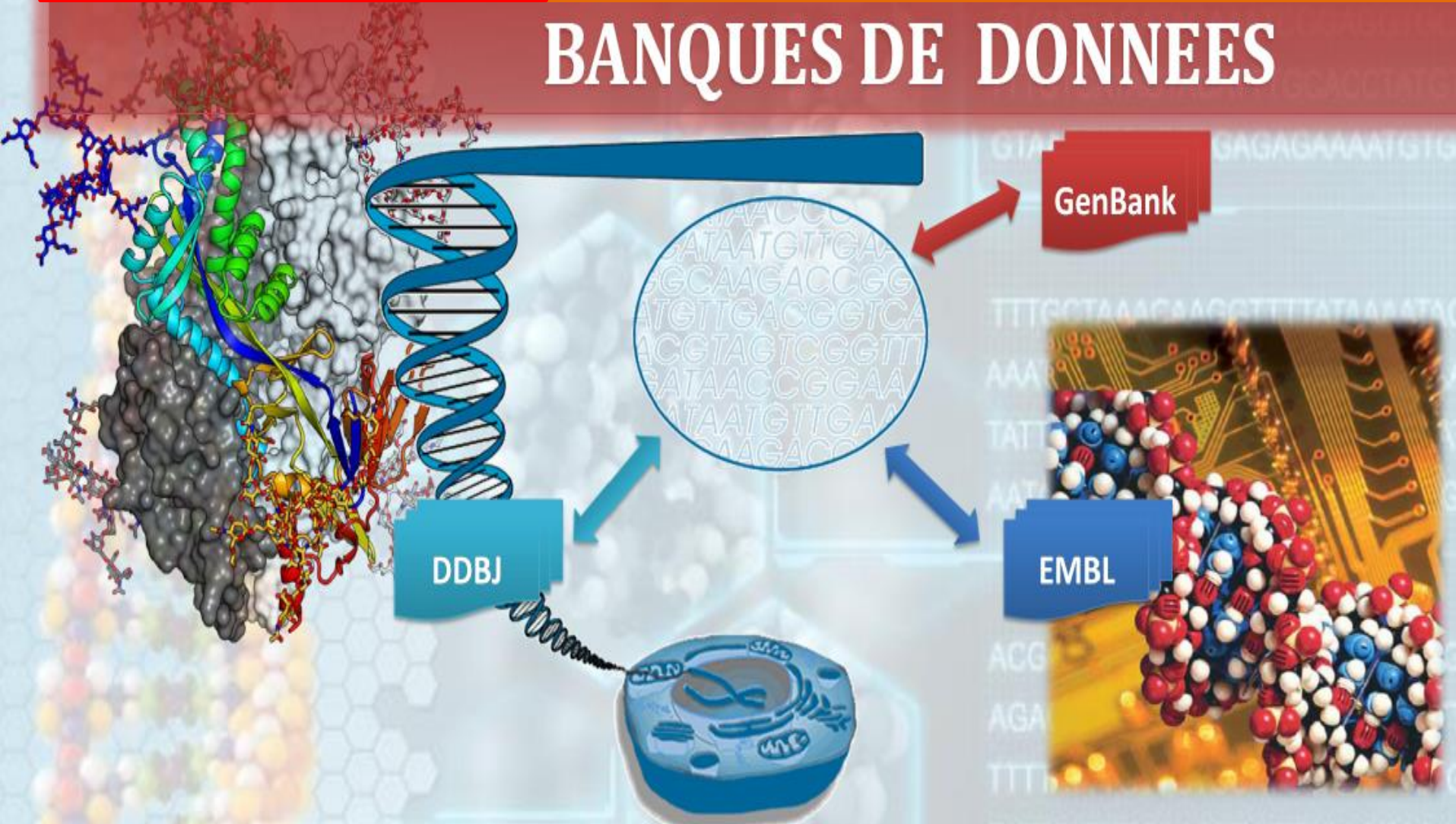
- **Stockage et Gestion des données** : Banques de données généralistes et spécialisées.
- **Structures moléculaires** : Visualisation, analyse, classification, prédiction.
- **Analyse de séquences** : Alignements, recherches de similarités, détection de motifs.
- **Génomique structurale** : Annotation des génomes,
- **Génomique fonctionnelle** : Transcriptome, protéome,
- **Phylogénie** : Relations évolutives entre gènes, entre génomes, entre organismes ; Inférence de scénarios évolutifs.

# Exemples d'applications

- Recherche en biologie
  - L'organisation moléculaire de la cellule / organisme
  - Développement
  - Mécanismes de l'évolution
- Médecine
  - Diagnostic de cancers
  - Détection des gènes impliqués dans le cancer
- La recherche pharmaceutique
  - mécanismes d'action des médicaments
  - identification de cibles pharmaceutiques
- Biotechnologie
  - La thérapie génique

# CHAPITRE I :

## BANQUES DE DONNEES



# 1. Définition

Les bases de données biologiques sont des **bibliothèques électronique et informatisé** qui contiennent des informations sur les sciences de la vie, collectées grâce à des expériences scientifiques, à la littérature publiée, aux technologies expérimentales à haut débit, et aux analyses informatiques.

## **2. Rôle des banques et bases de données biologiques**

Leur principale mission est de rendre publiques les séquences qui ont été déterminées, ainsi un des premiers intérêts de ces banques est la masse de séquences qu'elles contiennent. Entre autres ils ont pour mission l'archivage, le stockage, la diffusion et l'exploitation des données biologiques.

### 3. Contenus des bases de données biologiques

Ces bases de données peuvent contenir des informations : (**ADN, protéines, gènes et génomes, taxonomie**, autres, ...etc.). On y trouve également une bibliographie et une expertise biologique directement liées aux séquences traitées.

## 4- Types de Bases des données

### ❑ *Banques de données généralistes :*

- Celles qui correspondent à une collecte des données la plus exhaustive possible et la plus large possible.

### ❑ *Banque ou base de données spécialistes :*

- Celles qui correspondent à des données plus homogènes établies autour d'une thématique

# 1. Banques de séquences généralistes

– Données globales (pas de focus sur une application ou organisme particulier) informations hétérogènes

➤ Banques de séquences nucléiques (ADN et ARN)

➤ Banques de séquences protéiques

➤ Banques de structure 3D de macromolécules

➤ Banques d'articles scientifiques (Bibliographique)

□ **Avantage** : tout est consultable en une fois

□ **Inconvénients** : difficiles à maintenir, difficiles à interroger, problèmes de redondances



# 1.1. Les banques de séquences nucléiques

- **Origine des données**

- Séquençage d'ADN et d'ARN

- Les données stockées : séquences + annotations

- Fragments de génomes

- Un ou plusieurs gènes, un bout de gène, séquence intergénique, ...

- Génomes complets

- ARNm, ARNt, ARNr, ... (fragments ou entiers)

[Remarque 1] : toutes les séquences (ADN ou ARN) sont écrites avec des T

[Remarque 2] : les séquences sont toujours orientées 5' → 3'.

- Il existe sur internet 3 principales banques qui sont interconnectés

USA: *GenBank* (NCBI)

<https://www.ncbi.nlm.nih.gov/genbank/>

Europe: *ENA* (EMBL-EBI)

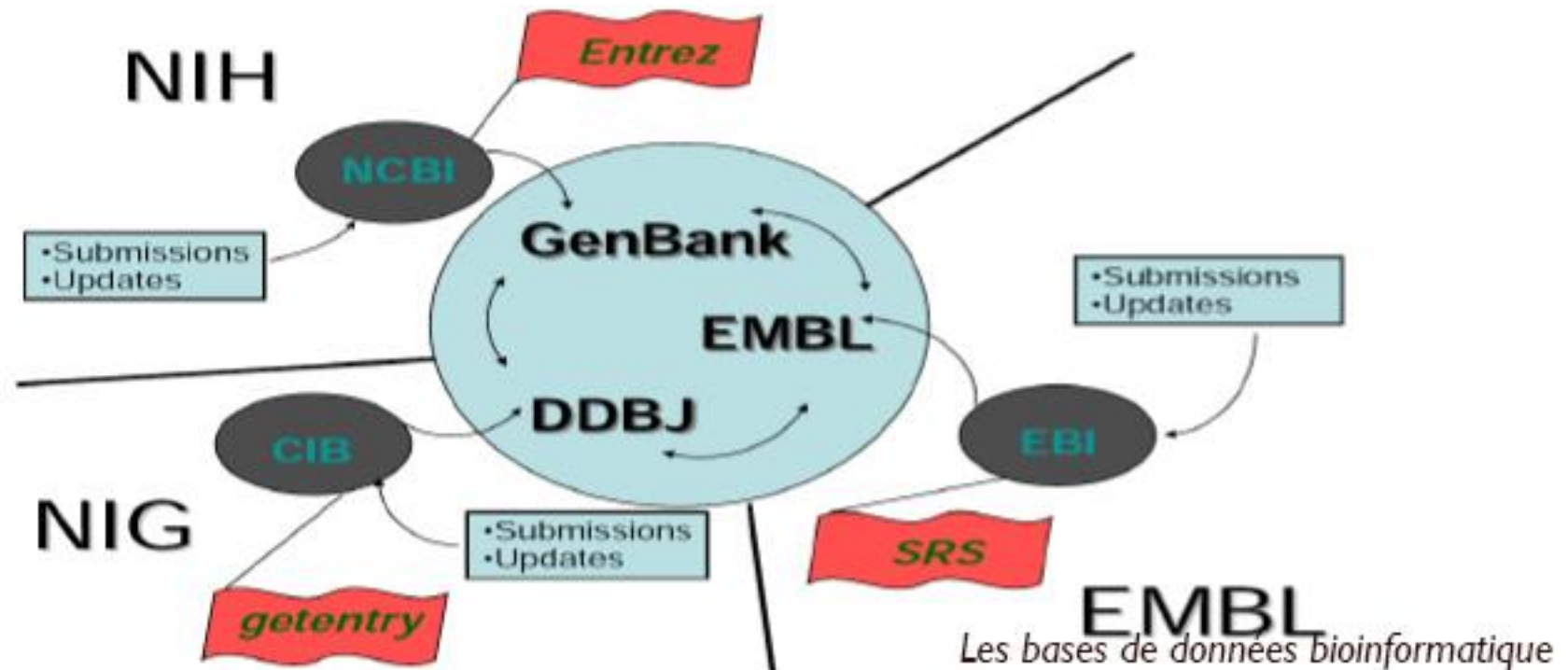
<https://www.ebi.ac.uk/ena>

Japon : *DDBJ*

<http://www.ddbj.nig.ac.jp/index-e.html>

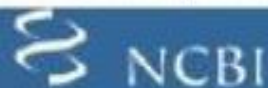


- ❑ 1988 Regroupement de ces 3 groupes (appelés International Nucleotide Sequence Database Collaboration, **INSDC**)
  - ✓ Accord pour un format commun
  - ✓ Echange des données journalier
  - ✓ Chaque groupe gère les mises a jour des séquences qu'il a créées.



# Genbank

<http://www.ncbi.nlm.nih.gov/Genbank/>



## GenBank Overview

PubMed

Entrez

BLAST

OMIM

Books

Taxonomy

Structure

Search Entrez for

Go

NCBI Home

NCBI Site Map

Submit to GenBank

Submit an update

Search GenBank

GenBank and RefSeq: a comparison

BLAST

### ► What is GenBank?

GenBank<sup>®</sup> is the NIH genetic sequence database, an annotated collection of all publicly available DNA sequences (*Nucleic Acids Research*, 2008 Jan;36(Database issue):D25-30). There are approximately 85,759,586,764 bases in 82,853,685 sequence records in the traditional GenBank divisions and 108,635,736,141 bases in 27,439,206 sequence records in the WGS division as of February 2008.

The complete [release notes](#) for the current version of GenBank are available on the NCBI ftp site. A new release is made every two months. GenBank is part of the [International Nucleotide Sequence Database Collaboration](#), which comprises the DNA DataBank of Japan (DDBJ), the European Molecular Biology Laboratory (EMBL), and GenBank at NCBI. These three organizations exchange data on a daily basis.

An example of a GenBank [record](#) may be viewed for a *Saccharomyces cerevisiae* gene.

### ► In The News: Platypus Genome

Explore Platypus Genome resources.

- [Platypus Genome Project](#)
- [Platypus Taxonomic and Sequence Resources](#)
- [Platypus Genome Resource Guide](#)
- [Duck-Billed Platypus Genome Sequence Published](#) (NIH Press Release)



# EMBL-EBI



The home for big data in biology

Our unique Search service helps you explore dozens of biological data resources.

[More about EBI Search](#)

All

map kinase



Example searches: [blast keratin bfl1](#)

[Find a tool](#)

[Deposit data](#)

## 18% webpage

Usage via web interfaces in 2016.

[More about EMBL-EBI's impact in our annual report](#)

Data from 2016

Find a tool for your data analysis.

Share your scientific data with the world.

## We are EMBL-EBI

The European Bioinformatics Institute (EMBL-EBI) is part of EMBL, Europe's flagship laboratory for the life sciences. [More about EMBL-EBI and our impact](#)

## Data resources

Explore our open data resources to enrich your research. Browse data, perform analyses or share your own results.

## Research

Find out about our research groups, postdoctoral schemes and PhD Programme



Accession [DNA](#) [Protein](#) [Taxonomy](#) [Site Search](#)

Accession numbers

[DDBJ](#)  [UnProt](#)  [PDB](#)  [DAD](#)  [PRF](#)  [Patent](#) [more](#)

HOME Submission How to Use **Search/Analysis** FTP/WebAPI Report/Statistics Contact Us [Japanese](#)

[HOME](#) > Search and Analysis

**Notice**

[Termination of a part of DDJ services \(2008.8.15\)](#)  
[Termination of providing SRS \(Sequence Retrieval System\) services \(2008.8.15\)](#)

**1 Database Search**

[getentry](#)  
Data retrieval by accession numbers, etc.

[ARSA](#)  
All-round Retrieval of Sequence and Annotation

[SRS](#)  
Data retrieval by key words  
**Terminated the service as of Dec. 26, 2008.**

[TXSearch](#)  
Retrieval of unified taxonomy database

[Homology Search](#)  
[FASTA](#) | [BLAST](#) | [PSI-BLAST](#) | [SEARCH](#)

[HMMPFAM](#)  
A motif search program against Pfam on the basis of Hidden Markov Model (HMM)

[DDBJ Vector Screening System](#)

[SQmatch](#)  
Matching a regular expression against sequences in Database

**2 Genome Analyses**

[GIB](#)  
Genome information broker

[GIB-V](#)  
GIB for Viruses

[GIPS](#)  
Reannotation of bacterial genomes using a new common protocol

[GTOP](#)  
Genome to protein structure and function

[H-Invitational Database CIB-DDJ Flat File Server](#)  
**suspended for the system maintenance**

[Mirroring H-InvDB](#)  
Database of full-length cDNAs assigned functional annotation as a result of H-Invitational

**3 Protein Database and Structure**

[PMD](#)  
Protein mutant database

# Banques de données de séquences protéiques

# Origine des informations sur les protéines

- Traduction de séquences d'ADN
  - \* **TrEMBL** traduction automatique de EMBL
  - \* **Genpept** traduction automatique de GenBank
- Séquençage de protéines (Rare car long et coûteux)
- Protéines dont la structure 3D est connue
- Les données stockées : séquences + annotations
  - Protéines entières
  - Fragments de protéines



# Les banques de séquences protéiques

Swiss-Prot + PIR + TrEMBL-EBI



UniProt

(Universal Protein Ressource)

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



# Uniprot



UniProtKB

atpase human

x Advanced

1

BLAST Align Retrieve/ID Mapping

Help Contact

Show help for UniProtKB

b Basket

## Results

Filter by<sup>1</sup>

S Reviewed  
(1,873)  
Swiss-Prot

t Unreviewed  
(40,176)  
TrEMBL

Popular  
organisms

Human (3,063)

Zebrafish (569)

Mouse (219)

S. cerevisiae  
(42)

Rat (16)

Other organisms

Go

Search  
terms

e Columns t BLAST i Align = Download b Add to basket 1 to 25 of 42,049 Show 25

| <input type="checkbox"/> | Entry  | Entry name  |   | Protein names                          | Gene names   | Organism             | Length | e |
|--------------------------|--------|-------------|---|--|--|----------------------|--------|---|
| <input type="checkbox"/> | P00846 | ATP6_HUMAN  | S | ATP synthase subunit a                 | MT-ATP6, ATP6, ATPASE6, MTATP6                     | Homo sapiens (Human) | 226    |   |
| <input type="checkbox"/> | P03928 | ATP8_HUMAN  | S | ATP synthase protein 8                 | MT-ATP8, ATP8, ATPASE8, MTATP8                     | Homo sapiens (Human) | 68     |   |
| <input type="checkbox"/> | P25685 | DNJB1_HUMAN | S | Dnaj homolog subfamily B member 1      | DNAJB1, DNAJ1, HDJ1, HSPF1                         | Homo sapiens (Human) | 340    |   |
| <input type="checkbox"/> | Q9UBS4 | DJB11_HUMAN | S | Dnaj homolog subfamily B member 11     | DNAJB11, EDJ, ERJ3, HDJ9, PSEC0121, UNQ537/PRO1080 | Homo sapiens (Human) | 358    |   |
| <input type="checkbox"/> | Q15645 | PCH2_HUMAN  | S | Pachytene checkpoint protein 2 homo... | TRIP13, PCH2                                       | Homo sapiens (Human) | 432    |   |
| <input type="checkbox"/> | Q9Y2G3 | AT11B_HUMAN | S | Probable phospholipid-transporting ... | ATP11B, ATPIF, ATPIR, KIAA0956                     | Homo sapiens         | 1,177  |   |

Protein  
Gene  
Organism  
Status

# Sodium/potassium-transporting ATPase subunit alpha-1

ATP1A1

*Homo sapiens (Human)*

**S** Reviewed - ●●●●● - Experimental evidence at protein level<sup>i</sup>

Display **None**

[t BLAST](#) [i Align](#) [b Format](#) [b Add to basket](#) [{ History](#)

[U Feedback](#) [V Help video](#)

- FUNCTION
- NAMES & TAXONOMY
- SUBCELL. LOCATION
- PATHOL./BIOTECH
- PTM / PROCESSING
- EXPRESSION

## Function<sup>i</sup>

This is the catalytic component of the active enzyme, which catalyzes the hydrolysis of ATP coupled with the exchange of sodium and potassium ions across the plasma membrane. This action creates the electrochemical gradient of sodium and potassium ions, providing the energy for active transport of various nutrients.

### Catalytic activity<sup>i</sup>



### Sites

- INTERACTION
- STRUCTURE
- FAMILY & DOMAINS
- SEQUENCES (4)
- CROSS-REFERENCING
- PUBLICATIONS
- ENTRY INFORMATION
- MISCELLANEOUS

| Feature key                | Position(s) | Length | Description   | Graphical view | Feature identifier | Actions |
|----------------------------|-------------|--------|---|----------------|--------------------|---------|
| Active site <sup>i</sup>   | 376 - 376   | 1      | 4-aspartylphosphate intermediate<br><a href="#">By similarity</a> |                |                    |         |
| Binding site <sup>i</sup>  | 487 - 487   | 1      | ATP <a href="#">By similarity</a>                                 |                |                    |         |
| Metal binding <sup>i</sup> | 717 - 717   | 1      | Magnesium<br><a href="#">By similarity</a>                        |                |                    |         |
| Metal binding <sup>i</sup> | 721 - 721   | 1      | Magnesium   |                |                    |         |

# Banques de structure 3D de macromolécules ( PDB ) ;

PDB (Protein DataBank)

<http://www.rcsb.org>

- Séquences et structures 3D des protéines et des acides nucléiques macromolécules .
- Visualisation en 3D .

RCSB PDB Deposit Search Visualize Analyze Download Learn More Documentation Careers MyPDB

183793 Biological Macromolecular Structures Enabling Breakthroughs In Research and Education

Enter search terms or PDB ID(s)

Advanced Search | Browse Annotations

Help

Celebrating 50 YEARS OF Protein Data Bank

PDB-101 PDB EMBL Data Resource Wellcome Open Research WorldWide Protein Data Bank Foundation

Facebook Twitter YouTube

Developers: Join the RCSB PDB Team Explore Open Positions

Welcome

Deposit

Search

Visualize

Analyze

Download

## A Structural View of Biology

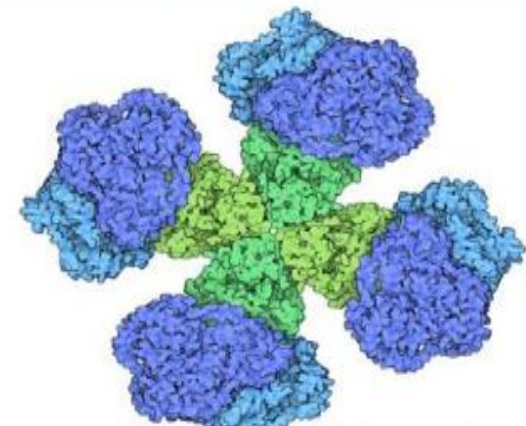
This resource is powered by the Protein Data Bank archive-information about the 3D shapes of proteins, nucleic acids, and complex assemblies that helps students and researchers understand all aspects of biomedicine and agriculture, from protein synthesis to health and disease.

As a member of the wwPDB, the RCSB PDB curates and annotates PDB data.

The RCSB PDB builds upon the data by creating tools and resources for research and education in molecular biology, structural biology, computational biology, and beyond.



## November Molecule of the Month



# PDB

- PDB (BrookHaven Protein DataBank)
  - <http://www.rcsb.org>
  - **Séquences et structures des protéines**
  - Visualisation en 3D
  - Les données proviennent de cristallographie, de RMN,...
  - Pour certaines protéines, plusieurs structures sont disponibles
    - Structure de la protéine seule ou avec ligand
    - Structure de la protéine dans différents milieux
    - Structure obtenue avec des méthodes expérimentales différentes

**WHAT'S NEW** | [HELP](#) | [PRINT](#)

PDB ID or keyword

map kinase

**Search**

**Advanced Search**

[Home](#) [Hide](#)

[News & Publications](#)  
[Policies](#)  
[FAQ](#)  
[Contact](#)  
[Feedback](#)  
[About Us](#)

[Deposition](#) [Hide](#)

[All Deposit Services](#)  
[Electron Microscopy](#)  
[NMR](#)  
[Validation Server](#)  
[BioSync Beamline](#)  
[Related Tools](#)

[Search](#) [Hide](#)

[Advanced Search](#)  
[Latest Release](#)  
[Latest Publications](#)  
[Sequence Search](#)  
[Ligand Search](#)  
[Unreleased Entries](#)  
[Browse Database](#)  
[Histograms](#)

Explorer:

**315 Structure Hits**

[9 Unreleased Structures](#)

[158 Citations](#)

[231 Ligand Hits](#)

[107 Web Page Hits](#)

[GO Hits](#)

[SCOP Hits](#)

[CATH Hits](#)

Advanced Keyword Query for: MAP KINASE

Query Options:

Display/Download:

Generate Reports:

Sort by:

Results per Page:

Displaying results 1 - 10 of 315 total | Page 1 of 32



**High resolution crystal structure of mitogen-activated protein kinase-activated protein kinase 3/inhibitor 2 complex**

**Characteristics**

Release Date: 26-Jan-2010 Exp. Method: X-RAY DIFFRACTION

Resolution: 2.00 Å

**Classification**

**Transferase**

**Compound**

**Molecule:** MAP kinase-activated protein kinase 3

**Polymer:** 1 **Type:** polypeptide(L)

**Length:** 336

**Chains:** A

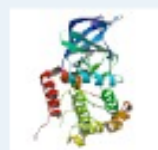
**EC#:** [2.7.11.1](#)

**Fragment:** Kinase domain, UNP residues 33-349

**Authors**

[Cheng, R.K.Y.](#), [Barker, J.](#), [Palan, S.](#), [Felicetti, B.](#),

[Whittaker, M.](#), [Hesterkamp, T.](#)



**Crystal structure of MAPKAP kinase 2 (MK2) complexed with a potent 3-aminopyrazole ATP site inhibitor**

Crystal structure of MAPKAP kinase 2 (MK2) complexed with a potent 3-aminopyrazole ATP site inhibitor

3KGA

- Display Files
- Download Files
- Share this Page

Sequence Display

The sequence display provides a graphical representation of the UniProtKB, PDB - ATOM and PDB - SEQRES sequences. Different 3rd party annotations can be graphically mapped on the sequence and displayed in the Jmol viewer.

The structure 3KGA has in total 1 chains.

Currently viewing unique chains only. [show all chains](#)

Chain A : MAP kinase-activated protein kinase 2

FASTA | [Sequence & DSSP](#) | [Image](#)

Polymer 1  
 Length: 299 residues  
 Chain Type: polypeptide(L)  
 Reference: [UniProtKB P49137](#)

Sequence & Structure Relationships

- Display Jmol
- Enable Jmol to view annotations in 3D.

Display Parameters

Currently displayed: SEQRES sequence.  
[Display external \(UniProtKB\) sequence](#)  
 Mouse over an annotation to see more details. Click annotation to enable jmol.

Annotations

Add Annotations

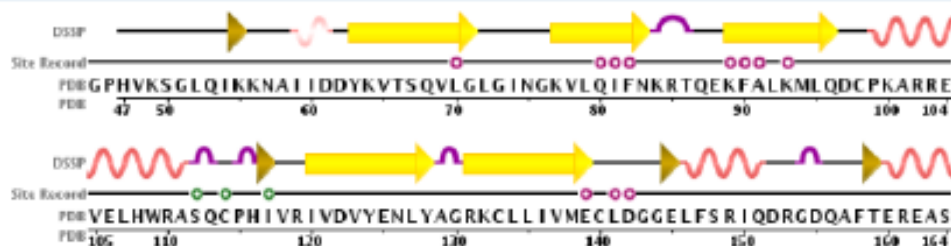
Select

Secondary Structure:DSSP  
[\[hide\]](#) [\[reference\]](#)

38% helical (15 helices: 116 residues)  
 19% beta sheet (14 strands: 57 residues)

Structural Feature:Site Record  
[\[hide\]](#) [\[reference\]](#)

**3KGA\_A\_AC2\_16** BINDING SITE FOR RESIDUE LX9 A 365 (SOFTWARE)  
**3KGA\_A\_AC1\_4** BINDING SITE FOR RESIDUE MG A 1 (SOFTWARE)



PDB :  
 structure  
 secondaires

# PDB : séquence des protéines

Summary **Sequence** Annotations Seq. Similarity 3D Similarity Literature Biol. & Chem. Methods Geometry Links

Crystal structure of MAPKAP kinase 2 (MK2) complexed with a potent 3-aminopyrazole ATP site inhibitor

**3KGA**

Display Files  
Download Files

Share this Page

## Sequence Display

The sequence display provides a graphical representation of the UniProtKB, PDB - ATOM and PDB - SEQRES sequences. Different 3rd party annotations can be graphically mapped on the sequence and displayed in the Jmol viewer.

The structure **3KGA** has in total **1** chains.

Currently viewing **unique chains** only. [show all chains](#)

**Chain A** : MAP kinase-activated protein kinase 2

**FASTA** | [Sequence & DSSP](#) | [Image](#)

Polymer 1

Length: 299 residues

Chain Type: polypeptide(L)

Reference: [UniProtKB P49137](#)

Sequence & Structure Relationships

Display Jmol

Enable Jmol to view annotations in 3D.

Display Parameters

Currently displayed: **SEQRES** sequence.  
[Display external \(UniProtKB\) sequence](#)

3KGA\_A.fasta.txt [Lecture seule] (/tmp) - gedit

Fichier Édition Affichage Rechercher Outils Documents Aide

Ouvrir Enregistrer Annuler

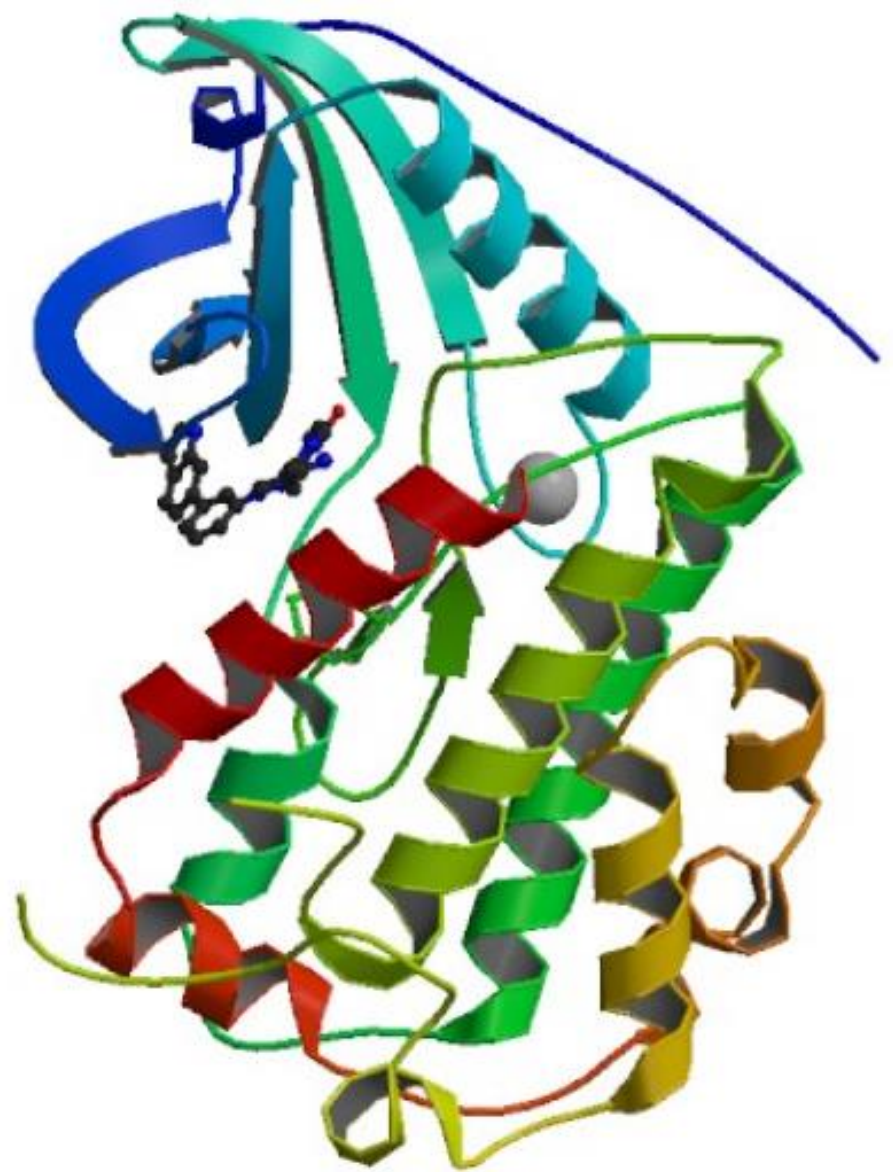
3KGA\_A.fasta.txt

3KGA : A | PDBID | CHAIN | SEQUENCE

```
GPHVKSG LQIKKNAI IDDYKYVTSQVLGLGINGKVLQIFNKRTQEKFALKMLQDCPKARREVELHWRASQCPHIVRIVDVY  
ENLYAGRKCLLIYMECLDGGELFSRIQDRGDQAFTEREASEIMKSIGEAIQYLHSINIAHRDVKPENLLYTSKRPNAILK  
LTDGFAKETTGEKYDKSCDMWSLGVIMYILLCGYPPFYSNHGLAISP GMKTRIRMGQYEFNPPEWSEVSEEVKMLIRNL  
LKTEPTQRMTITEFMNHPWIMQSTKVPQTPLHTSRVLKEDKERWEDVKEEMTSALATMR
```

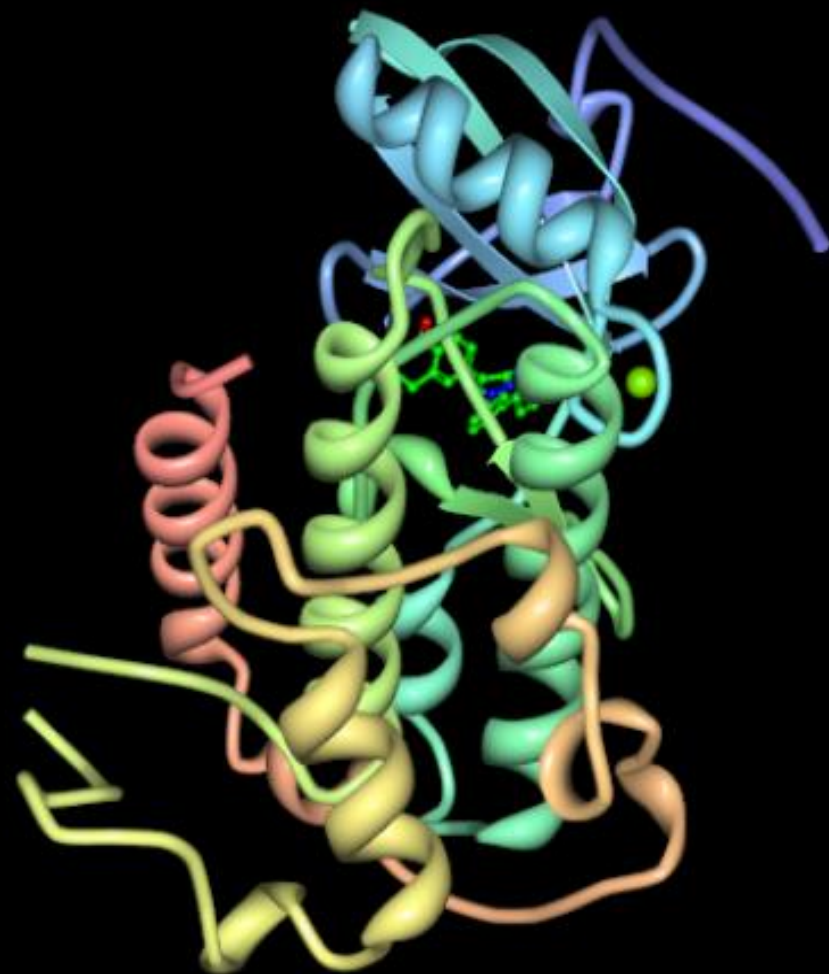


# PDB : structure tertiaires



## Biological Assembly Image for 3KGA

Crystal structure of MAPKAP kinase 2 (MK2) complexed with a potent 3-aminopyrazole ATP site inhibitor



**PDB 3KGA: Crystal structure of MAPKAP kinase 2 (MK2) complexed with a potent 3-aminopyrazole ATP site inhibitor**

Velcicky, J., Feifel, R., Hawtin, S., Heng, R., Huppertz, C., et al.  
Bioorg.Med.Chem.Lett. (2009)

## Crystal structure of MAPKAP kinase 2 (MK2) complexed with a potent 3-aminopyrazole ATP site inhibitor

# 3KGA

-  Display Files ▾
-  Download Files ▾
-  Print this Page
-  Share this Page

Primary Citation and Related Literature

### Primary Citation

**Novel 3-aminopyrazole inhibitors of MK-2 discovered by scaffold hopping strategy.**

[Velcicky, J.](#), [Feifel, R.](#), [Hawtin, S.](#), [Heng, R.](#), [Huppertz, C.](#), [Koch, G.](#), [Kroemer, M.](#), [Moebitz, H.](#), [Revesz, L.](#), [Scheufler, C.](#), [Schlapbach, A.](#)

(2009) Bioorg.Med.Chem.Lett.

**PubMed:** [20060294](#) 

**DOI:** [10.1016/j.bmcl.2009.10.138](#) 

[Search Related Articles in PubMed](#) 



### PubMed Abstract:

New, selective 3-aminopyrazole based MK2-inhibitors were discovered by scaffold hopping strategy. The new derivatives proved to inhibit intracellular phosphorylation of hsp27 as well as LPS-induced TNFalpha release in cells. In addition, selected derivative 14e also inhibited LPS-induced TNFalpha release in ... [\[ Read More & Search PubMed Abstracts \]](#)

### Publication Details

#### MeSH Terms (Primary Citation)

#### Literature Network

iHOP: [9261](#)  GenerIF: [9261](#) 

Information provided by [BioLit](#):

### PubMedCentral articles found to contain 3KGA

No related articles could be found in PubMedCentral.

# Bases des données sur les génomes

The image shows the NCBI Entrez Genome Project database homepage. At the top, the NCBI logo is on the left, and the 'ENTREZ Genome Project' title is in the center. To the right, there is a 'connection information discovery' logo and a 'My NCBI' button with 'Sign In' and 'Register' links. Below the header is a navigation bar with tabs for 'All Databases', 'PubMed', 'Nucleotide', 'Protein', 'Genome', 'Structure', 'PMC', 'Taxonomy', and 'Books'. A search bar contains the text 'Genome Project' and has 'Go' and 'Clear' buttons. Below the search bar are buttons for 'Limits', 'Preview/Index', 'History', 'Clipboard', and 'Details'. On the left side, there is a vertical menu with links for 'About Entrez', 'Entrez Genome Project', 'Home', 'Overview', 'Help', 'Statistics', 'Sequencing Centers', 'Submitting', 'Project Submissions', 'Project Instructions', 'General Genome Submissions', 'Feature Tables', 'Bacterial Genome Submissions', 'Metagenome Submissions', and 'Whole Genome Shotgun Sequences'. The main content area features a welcome message: 'Welcome to the NCBI Entrez Genome Project database. This searchable database is a collection of complete and incomplete large-scale sequencing, assembly, annotation, and mapping projects for cellular organisms. The database is organized into organism-specific overviews that function as portals from which all projects in the database pertaining to that organism can be browsed and retrieved. [Read more...](#)' Below this is a large graphic showing a tree-like structure of biological groups: 'ANIMALS' (Mammals, Insects, Amphibians, Birds, Flatworms, Reptiles, Fishes, Roundworms, Other), 'PLANTS' (Green Algae, Land Plants), 'FUNGI' (Ascomycetes, Basidiomycetes, Other), 'PROTISTS' (Apicomplexans, Kinetoplasts, Other), and 'EUKARYOTES'. To the right of the main content is a 'NCBI Resources' section with links to various databases and tools: Entrez Gene, Entrez Genome, Entrez Protein Clusters, Metagenomic Projects, Eukaryotic Projects, Genomic Biology, Prokaryotic Projects, Organellar Genomes, Plant Genomes, and RefSeq.

NCBI

ENTREZ Genome Project

connection information discovery

My NCBI  
[Sign In](#) [Register](#)

All Databases PubMed Nucleotide Protein Genome Structure PMC Taxonomy Books

Search Genome Project for

Limits Preview/Index History Clipboard Details

About Entrez

Entrez Genome Project

Home Overview Help Statistics Sequencing Centers

Submitting

Project Submissions Project Instructions General Genome Submissions Feature Tables Bacterial Genome Submissions Metagenome Submissions Whole Genome Shotgun Sequences

Welcome to the NCBI Entrez Genome Project database. This searchable database is a collection of complete and incomplete large-scale sequencing, assembly, annotation, and mapping projects for cellular organisms. The database is organized into organism-specific overviews that function as portals from which all projects in the database pertaining to that organism can be browsed and retrieved. [Read more...](#)

**NCBI Resources**

- [Entrez Gene](#) gene-related information
- [Entrez Genome](#) sequence and map data from whole genomes
- [Entrez Protein Clusters](#) a collection of related protein sequences
- [Metagenomic Projects](#) metagenomic-specific genome projects
- [Eukaryotic Projects](#) eukaryotic-specific genome projects
- [Genomic Biology](#) organism-specific links
- [Prokaryotic Projects](#) prokaryotic-specific genome projects
- [Organellar Genomes](#) organellar reference sequences and tools
- [Plant Genomes](#) major plant genome projects
- [RefSeq](#) the reference sequence project
- [Viral Genomes](#) viral reference sequences and tools

# Banques de données spécialisées

# Les banques de données spécialisées

- ❑ Ces banques contiennent des données **homogènes**
  - Collecte, établie autour d'une thématique particulière
- ❑ **Avantages** : facilité pour mettre à jour les données, vérifier leur intégrité, offrir une interface adaptée, ...
- ❑ **Inconvénients** : ne cible pas toujours ce que l'on veut; toutes les banques possibles n'existent pas
- ❑ **Exemples** : banques spécialisées pour un génome, banques de séquences d'immunologies, banques sur des séquences validées, ...

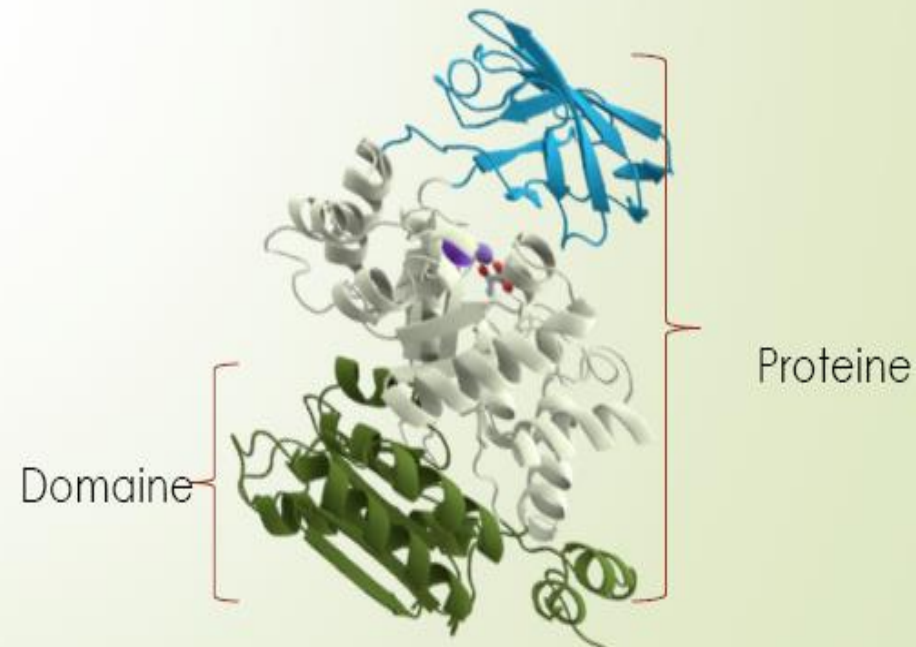
## Quelques exemples de Bases de données spécialisés

## 1- Base de données PROSITE :

<http://www.expasy.ch/prosite/>

**définition** : base de données sur les domaines des protéines, les familles protéiques et les fonctions biologiques associées .

→ Un domaine = une région d'une protéine ayant une fonction biologique propre, que l'on retrouve sur plusieurs protéines.







## Database of protein domains, families and functional sites



SARS-CoV-2 relevant PROSITE motifs

PROSITE consists of documentation entries describing protein domains, families and functional sites as well as associated patterns and profiles to identify them [[More...](#) / [References](#) / [Commercial users](#)].

PROSITE is complemented by [ProRule](#), a collection of rules based on profiles and patterns, which increases the discriminatory power of profiles and patterns by providing additional information about functionally and/or structurally critical amino acids [[More...](#)].

**Release 2021\_04 of 29-Sep-2021 contains 1895 documentation entries, 1311 patterns, 1326 profiles and 1338 ProRule.**

### Search

 e.g. PDOC00022, PS50089, SH3, zinc finger

### Browse

- by documentation entry
- by ProRule description
- by taxonomic scope
- by number of positive hits

[Quick Scan mode of ScanProsite](#)[Other tools](#)

# PROSITE



Swiss Institute of  
Bioinformatics



Search PROSITE

for

Go

Clear

ExPASy Proteomics Server

Databases Tools Services Mirrors About Contact

You are here: ExPASy CH > Databases > PROSITE

Search in PROSITE for: atpase

(Release 20.60, of 09-Feb-2010 )

Enter search terms:

atpase

Prefix and append wildcard '\*' to words.

new search

clear

By default, this search engine searches for complete words only. If you did not find what you expected, and would try to do a substring match, you should perform a new search and select 'prefix and append wildcard to words'.

Number of documents in PROSITE containing the search term:47

- [PDOC00572](#) AAA-protein family signature
- [PDOC00364](#) ABC transporter integral membrane type-1 domain profiles
- [PDOC00420](#) ATP synthase a subunit signature
- [PDOC00137](#) ATP synthase alpha and beta subunits signature
- [PDOC00526](#) ATP synthase c subunit signature
- [PDOC00327](#) ATP synthase delta (OSCP) subunit signature
- [PDOC00138](#) ATP synthase gamma subunit signature
- [PDOC00185](#) ATP-binding cassette, ABC transporter-type, signature and profile
- [PDOC00017](#) ATP/GTP-binding site motif A (P-loop)
- [PDOC00661](#) ArsR-type HTH domain signature and profile
- [PDOC00610](#) Chaperonins TCP-1 signatures
- [PDOC00576](#) Chaperonins cpn10 signature
- [PDOC00268](#) Chaperonins cpn60 signature
- [PDOC00196](#) Clathrin light chains signatures
- [PDOC51413](#) DBINO domain profile

Identifiants :

PDOC... => domaine

P... => protéine

# PROSITE



Swiss Institute of  
Bioinformatics



Search  for

ExpASY Proteomics Server

[Databases](#) [Tools](#) [Services](#) [Mirrors](#) [About](#) [Contact](#)

You are here: [ExpASY](#) [CH](#)

## ABC transporter integral membrane type-1 domain profiles

### Description:

ABC transporters belong to the ATP-Binding Cassette (ABC) superfamily which uses the hydrolysis of ATP to energize diverse biological import and export systems (see <[PDOC00185](#)>). ABC transporters are minimally constituted of two conserved regions: a highly conserved ATP binding cassette (ABC) and a less conserved transmembrane domain (TMD). These regions can be found on the same protein (mostly in eukaryotes and bacterial exporters) or on two different ones (mostly bacterial importers) [[1,2,3](#)]. The function of the integral inner-membrane protein is to translocate the substrate across the membrane. Studies of P-glycoprotein function indicate that residues lining the proposed chamber opening (residues of TM2, TM5 and TM6) play an important role in substrate recognition [[4](#)].

In exporters and eukaryotes, ABC transporters consist of a single polypeptide composed of an N-terminal domain of approximately 320 residues, apparently containing six transmembrane segments, fused to a highly conserved ABC-ATPase domain of approximately 260 residues [[5,6,7](#)]. In some cases an N-terminal peptidase domain of 130-150 residues appended to the TMD is also found, which may contain additional transmembrane segments as in the HlyB subfamily [[8,9](#)].

The 3D structure of the E. coli lipid A flippase MsbA homodimer reveals that association of the two transmembrane domains forms one chamber that adopt a cone-shape which extends along a pseudo two-fold axis perpendicular to the cell membrane (see <[PDB:1JSQ](#)>) [[10](#)]. The chamber has an opening on either side of the membrane to provide free access for the lipid substrate from the cytoplasmic leaflet of the lipid bilayer, while excluding molecules from the outer leaflet. The chamber openings are defined by intramolecular interactions between TM2 of one monomer and TM5 of the other. The residues lining the chamber are contributed by all 12 transmembrane  $\alpha$ -helices [[10,11](#)].

In importers (found only in prokaryotes or archaea) most ABC transporters consist of four domains usually encoded by independent polypeptides, two ABC modules and two TMDs which are thought to contain six transmembrane regions [[12,13](#)]. The approximately 30 kD TMD displays a distinctive signature, the EAA motif, a 20 amino acid conserved sequence located about 100 residues from the C-terminus. The motif is hydrophilic and has been found to reside in a cytoplasmic loop located between the penultimate and the antepenultimate transmembrane segment in all proteins with a known topology [[14,15,16](#)]. It appears to play an important role in ensuring the correct assembly of the prokaryotic ABC transport complex [[17](#)] and constituting an interaction site with the so-called helical domain of the ABC module [[18,19](#)]. The TMDs form either homo-oligomeric channels or associate with another TMD to form hetero-oligomers.

**InterPro** : est une base de données utilisée pour la classification et l'annotation automatique de protéines!

# Base spécialisées dans les domaines de proteins

## INTERPRO

<http://www.ebi.ac.uk/interpro/>

The screenshot shows the InterPro website interface. At the top, there is a navigation bar with the EMBL-EBI logo, an "EB-eye Search" dropdown menu set to "All Databases", a search input field with the placeholder "Enter Text Here", and buttons for "Go", "Reset", "Advanced Search", and "Give us feedback". Below the navigation bar are tabs for "Databases", "Tools", "EBI Groups", "Training", "Industry", "About Us", and "Help", along with a "Site Index" link and social media icons.

The main content area is titled "InterPro: Home" and includes a search bar with the text "Search InterPro:". Below the search bar, there is a description of InterPro: "InterPro is a database of protein families, domains, repeats and sites in which identifiable features found in known proteins can be applied to new protein sequences."

On the left side, there is a grid of logos for various databases: UniProt (Universal Protein Resource), proSite, Pfam, PRINTS (Protein Family and Repeat Identification), ProDom, SMART, TIGR (The Institute for Genome Research and Biotechnology), and PANTHER (Protein Classification System).

Under the "Release News" section, there is an "Announcement:" with the following bullet points:

- **InterPro 18.0 is released** and covers 75.6% of UniProtKB, with new methods from PROSITE, GENE3D and SUPERFAMILY.
- **PROSITE pattern matches** are now evaluated to either TRUE (T) or UNKNOWN (?) using miniprofiles or associated existing PROSITE profiles.

# Base de données dédiée au métabolisme

## KEGG

<http://www.genome.jp/kegg/>



### KEGG PATHWAY Database

Wiring diagrams of molecular interactions, reactions, and relations

KEGG2

ATLAS

PATHWAY

BRITE

GENES

SSDB

LIGAND

DBGET

#### Pathway Maps

**KEGG PATHWAY** is a collection of manually drawn pathway maps representing our knowledge on the molecular interaction and reaction networks for:

**1. Metabolism**

Carbohydrate Energy Lipid Nucleotide Amino acid Other amino acid  
Glycan PK/NRP Cofactor/vitamin Secondary metabolite Xenobiotics

**2. Genetic Information Processing**

**3. Environmental Information Processing**

**4. Cellular Processes**

**5. Human Diseases**

and also on the structure relationships (KEGG drug structure maps) in:

**6. Drug Development**

#### Pathway Modules

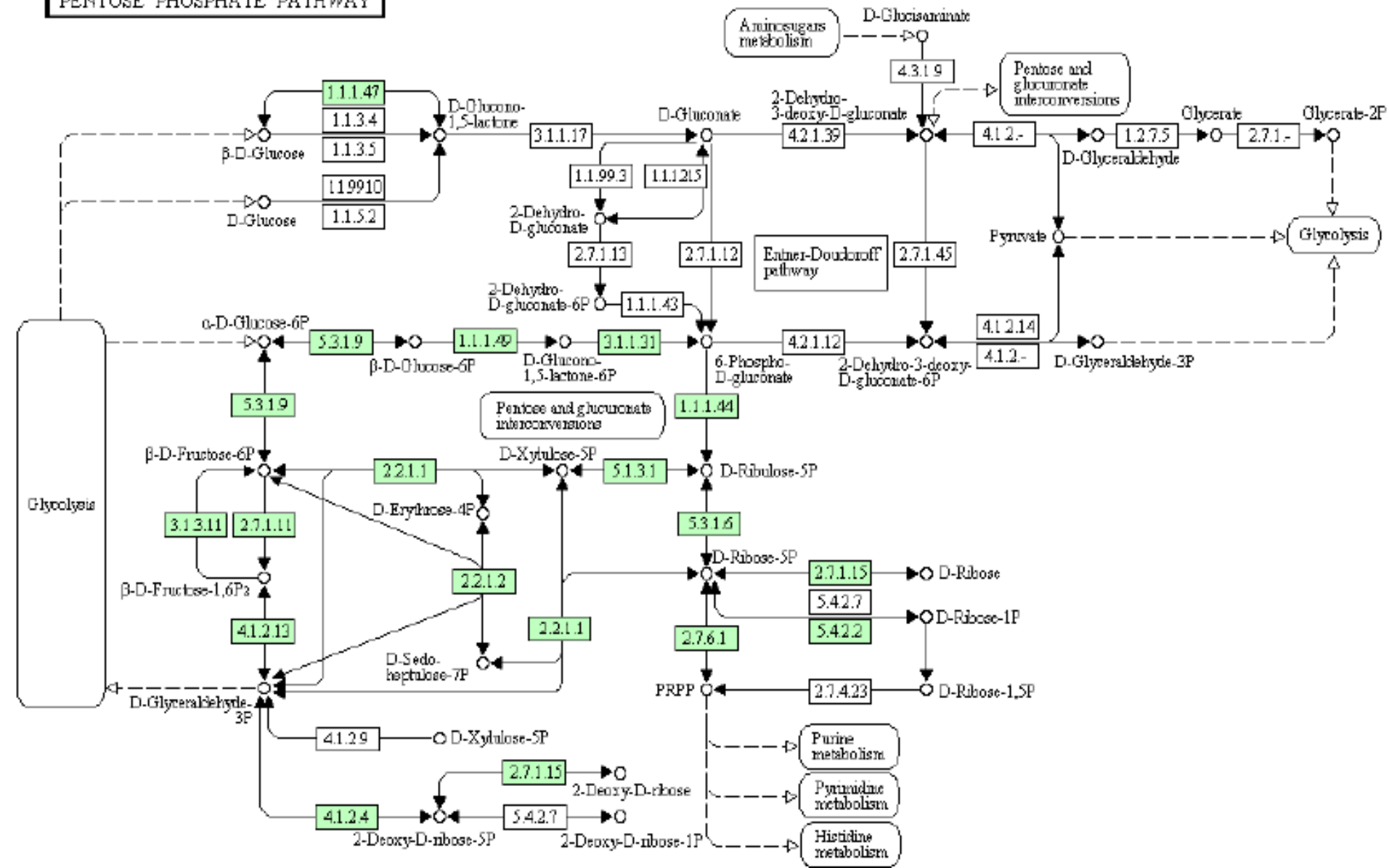
**KEGG MODULE** is a new collection of pathway modules, molecular complexes, and other functional units, each represented as a list of KEGG Orthology (KO) identifiers. KEGG MODULE can be accessed through a BRITE hierarchy:

**KEGG pathway modules**

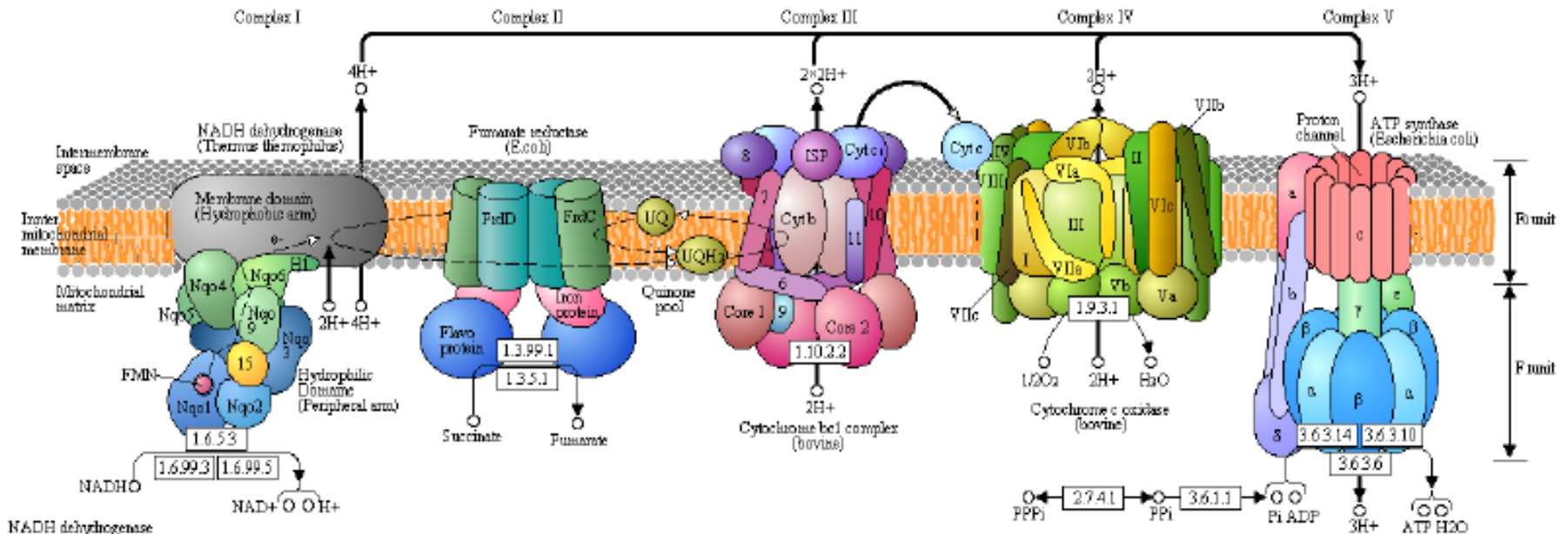
or by the DBGET search.

# Exemple de voies métaboliques dans KEG

## PENTOSE PHOSPHATE PATHWAY



# OXIDATIVE PHOSPHORYLATION



## NADH dehydrogenase

|   |     |     |     |     |      |     |     |
|---|-----|-----|-----|-----|------|-----|-----|
| B | ND1 | ND2 | ND3 | ND4 | ND4L | ND5 | ND6 |
|---|-----|-----|-----|-----|------|-----|-----|

|   |        |        |        |        |        |        |        |        |        |         |         |         |
|---|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|
| B | Ndubf1 | Ndubf2 | Ndubf3 | Ndubf4 | Ndubf5 | Ndubf6 | Ndubf7 | Ndubf8 | Ndubf9 | Ndubf10 | Ndubf11 | Ndubf12 |
|---|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|

|     |      |      |      |      |      |      |      |      |      |      |      |      |      |      |
|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| B/A | NunA | NunB | NunC | NunD | NunE | NunF | NunG | NunH | NunI | NunJ | NunK | NunL | NunM | NunN |
|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|

|     |      |      |      |      |      |      |      |      |      |      |      |      |      |      |      |      |      |
|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|
| B/A | NdaC | NdaK | NdaN | NdaH | NdaA | NdaI | NdaG | NdaE | NdaP | NdaD | NdaB | NdaL | NdaM | NdaN | HoxE | HoxF | HoxU |
|-----|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|------|

|   |        |        |        |        |        |        |        |        |        |         |         |         |
|---|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|
| E | Ndubf1 | Ndubf2 | Ndubf3 | Ndubf4 | Ndubf5 | Ndubf6 | Ndubf7 | Ndubf8 | Ndubf9 | Ndubf10 | Ndubf11 | Ndubf12 |
|---|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|

|   |        |        |        |        |        |        |        |        |        |         |         |         |
|---|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|
| B | Ndubf1 | Ndubf2 | Ndubf3 | Ndubf4 | Ndubf5 | Ndubf6 | Ndubf7 | Ndubf8 | Ndubf9 | Ndubf10 | Ndubf11 | Ndubf12 |
|---|--------|--------|--------|--------|--------|--------|--------|--------|--------|---------|---------|---------|

## Succinate dehydrogenase / Fumarate reductase

|   |      |      |      |      |
|---|------|------|------|------|
| E | SDHC | SDHD | SDHA | SDHB |
|---|------|------|------|------|

|     |      |      |      |      |
|-----|------|------|------|------|
| B/A | SdhC | SdhD | SdhA | SdhB |
|     | FcdA | FcdB | FcdC | FcdD |

## Cytochrome c oxidase

|   |       |      |      |      |      |       |       |       |       |       |       |       |       |      |       |       |       |       |
|---|-------|------|------|------|------|-------|-------|-------|-------|-------|-------|-------|-------|------|-------|-------|-------|-------|
| E | COX10 | COX3 | COX1 | COX2 | COX4 | COX5A | COX5B | COX6A | COX6B | COX6C | COX7A | COX7B | COX7C | COX8 | E/B/A | COX11 | COX15 | COX17 |
|---|-------|------|------|------|------|-------|-------|-------|-------|-------|-------|-------|-------|------|-------|-------|-------|-------|

|     |      |      |      |      |      |
|-----|------|------|------|------|------|
| B/A | CyoE | CyoD | CyoC | CyoB | CyoA |
|     | CoxD | CoxC | CoxA | CoxB |      |
|     | QoxD | QoxC | QoxB | QoxA |      |

## Cytochrome c reductase

|       |      |      |      |      |      |      |       |
|-------|------|------|------|------|------|------|-------|
| E/B/A | ISP  | Cytb | Cyt1 |      |      |      |       |
| B     | COR1 | QCR2 | QCR6 | QCR7 | QCR8 | QCR9 | QCR10 |

## Cytochrome c oxidase, cbb3-type

|   |   |    |    |     |
|---|---|----|----|-----|
| B | I | II | IV | III |
|---|---|----|----|-----|

## Cytochrome bd complex

|     |      |      |
|-----|------|------|
| B/A | CydA | CydB |
|-----|------|------|

## F-type ATPase (Bacteria)

|      |       |       |       |         |   |   |   |
|------|-------|-------|-------|---------|---|---|---|
| beta | alpha | gamma | delta | epsilon | c | a | b |
|------|-------|-------|-------|---------|---|---|---|

## F-type ATPase (Eukaryotes)

|      |       |       |      |       |         |   |   |
|------|-------|-------|------|-------|---------|---|---|
| beta | alpha | gamma | OSCP | delta | epsilon | c | a |
| b    | e     | fb    | f    | s     |         |   |   |
| d    | f     | h     | j    | k     | g       |   |   |

## F-type ATPase (Prokaryotes)

|   |   |   |   |   |   |   |   |
|---|---|---|---|---|---|---|---|
| A | B | C | D | E | F | I | K |
|---|---|---|---|---|---|---|---|

## F-type ATPase (Eukaryotes)

|   |      |      |    |       |   |   |   |
|---|------|------|----|-------|---|---|---|
| A | B    | C    | D  | E     | F | G | H |
| I | AC39 | 54kD | S1 | lipid |   |   |   |

# Base consacrée uniquement sur *Escherichia coli*

## ECOCYC

<http://ecocyc.org/>



### EcoCyc™

Encyclopedia of *Escherichia coli* K-12 Genes and Metabolism

#### EcoCyc Home

#### Quick Search

[Database Search](#)  
[Advanced Database Search](#)  
[BLAST](#)

#### Browse

[Pathways](#)  
[Genes](#)  
[Genome Browser](#)  
[Reactions](#)  
[Compounds](#)  
[Metabolic Chart](#)  
[Omics Viewer](#)

#### About EcoCyc

[Project Overview](#)  
[Guided Tour](#)  
[Webinars \(Instructional Videos\)](#)  
[Publications](#)  
[Update History](#)  
[Steering Committee](#)  
[Credits](#)

#### Services

#### Project Overview

EcoCyc is a scientific database for the bacterium *Escherichia coli* K-12 MG1655. The EcoCyc project performs [literature-based curation](#) of the entire genome, and of transcriptional regulation, transporters, and metabolic pathways. [\[project overview\]](#)

#### New Users

Take the [guided tour](#) of the EcoCyc Web site, watch our [free online instructional videos](#), or read our 2007 article: ["Multidimensional annotation of the \*Escherichia coli\* K-12 Genome."](#)

#### New in EcoCyc

**Now available: Gene Ontology annotation file for *E. coli* K-12:** [\[Download from GO site\]](#)

#### A highlight from the 12.5 release:

- [NADH to cytochrome \*bo\* oxidase electron transfer](#). The proton-motive force across the cytoplasmic membrane is essential for life, powering ATP synthesis and the action of proton-driven symporters. As shown in this pathway, two NADH:ubiquinone oxidoreductase and cytochrome *bo* terminal oxidase work together to transfer electrons from NADH to oxygen, using the energy from those electrons to pump protons across the cytoplasmic membrane and generate the proton-motive force. This pathway is one of eleven new electron transfer pathways in the 12.5 release, capitalizing on our recently added ability to represent electron transfer half reactions and combine them to generate pathways. [Click here](#) to learn more about this fundamental pathway of energy generation.

The full EcoCyc release history is available [here](#). You can read past highlights pieces by [clicking here](#).

#### Update Frequency

The EcoCyc Web site and downloadable files are updated quarterly. A faster, more powerful EcoCyc that you can [install locally](#) on your computer (Macintosh, PC/Windows, PC/Linux) is released semiannually. [\[Full EcoCyc release history\]](#)



**Late Embryogenesis Abundant Proteins database (G. Hunault & E. Jaspard) :** cette base de données contient un grand nombre d'informations sur les protéines LEA impliqués dans la tolérance à de nombreux stress, notamment la déshydratation et le froid. Pour l'instant, on les a mises en évidence principalement chez les plantes.

**LEAPDB**

1642 proteins

[Home](#)

[Browse](#)

[Search](#)

[Blast](#)

[Statistical analysis](#)

[Export](#)

[Submit](#)

[Help](#)

[Contacts](#)



AMBIOR

**Late Embryogenesis Abundant Proteins**

LEA proteins have been discovered in 1981. Although, they are almost associated with abiotic stress tolerance (particularly dehydration and cold stress), their actual function remains unknown.

The LEAP database provides usefull data about the different CLASSES of LEA proteins for the analysis of their structure - function relationships. More...

Multi-classification of LEAP with three - 4th reduction of CLASS nomenclature

| #FAI   | Gen. nr. | Src.    | Taxonomy and size | Subtype #1 | Sub-Subtype #1 | Hydrophobic residues | 1.3.1 class  |
|--------|----------|---------|-------------------|------------|----------------|----------------------|--------------|
|        | 198      | 198     | 207               | 208        | 208            | 208                  | 210          |
| PF0281 | 01       | Group 2 | Group 2           | Group 2    | Group 2        | Dehydrin             | Class 1 to 4 |
| PF0287 | 018      | Group 1 | Group 1           | Group 1    | Group 1        | LEA_3                | Class 5      |
|        | 018      | Group 1 | Group 1           | Group 1    | Group 1        | LEA_3                | Class 5      |
| PF0281 | 02       | Group 2 | Group 2           | Group 2    | Group 2        | LEA_4                | Class 6      |
|        | 02       | Group 2 | Group 2           | Group 2    | Group 2        | LEA_4                | Class 6      |
| PF0284 | 04       | Group 3 | Group 3           | Group 3    | Group 3        | LEA_3                | Class 7 to 8 |
|        | 04       | Group 3 | Group 3           | Group 3    | Group 3        | LEA_3                | Class 7 to 8 |
| PF0284 | 07       | Group 4 | Group 4           | Group 4    | Group 4        | LEA_2                | Class 9      |
|        | 07       | Group 4 | Group 4           | Group 4    | Group 4        | LEA_2                | Class 9      |
| PF0284 | 08       | Group 4 | Group 4           | Group 4    | Group 4        | LEA_1                | Class 10     |
|        | 08       | Group 4 | Group 4           | Group 4    | Group 4        | LEA_1                | Class 10     |
| PF0287 | 09       | Group 5 | Group 5           | Group 5    | Group 5        | SNP                  | Class 11     |
|        | 09       | Group 5 | Group 5           | Group 5    | Group 5        | SNP                  | Class 11     |
| PF0284 | 10       | Group 6 | Group 6           | Group 6    | Group 6        | PALEA                | Class 12     |
|        | 10       | Group 6 | Group 6           | Group 6    | Group 6        | PALEA                | Class 12     |

Gilles Hunault and Emmanuel Jaspard,  
LEAPdb: a database for the late embryogenesis abundant proteins  
*BMC Genomics* 2010, **11**:221

Emmanuel Jaspard, David Macherel and Gilles Hunault,  
Computational and Statistical Analyses of Amino Acid Usage and Physico-Chemical Properties of the Twelve Late Embryogenesis Abundant Protein Classes  
*PLoS ONE* 2012, **7**(5)

**Les protéines abondantes de l'embryogenèse tardive (protéines LEA) sont des protéines des plantes et de certaines bactéries et invertébrés qui protègent contre l'agrégation des protéines due à la dessiccation ou aux stress osmotiques associés aux basses températures. Les protéines LEA ont été initialement découvertes s'accumulant tard dans l'embryogenèse des graines de coton.**

**Les protéines LEA fonctionnent par des mécanismes distincts de ceux présentés par les chaperons moléculaires de choc thermique. Bien que les causes de l'induction de la protéine LEA n'aient pas encore été déterminées, des changements conformationnels dans les facteurs de transcription ou les protéines membranaires intégrales dus à la perte d'eau ont été suggérés. Les protéines LEA protègent particulièrement les membranes mitochondriales contre les dommages dus à la déshydratation.**

**A****OSMOTIC POTENTIAL (BARS)**

-9.5 to -30.5

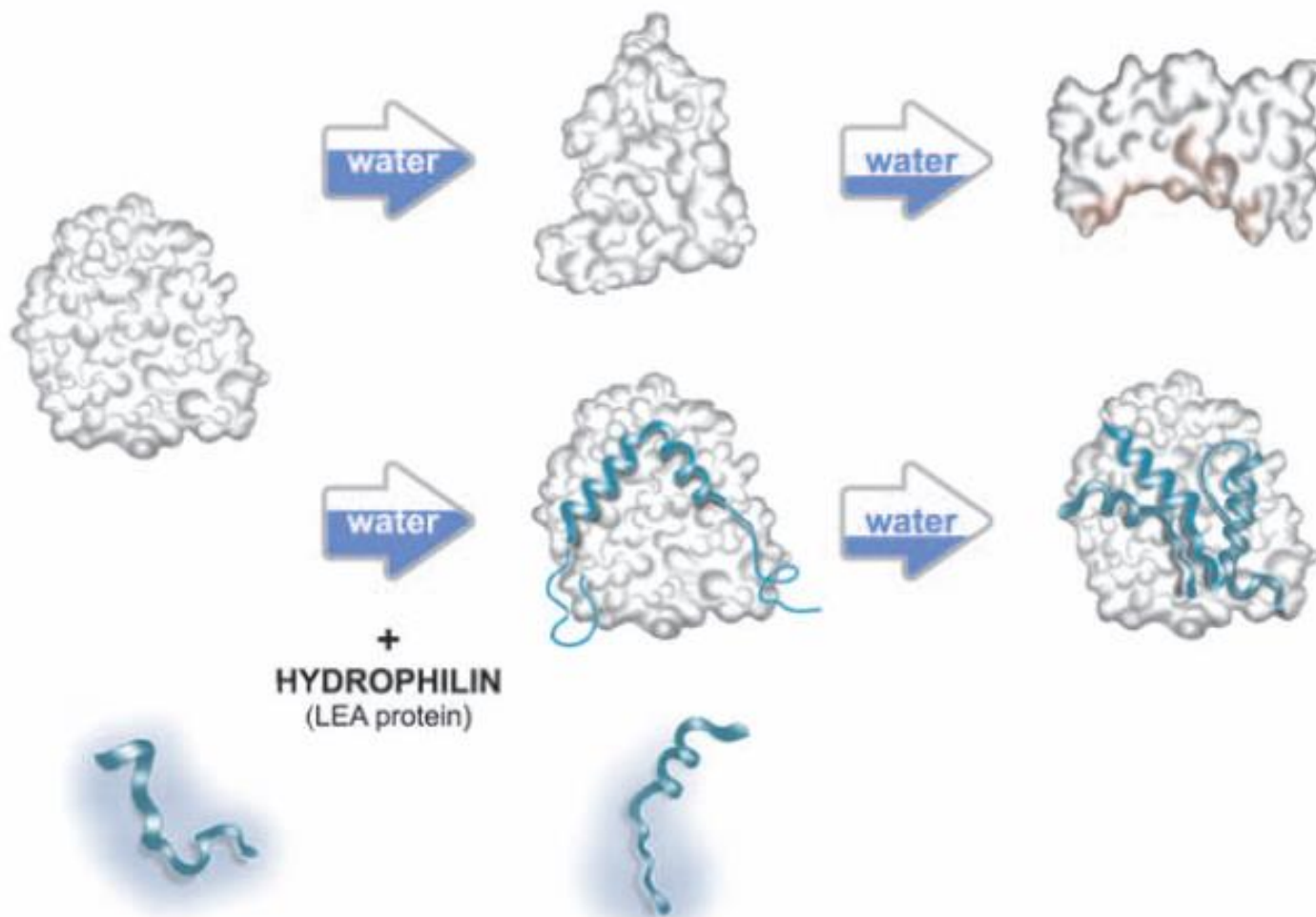
-58 or LOWER

**WATER LOSS (%)** 0

85 - 98

99 - 99.4

Enzyme activity

**ENZYME**

Enzyme activity



**Ce schéma illustre un modèle hypothétique pour la fonction des protéines LEA et d'autres hydrophilines. Dans cet exemple, sous modération déficit hydrique, une enzyme subit des changements conformationnels qui conduisent à une diminution de son activité et, dans des conditions de stress plus sévères, plus critiques les modifications structurelles entraînent l'exposition de résidus hydrophobes (ombrage rouge). La présence de protéines LEA (hydrophilines) (brin vert) empêche les modifications de la conformation de l'enzyme, à la suite desquelles l'enzyme conserve son activité, dans des conditions de limitation en eau. Cet effet peut être atteint à un rapport hydrophiline:enzyme de 1:1 sous un stress hydrique modéré; cependant, en cas de déshydratation sévère, l'action de plus d'un l'hydrophiline par molécule d'enzyme pourrait éviter d'autres changements conformationnels pouvant conduire à l'agrégation des protéines.**

## LA BASE DE DONNÉES OMIM (ONLINE MENDELIAN INHERITANCE IN MAN)

Donne de nombreuses informations sur la classification des maladies génétiques, des présentations cliniques et la cartographie génomique de la localisation de la maladie.

La base de données est mise à jour continuellement et offre probablement le meilleur lors de la recherche d'information sur les maladies héréditaires.



**OMIM**® Online Mendelian Inheritance in Man®  
An Online Catalog of Human Genes and Genetic Disorders  
Updated 11 March 2016

**Advanced Search :** OMIM, Clinical Synopses, Gene Map

**Need help? :** Example Searches, OMIM Search Help, OMIM Tutorial

**Mirror sites :** [us-east.omim.org](http://us-east.omim.org), [europe.omim.org](http://europe.omim.org)

OMIM is supported by a grant from NHGRI, licensing fees, and generous contributions from people like you.



**Advanced Search:** [OMIM](#), [Clinical Synopses](#), [OMIM Gene Map](#)**Search History:** [View](#), [Clear](#)

#219700

ICD+

CYSTIC FIBROSIS; CF

# OMIM : maladie

*Alternative titles; symbols*

MUCOVISCIDOSIS

[▶ Table of Contents - #219700](#)

External Links:

[▶ Clinical Resources](#)[▶ Animal Models](#)[▶ Cell Lines](#)[▶ Cellular Pathways](#)

## Phenotype Gene Relationships

| Location                | Phenotype                                   | Phenotype MIM number   | Gene/Locus | Gene/Locus MIM number  |
|-------------------------|---|------------------------|------------|------------------------|
| <a href="#">7q31.2</a>  | Cystic fibrosis                             | <a href="#">219700</a> | CFTR       | <a href="#">602421</a> |
| <a href="#">19q13.2</a> | {Cystic fibrosis lung disease, modifier of} | <a href="#">219700</a> | TGFB1      | <a href="#">190180</a> |

[Clinical Synopsis](#)

## TEXT

A number sign (#) is used with this entry because the disorder is caused by mutations in the cystic fibrosis conductance regulator gene (CFTR; [602421](#)), located on chromosome 7.

## Description

Formerly known as cystic fibrosis of the pancreas, this entity has increasingly been labeled simply 'cystic fibrosis.' Manifestations relate not only to the disruption of exocrine function of the pancreas but also to intestinal glands (meconium ileus), biliary tree (biliary cirrhosis), bronchial glands (chronic bronchopulmonary infection with emphysema), and sweat glands (high sweat electrolyte with depletion in a hot environment). Infertility occurs in males and females.

Lien vers les gènes ou les portions de chromosome responsables de la maladie

# Tree of life

- **Tree of life :**
  - <http://tolweb.org>
  - Base de données de taxonomie
    - Classification des êtres vivants
  - Avec des photos !

# Tree of life

## Fungi

Eumycota: mushrooms, sac fungi, yeast, molds, rusts, smuts, etc.

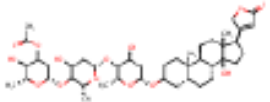
↗ Meredith Blackwell, Rytas Vilgalys, Timothy Y. James, and John W. Taylor





# DrugBank

- Base de données sur les médicaments
  - <http://www.drugbank.ca>
  - Information sur les cibles des médicaments
  - Attention : base américaine
    - => médicaments américains !

| Identification      |   |
|---------------------|---|
| Name                | <b>Acetyldigitoxin</b>  |
| Accession Number    | <b>DB00511</b> (APRD01334)  |
| Type                | small molecule  |
| Groups              | approved  |
| Description         | Cardioactive derivative of lanatoside A or of digitoxin used for fast digitalization in congestive heart failure.   |
| Structure           |  <p>Download: <a href="#">MOL</a>   <a href="#">SDF</a>   <a href="#">SMILES</a>   <a href="#">InChI</a><br/>                     Display: <a href="#">2D Structure</a>   <a href="#">3D Structure</a></p> |
| Synonyms            | Not Available   |
| Brand names         | <ul style="list-style-type: none"> <li>• Acigoxin</li> <li>• Acylanid</li> <li>• Crystodigin</li> </ul>   |
| Brand name mixtures | Not Available   |
| Categories          | <ul style="list-style-type: none"> <li>• Enzyme Inhibitors</li> <li>• Cardiotonic Agents</li> <li>• Anti-Arrhythmia Agents</li> </ul>   |

# DrugBank

## Targets


### 1. [Sodium/potassium-transporting ATPase alpha-1 chain](#)

Pharmacological action: **yes**

Actions: **inhibitor**

This is the catalytic component of the active enzyme, which catalyzes the hydrolysis of ATP coupled with the exchange of sodium and potassium ions across the plasma membrane. This action creates the electrochemical gradient of sodium and potassium ions, providing the energy for active transport of various nutrients

Organism class: **human**

UniProt ID: [P05023](#) 

Gene: [ATP1A1](#) 

Protein Sequence: [FASTA](#)

Gene Sequence: [FASTA](#)

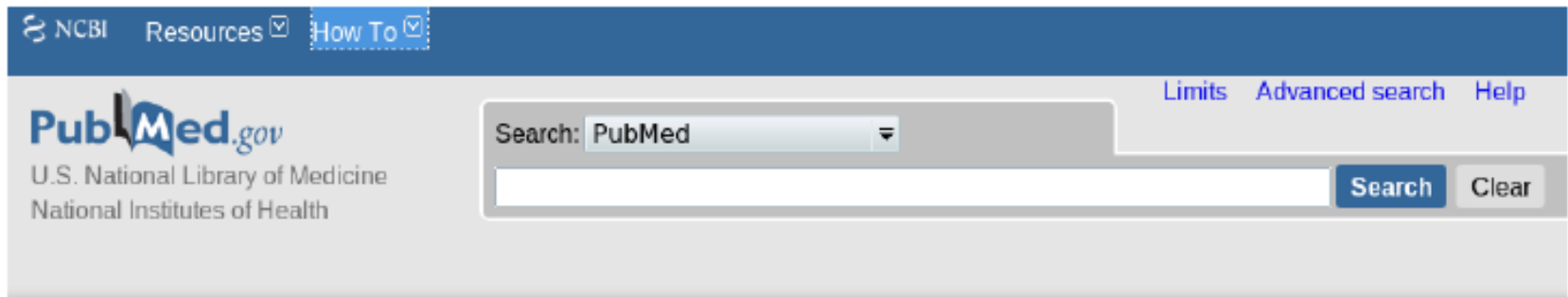
SNPs: [SNPJam Report](#) 

References:

1. Gonzalez-Garcia C, Cena V, Klein DC: Characterization of the alpha -like Na,K+-ATPase which mediates ouabain inhibition of adrenergic induction of N-acetyltransferase (EC 2.3.1.87) activity: studies with isolated pinealocytes. Mol Pharmacol. 1987 Dec;32(6):792-7. [Pubmed](#)

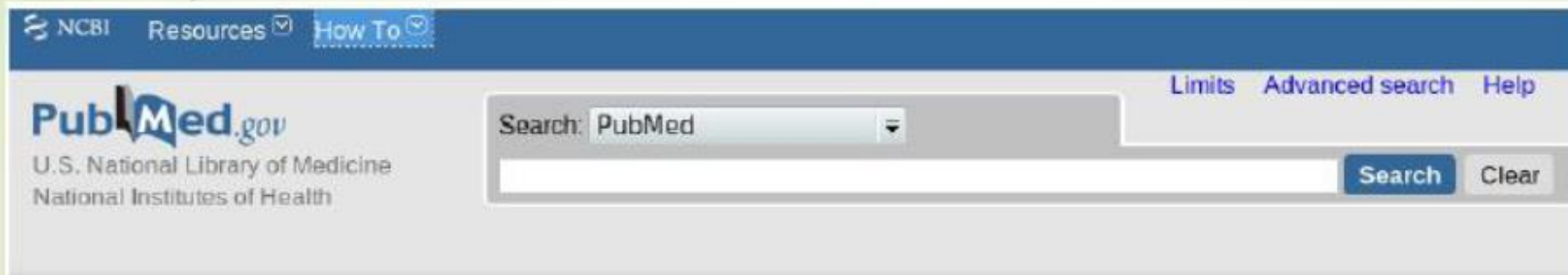
# Pubmed / Medline

- <http://www.ncbi.nlm.nih.gov/pubmed>
- Medline : base de données bibliographiques en médecine / biologie
- Pubmed : interface permettant de consulter la base
- Lien vers le texte des articles disponibles en ligne



The screenshot shows the top navigation bar of the PubMed website. It includes the NCBI logo, a 'Resources' dropdown menu, and a 'How To' dropdown menu. Below the navigation bar, the PubMed.gov logo is displayed on the left, along with the text 'U.S. National Library of Medicine' and 'National Institutes of Health'. On the right side of the page, there are links for 'Limits', 'Advanced search', and 'Help'. The main search area features a search bar with a dropdown menu set to 'PubMed', a large empty search input field, and 'Search' and 'Clear' buttons.

**PubMed** est une base de données bibliographiques, développé par le National Center for Biotechnology Information (NCBI), centrée sur la **documentation en sciences biologiques**.



The screenshot shows the top navigation bar of the PubMed website. On the left, there are links for 'NCBI', 'Resources', and 'How To'. The main logo 'PubMed.gov' is displayed, along with the text 'U.S. National Library of Medicine' and 'National Institutes of Health'. On the right side of the header, there are links for 'Limits', 'Advanced search', and 'Help'. Below the header is a search bar with a dropdown menu currently set to 'PubMed'. To the right of the search bar are 'Search' and 'Clear' buttons.

### Welcome to PubMed

PubMed comprises more than 19 million citations for biomedical articles from MEDLINE and life science journals. Citations may include links to full-text articles from PubMed Central or publisher web sites.

## Database Categories List

- ▶ Major Sequence Repositories
- ▶ Comparative Genomics
- ▶ Gene Expression
- ▶ Gene Identification and Structure
- ▶ Genetic and Physical Maps
- ▶ Genomic Databases
- ▶ Intermolecular Interactions
- ▶ Metabolic Pathways and Cellular Regulation
- ▶ Mutation Databases
- ▶ Pathology
- ▶ Protein Databases
- ▶ Protein Sequence Motifs
- ▶ Proteome Resources
- ▶ RNA Sequences
- ▶ Retrieval Systems and Database Structure
- ▶ Structure
- ▶ Transgenics
- ▶ Varied Biomedical Content

# Formats de fichiers en bioinformatique

# Définition d'un format:

- Les séquences sont stockées en général sous forme de fichiers texte, accessibles par des systèmes d'interrogations (Entrez pour GenBank, ACNUC pour EMBL, SRS pour UniProt, ...)
- Le Format: correspond à l'ensemble des règles de présentation auxquelles sont soumises la ou les séquences dans un fichier donné.
- Le Format permet:
  - Une mise en forme automatisée des données;
  - Le stockage et la gestion homogène de l'information;
  - Le traitement informatique ultérieur de l'information.



# Syntaxe d'une Entrée:

- Une Entrée: est une fiche signalétique d'une séquence donnée.

\* Contient 3 parties:

☐ Entête (header) : Description générale de la séquence;

❖ Les caractéristiques (Features): Description des objets biologiques présents sur la séquences

➤ La séquence

- ❑ Les formats de DDBJ et de GenBank sont très similaires
- ❑ Chaque ligne commence par un mot clé
  - Deux lettres pour EMBL
  - Maximum 12 lettres pour GenBank et DDBJ
- ❑ Fin d'une entrée par //

# Partie 01: Header = entête

## Thermus aquaticus DNA polymerase (PolI) gene

GenBank: J04639.1

[FASTA](#) [Graphics](#)

Go to:

LOCUS TTHTAQP1A 2626 bp DNA linear BCT 26-APR-1993  
DEFINITION Thermus aquaticus DNA polymerase (PolI) gene.  
ACCESSION J04639 M26480  
VERSION J04639.1  
KEYWORDS DNA polymerase.  
SOURCE Thermus aquaticus  
ORGANISM [Thermus aquaticus](#)  
Bacteria; Deinococcus-Thermus; Deinococci; Thermales; Thermaceae;  
Thermus.  
REFERENCE 1 (bases 1 to 2626)  
AUTHORS Lawyer,F.C., Stoffel,S., Saiki,R.K., Myambo,K., Drummond,R. and  
Gelfand,D.H.  
TITLE Isolation, characterization, and expression in Escherichia coli of  
the DNA polymerase gene from Thermus aquaticus  
JOURNAL J. Biol. Chem. 264 (11), 6427-6437 (1989)  
PUBMED [2649500](#)  
COMMENT Original source text: T.aquaticus (strain YT1; ATCC 25104) DNA.  
Draft entry and computer readable copy of sequence [1] kindly  
provided by D.H.Gelfand, 23-JUN-1989.

# Partie 02: features = caractéristiques

```
FEATURES             Location/Qualifiers
    source             1..2626
                        /organism="Thermus aquaticus"
                        /mol_type="genomic DNA"
                        /db_xref="taxon:271"
    CDS                121..2619
                        /EC_number="2.7.7.7"
                        /standard_name="Taq polymerase"
                        /codon_start=1
                        /transl_table=11
                        /product="DNA polymerase"
                        /protein_id="AAA27507.1"
                        /translation="MRGMLPLFEPKGRVLLVDGHHLAYRTFHALKGLTTSRGEPVQAV
YGFAKSLKALKEDGDAVIVVFDAAKAPSRHEAYGGYKAGRPTPEDFPRQLALIKEL
VDLLGLARLEVPGYEADDVLASLAKKAEKEGYEVRI LTADKDLYQLLSDRIHVLHPEG
YLITPAWLWEKYGLRPDQWADYRALTGDESDNLPGVKIGIEKTARKLLEEWGSLEALL
KNLDR LKPAIREKILAHMDDLKLSWDLAKVRTDLPLEVDFAKRREPDRERLRAFLERL
EFGSLLHEFGLLESPKALEEAPWPPPEGAFVGFVLSRKEPMWADLLALAAARGGRVHR
APEPYKALRDLKEARGLLAKDLSVLALREGLGLPPGDDPMLLAYLLDPSNTTPEGVAR
RYGGEWTEEAGERAALSERLFANLWGRLEGEERLLWLYREVERPLSAVLAHMEATGVR
LDVAYLRALSLEVAEEIARLEAEVFRLAGHPFNLNSRDQLERVL FDELGLPAIGKTEK
TGKRSTSAAVLEALREAHPIVEKILQYRELTKLKSTYIDPLPDLIHPRTGRLHTRFNQ
TATATGRLSSSDPNLQNI PVRTPLGQRIRRAFIAEEGWLLVALDYSQIELRVLAHLSG
DENLIRVFQEGRDIHTETASWMFGVPREAVDPLMRRAAKTINFGVLYGMSAHRLSQEL
AIPYEEAQAFIERYFQSFPKVRAWIEKTLEEGRRRGYVETLFGRRRYVPDLEARVKS
VREAAERMAFNMVQGTAA DLMKLAMVKLFPRLEEMGARMLLQVHDELVLEAPKERA
VARLAKEVMEGVYPLAVPLEVEVGIGEDWLSAKE"
```

# Partie 03: origin= séquence

ORIGIN

```
1  ggcatgaaag tcagggcaga gccatctatt gcttacattt gcttctgaca caactgtgtt
61  cactagcaac ctcaaacaga caccatgggtg cacctgactc ctgaggagaa gtctgcccgtt
121 actgccctgt ggggcaaggt gaacgtggat gaagttgggtg gtgaggccct gggcaggttg
181 gtatcaaggt tacaagacag gtttaaggag accaatagaa actgggcatg tggagacaga
241 gaagactctt gggtttctga taggcactga ctctctctgc ctattggtct attttcccac
301 ccttaggctg ctgggtggtct acccttggac ccagaggttc tttgagtcct ttgggggatct
361 gtccactcct gatgctgtta tgggcaacc taaggtgaag gctcatggca agaaagtgct
421 cggtgccttt agtgatggcc tggctcacct ggacaacctc aagggcacct ttgccacact
481 gagtgagctg cactgtgaca agctgcacgt ggatcctgag aacttcaggg tgagtctatg
541 ggacccttga tgttttcttt ccccttcttt tctatggtta agttcatgtc ataggaaggg
601 gagaagtaac agggtacagt ttagaatggg aaacagacga atgatt
```

//

Définition et numéro d'accèsion

LOCUS AAB59426 490 aa  
DEFINITION cytochrome.  
ACCESSION AAB59426  
VERSION AAB59426.1 GI:181344  
DBSOURCE locus HMCYPC219 accession [M61854.1](#)

linear

PRI 31-DEC-1994

Date de soumission

Type de molécule: linéaire ou circulaire

Organisme et taxonomie

SOURCE Homo sapiens (human)  
ORGANISM [Homo sapiens](#)  
Eukaryota; Metazoa; Chordata; Craniota; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo.

Référence bibliographique

REFERENCE 1 (residues 1 to 490)  
AUTHORS Romkes, M., Yaletto, K.B., Blaisdell, J.A., Raucy, J.L. and Goldstein, J.A.  
TITLE Cloning and expression of complementary DNAs for multiple members of the human cytochrome P45011C subfamily  
JOURNAL Biochemistry 30 (13), 3247-3255 (1991)  
PubMed [3009263](#)

Caractéristiques de la séquence: Domaines, Motifs, RBS, CDS...

COMMENT Method: conceptual translation.  
FEATURES  
Location/Qualifiers  
source 1..490  
/organism="Homo sapiens"  
/db\_xref="taxon:[9606](#)"  
[Protein](#) 1..490  
/product="cytochrome"  
[Region](#) 30..487  
/region\_name="p450"  
/note="Cytochrome P450; pfam00067"  
/db\_xref="CDD:[40168](#)"  
[Region](#) 30..480  
/region\_name="CypX"  
/note="Cytochrome P450 [Secondary metabolites biosynthesis, transport, and catabolism]; COG2124"  
/db\_xref="CDD:[32307](#)"  
[CDS](#) 1..490  
/gene="CYP2C19"  
/coded\_by="M61854.1:6..1470"

GenBank/DDBJ

Séquence

ORIGIN  
1 mdprvvlvic lscilllisv rqsagrgkip pgptlpvvg nllqidikdv sksitniski  
61 ygpvftlyfg lernvvihg y evvkealidl geefagrgfh plaeranrgi givisngkrv  
121 keirrfslmt lrfngmgkrs iedrvqeear clveelrktk aspedptfil qcapcnvics  
181 iifqkrfdyk dqqlnlsek lnenirivst pwiqcanfp tiidyfpgth nkliknlafm  
241 esdilekvke hqemadinp rdfidoflik wekekqgqs eftienlvit eadilgagte  
301 ttsttiryal lllkhpevt akvqecierv igrnrpcmq drqlmptyda vvhcvgryd  
361 liptslphav cedvkfrnyl ipkgttlits lcevlhdnke fpapemfdpr hfldeggmfk  
421 ksnymfpea gkricvgegl armelfifit filqnfalks lidpkldtct pvvngfasvp  
481 pfvqlcfiyp

Descripteur de fin de fichier

///

# EMBL Flat File

## Header

- Title
- Taxonomy
- Citation

## Features (AA seq)

## DNA Sequence

```
ID AF115338 standard: DNA: PRO: 591 BP.
AC AF115338;
SV AF115338.1
DT 03-JUN-1999 (Rel. 59, Created)
DT 23-AUG-1999 (Rel. 60, Last updated, Version 2)
DE Pseudomonas fluorescens ECF sigma factor SigX (sigX) gene, complete cds.
KW .
OS Pseudomonas fluorescens
OC Bacteria: Proteobacteria: gamma subdivision: Pseudomonadaceae: Pseudomonas.
RN [1]
RP 1-591
RX MEDLINE; 99369842.
RA Brinkman F.S., Schoofs G., Hancock R.E., De Mot R.;
RT "Influence of a putative ECF sigma factor on expression of the major outer
RT membrane protein, OprF, in Pseudomonas aeruginosa and Pseudomonas
RT fluorescens";
RL J. Bacteriol. 181(16):4746-4754(1999).
RN [2]
RP 1-591
RA De Mot R.;
RT ;
RL Submitted (04-DEC-1998) to the EMBL/GenBank/DDBJ databases.
RL F.A. Janssens Laboratory of Genetics, Applied Plant Sciences, K.
RL Mercierlaan 92, Heverlee B-3001, Belgium
DR SPTREMBL: Q9X4L7: Q9X4L7.
FH Key Location/Qualifiers
FH
FT source 1..591
FT /db_xref="taxon:294"
FT /organism="Pseudomonas fluorescens"
FT /strain="M114"
FT CDS 1..591
FT /codon_start=1
FT /db_xref="SPTREMBL:Q9X4L7"
FT /transl_table=11
FT /gene="sigX"
FT /product="ECF sigma factor SigX"
FT /protein_id="AAD34329.1"
FT /translation="MNKAQTLSTRYDPRELSDEELVARSHTELFHVTRAYEELMRRYQR
FT ILFNVCARYLGNDRDADDVQVEVMLKVLYGLKNLEGKSKFKTWWLYSITYNECITQYRKE
FT RRRRLMDALSLDPLEEASEEKALQPEEKGGLDRLWLVVNPIDRGIIVLRFVAE LEFQE
FT IADIMRMGLSATKMRYKRALDKLREKTFAGE TET"
SQ Sequence 591 BP; 157 A; 133 C; 170 G; 131 T; 0 other:
atgaataaag cccaaacgct atccacgcgc tacgaccccc gccagctctc tgatgaggag 60
ttggtcgcgc gotcgcatac cgagcttttt cacgtaacgc gccctatga agaactgatg 120
cggcgttacc agcgaacatt atttaacggt tgtgcgagat atcttgggaa cgatcgcgac 180
gcagacgatg totgtcagga agtcatggtt aaggtgctgt atggcctgaa gaaacctcag 240
gggaaatcga agttcaaaac guggctctac agcatcacgt acaacgaatg tattacgcag 300
tatcgggaag aacggcgaaa gogtgccttg atggacgcat tgagtcttga cccctcag 360
```

# Les flatfiles en détails (1/8)

## 1. Les différentes informations de l'entête (*header*) :

- Première ligne : Locus/ID

→ Embl

|    |         |   |
|----|---------|---|
| ID | U49845; | SV 1; linear; genomic DNA; STD; FUN; 5028 BP. |
|----|---------|---|

→ DDBJ/GenBank

|       |          |         |     |        |                 |
|-------|----------|---------|-----|--------|-----------------|
| LOCUS | SCU49845 | 5028 bp | DNA | linear | PLN 21-JUN-1999 |
|-------|----------|---------|-----|--------|-----------------|

- Le champ date (chez EMBL uniquement)

|    |                                |
|----|--------------------------------|
| DT | 07-MAY-1996 (Rel. 47, Created) |
|----|--------------------------------|

|    |  |
|----|--|
| DT | 17-APR-2005 (Rel. 83, Last updated, Version 4) |
|----|--|

- Lignes de définition : synthèse du contenu biologique

|    |   |
|----|---|
| DE | Saccharomyces cerevisiae TCP1-beta gene, partial cds; and Axl2p (AXL2)<br>and |
|----|---|

|    |                                   |
|----|-----------------------------------|
| DE | Rev7p (REV7) genes, complete cds. |
|----|-----------------------------------|

|            |  |
|------------|--|
| DEFINITION | Saccharomyces cerevisiae TCP1-beta gene, partial cds; and Axl2p (AXL2) and Rev7p (REV7) genes, complete cds. |
|------------|--|



# Division en organisme

Utilisée dans les 3 principales banques de séquences nucléiques

| Sequence division                          | Database            |
|--|---------------------|
| <i>Organismal</i>                          |                     |
| BCT Bacterial                              | DCBJ, GenBank       |
| PRO Prokaryotic                            | EMBL                |
| FUN Fungal                                 | EMBL                |
| HUM Human                                  | DCBJ, EMBL          |
| PRI Primate                                | DCBJ, EMBL, GenBank |
| ROD Rodent                                 | DCBJ, EMBL, GenBank |
| MAM Other mammalian                        | DCBJ, EMBL, GenBank |
| VRT Other vertebrate                       | DCBJ, EMBL, GenBank |
| INV Invertebrate                           | DCBJ, EMBL, GenBank |
| PLN Plant                                  | DCBJ, EMBL, GenBank |
| ORG Organelle                              | EMBL                |
| VRL Viral                                  | DCBJ, EMBL, GenBank |
| PHG Phage                                  | DCBJ, EMBL, GenBank |
| RNA Structural RNA                         | DCBJ, EMBL, GenBank |
| SYN Synthetic and chimeric                 | DCBJ, EMBL, GenBank |
| UNA Unannotated                            | DCBJ, EMBL, GenBank |
|  |                     |
| <i>Functional</i>                          |                     |
| EST Expressed sequence tag                 | DCBJ, EMBL, GenBank |
| STS Sequence tagged site                   | DCBJ, EMBL, GenBank |
| GSS Genome survey                          | DCBJ, EMBL, GenBank |
| HTG High-throughput genomic                | DCBJ, EMBL, GenBank |
| PAT Patent                                 | DCBJ, EMBL, GenBank |
| CON Virtual contigs of segmented sequences | DCBJ, EMBL, GenBank |

# Les flatfiles en détails (2/8)

- Le numéro d'accèsion : un id unique (commun aux 3 banques)

|    |         |
|----|---------|
| AC | U49845; |
|----|---------|

|           |        |
|-----------|--------|
| ACCESSION | U49845 |
|-----------|--------|

- La version (équivalent à SV dans la 1re ligne d'EMBL)

|         |          |             |
|---------|----------|-------------|
| VERSION | U49845.1 | GI :1293613 |
|---------|----------|-------------|

- Lignes avec des mots-clés (KEYWORDS ou KW)

- Lignes de taxonomie

|    |  |
|----|--|
| OS | Saccharomyces cerevisiae (baker's yeast) |
|----|--|

|    |  |
|----|--|
| OC | Eukaryota ; Fungi ; Dikarya ; Ascomycota ; Saccharomycotina ;<br>Saccharomycetes ; |
|----|--|

|    |   |
|----|---|
| OC | Saccharomycetales ; Saccharomycetaceae ; Saccharomyces. |
|----|---|

|        |  |
|--------|--|
| SOURCE | Saccharomyces cerevisiae (baker's yeast) |
|--------|--|

|          |                          |
|----------|--------------------------|
| ORGANISM | Saccharomyces cerevisiae |
|----------|--------------------------|

|  |   |
|--|---|
|  | Eukaryota ; Fungi ; Dikarya ; Ascomycota ; Saccharomycotina ;<br>Saccharomycetes ; Saccharomycetales ; Saccharomycetaceae ;<br>Saccharomyces. |
|--|---|

# Les flatfiles en détails (3/8)

- Les références : publication ou origine de la soumission

RN [1]  
RP 1-5028  
RX PUBMED; 7871890.  
RA Torpey L.E., Gibbs P.E., Nelson J., Lawrence C.W. ;  
RT "Cloning and sequence of REV7, a gene whose function is required for DNA  
RT damage-induced mutagenesis in *Saccharomyces cerevisiae*" ;  
RL Yeast 10(11) :1503-1509(1994).

XX  
RN [2]  
RP 1-5028  
RX PUBMED; 8846915.  
RA Roemer T., Madden K., Chang J., Snyder M. ;  
RT "Selection of axial growth sites in yeast requires Axl2p, a novel plasma  
RT membrane glycoprotein" ;  
RL Genes Dev. 10(7) :777-793(1996).

XX  
RN [3]  
RP 1-5028  
RA Roemer T. ;  
RT ;  
RL Submitted (22-FEB-1996) to the EMBL/GenBank/DDBJ databases.  
RL Terry Roemer, Biology, Yale University, New Haven, CT, USA

# Les flatfiles en détails (4/8)

REFERENCE 1 (bases 1 to 5028)  
AUTHORS Torpey,L.E., Gibbs,P.E., Nelson,J. and Lawrence,C.W.  
TITLE Cloning and sequence of REV7, a gene whose function is required for DNA damage-induced mutagenesis in *Saccharomyces cerevisiae*  
JOURNAL Yeast 10 (11), 1503-1509 (1994)  
PUBMED 7871890

REFERENCE 2 (bases 1 to 5028)  
AUTHORS Roemer,T., Madden,K., Chang,J. and Snyder,M.  
TITLE Selection of axial growth sites in yeast requires Axl2p, a novel plasma membrane glycoprotein  
JOURNAL Genes Dev. 10 (7), 777-793 (1996)  
PUBMED 8846915

REFERENCE 3 (bases 1 to 5028)  
AUTHORS Roemer,T.  
TITLE Direct Submission  
JOURNAL Submitted (22-FEB-1996) Terry Roemer, Biology, Yale University, New Haven, CT, USA

# Les flatfiles en détails (5/8)

---

## 2. Les caractéristiques (*features*)

- Features table : annotation des séquences, gestion du transfert d'information, localisation des éléments biologiques
- Mise à disposition d'un vocabulaire étendu contrôlé (**Gene Ontology**) pour décrire les caractéristiques biologiques des séquences (annotations)
- Annotation = information additionnelle venant enrichir un document d'intérêt
  - Annotation structurale : caractériser les gènes au travers, notamment, de leur structure intron-exon, signaux de régulation, sites d'épissage, ...
  - Annotation fonctionnelle : anticiper sur les motifs fonctionnels qui seront retrouvés au niveau protéiques, caractérisation des fonctions biologiques des produits d'expression des gènes

# Les flatfiles en détails (6/8)

```
FT source 1..5028
FT /organism="Saccharomyces cerevisiae"
FT /chromosome="IX"
FT /map="9"
FT /mol_type="genomic DNA"
FT /db_xref="taxon:4932"
FT CDS <1..206
FT /codon_start=3
FT /product="TCP1-beta"
FT /db_xref="GOA:P39076"
FT /db_xref="UniProtKB/Swiss-Prot:P39076"
FT /protein_id="AAA98665.1"
FT /translation="SSIYNGISTSGLDLNNGTIADMRQLGIVESYKLGKRAVVSSASEAA
FT EVLLRVDNIIRARPRTANRQHM"
FT [...]
FT CDS complement(3300..4037)
FT /codon_start=1
FT /gene="REV7"
FT /product="Rev7p"
FT /db_xref="GOA:P38927"
FT /db_xref="InterPro:IPR003511"
FT /db_xref="SGD:S000001401"
FT /db_xref="UniProtKB/Swiss-Prot:P38927"
FT /protein_id="AAA98667.1"
FT /translation="MNRWVEKWLRVYLKCYINLILFYRNVYPPQSFDTTYQSFNLPQF
FT VPINRHPALIDYIEELILDVLSKLT HVYRFSICIINKKNDLCIEKYVLD FSELQHVDKD
FT DQIITETEVFDEFRSSLNSLIMHLEKLPKVND D TITFEAVINAIELELGHKLDRNRRVD
FT SLEEKAEIERDSNWVKCQEDENLPDNNGFQPPKIKL TSLVGSDVGPLIIHQFSEKLISG
FT DDKILNGVYSQYEEGESIFGSLF"
```

# Les flatfiles en détails (8/8)

## 3. La séquence de nucléotides

### → EMBL

```
50 Sequence 5028 BP; 1510 A; 1074 C; 835 G; 1609 T; 0 other;
gatcctccat atacaacggg atctccacct caggtttaga tctcaacaac ggaaccattg      60
ccgacatgag acagttaggg atcgtcgaga gttacaagct aaaacgagca gtagtcagct      120
ctgcatctga agccgctgaa gttctactaa ggggggataa catcatccgt gcaagaccaa      180
gaaccgccaa tagacaacat atgtaacata tttaggatat acctcgaaaa taataaacgg      240
ccacactgtc attattataa ttagaaacag aacgcaaaaa ttatccacta tataattcaa      300
agacgcgaaa aaaaaagaac aacgcgtcat agaactttg gcaattcgcg tcacaaataa      360
atittggcaa cttatgttcc ctcttcgagc agtactcgag cctgtctca agaatgtaat      420
aataccatc  gtaggtatgg ttaaagatag catctccaca acctcaaaagc tctttgccga      480
gagtcgcct  cttttgtcga gtaattttca cttttcatat gagaacttat ttctttatc      540
```

### → GenBank

ORIGIN

```
1 gatcctccat atacaacggg atctccacct caggtttaga tctcaacaac ggaaccattg
61 ccgacatgag acagttaggg atcgtcgaga gttacaagct aaaacgagca gtagtcagct
121 ctgcatctga agccgctgaa gttctactaa ggggggataa catcatccgt gcaagaccaa
181 gaaccgccaa tagacaacat atgtaacata tttaggatat acctcgaaaa taataaacgg
241 ccacactgtc attattataa ttagaaacag aacgcaaaaa ttatccacta tataattcaa
301 agacgcgaaa aaaaaagaac aacgcgtcat agaactttg gcaattcgcg tcacaaataa
361 atittggcaa cttatgttcc ctcttcgagc agtactcgag cctgtctca agaatgtaat
421 aataccatc  gtaggtatgg ttaaagatag catctccaca acctcaaaagc tctttgccga
481 gagtcgcct  cttttgtcga gtaattttca cttttcatat gagaacttat ttctttatc
541 tttactctca catctgttag tgattgacac tgcaacagcc acctccacta gaagaacaga
```

ID TCPB\_YEAST Reviewed; 527 AA.  
AC P39076; D6VVE5;  
DT 01-FEB-1995, integrated into UniProtKB/Swiss-Prot.  
DT 01-FEB-1995, sequence version 1.  
DT 12-SEP-2018, entry version 165.  
DE RecName: Full=T-complex protein 1 subunit beta;  
DE Short=TCP-1-beta;  
DE AltName: Full=CCT-beta;  
GN Name=CCT2; Synonyms=BIN3, TCP2; OrderedLocusNames=YIL142W;  
OS *Saccharomyces cerevisiae* (strain ATCC 204508 / S288c) (Baker's yeast).  
OC Eukaryota; Fungi; Dikarya; Ascomycota; Saccharomycotina;  
OC Saccharomycetes; Saccharomycetales; Saccharomycetaceae; *Saccharomyces*.  
OX NCBI\_TaxID=559292;  
RN [1]  
RP NUCLEOTIDE SEQUENCE [GENOMIC DNA].  
RC STRAIN=ATCC 204511 / S288c / AB972;  
RX PubMed=7908441; DOI=10.1073/pnas.91.7.2743;  
RA Miklos D., Caplan S., Mertens D., Hynes G., Pitluk Z., Kashi Y.,  
RA Harrison-Lavoie K., Stevenson S., Brown C., Barrell B.G.,  
RA Horwich A.L., Willison K.;  
RT "Primary structure and function of a second essential member of the  
RT heterooligomeric TCPI chaperonin complex of yeast, TCPI beta."  
RL Proc. Natl. Acad. Sci. U.S.A. 91:2743-2747(1994).



# Formats de fichiers

Il existe plusieurs formats de stockage du texte :

☐ Spécifiques aux séquences :

➤ Séquence brut (plain raw sequence)

➤ Format FASTA

➤ Format EMBL

➤ Format GenBank

➤ Format PDB

## Formats de fichiers : Texte brut

- ACAAGATGCCATTGTCCCCGGCCTCCTGCTGCTGCTGCTCTCCGGGGC  
CACGGCCACCGCTGCCCTGCCCTGGAGGGTGGCCCCACCGGCCGAGAC  
AGCGAGCATATGCAGGAAGCGGCAGGAATAAGGAAAAGCAGCCTCCTGA  
CTTTCCTCGCTTGGTGGTTTGAGTGGACCTCCCAGGCCAGTGCCGGGGCC  
CCTCATAGGAGAGGAAGCTCGGGAGGTGGCCAGGCGGCAGGAAGGCGCA  
CCCCCCCAGCAATCCGCGCGCCGGGACAGAATGCCCTGCAGGAACTTCT  
TCTGGAAGACCTTCTCCTCCTGCAAATAAAACCTCACCCATGAATGCTC  
ACGCAAGTTTAATTACAGACCTGAA

- ❑ Ne contient que des lettres désignant la séquence (acides aminés ou ADN)
- ❑ Une seule séquence est représentée

# Formats de fichiers : FASTA

- Format commun de manipulation des données :
  - Objectif : **manipuler facilement** des séquences dans les bases de données, à l'aide d'un **format universel**, compatibles avec les traitements de texte (sous forme de fichier texte), ou par copier – coller.
  - Le format FASTA : une 1<sup>re</sup> ligne de définition introduite par un >, contenant un identifiant sans espace et une description suivie de la séquence elle-même, sur plusieurs lignes (de taille 60 ou 80 classiquement)

```
>embl|U49845|U49845 Saccharomyces cerevisiae TCPI-beta  
gene, partial cds; and Axl2p (AXL2) and Rev7p (REV7) genes,  
complete
```

```
gatcctccatatacaacggtatctccacctcaggtttagatctcaacaacggaaccattg  
ccgacatgagacagtttaggtatcgctcgagagttacaagctaaaacgagcagtagtcagct  
ctgcatctgaagccgctgaagttctactaagggtggataacatcatccgtgcaagaccaa  
gaaccgccaatagacaacatatgtaacatatttaggatatacctcgaaaataataaacg  
ccacactgtcattattataattagaaacagaacgcaaaaattatccactatataattcaa  
agacgcgaaaaaaaaagaacaacgcgctcatagaacttttggcaattcgcgtcacaaataa  
atthttggcaacttatgtttcctcttcgagcagtagctcgagccctgtctcaagaatgtaat  
aataccatcgtaggtatggttaaagatagcatctccacaacctcaaagctccttgccga
```

# Formats de fichiers : EMBL

ID AB000263 standard; RNA; PRI; 368 BP.

XX

o AC AB000263;

XX

DE Homo sapiens mRNA for prepro cortistatin like peptide, complete cds.

XX

SQ Sequence 368 BP;

```
acaagatgcc attgtccccc ggccctctgc tgctgctgct ctccggggcc acggccaccg 60
ctgccctgcc cctggagggt ggccccaccg gccgagacag cgagcatatg caggaagcgg 120
caggaataag gaaaagcagc ctctgactt tcctcgcttg gtggtttgag tggacctccc 180
aggccagtgc cgggcccctc ataggagagg aagctcggga ggtggccagg cggcaggaag 240
gcgcaccccc ccagcaatcc gcgcgccggg acagaatgcc ctgcaggaac ttcttctgga 300
agaccttctc ctctgcaaa taaaacctca cccatgaatg ctcacgcaag ttttaattaca 360
gacctgaa 368
//
```

- Peut contenir plusieurs séquences
- Chaque séquence commence par « ID » suivi par des descriptions.
- La séquence suit une ligne marquée de « SQ » et est terminée par « // »

# Formats de fichiers : GenBank

LOCUS AB000263 368 bp mRNA linear PRI 05-FEB-1999

DEFINITION Homo sapiens mRNA for prepro cortistatin like peptide, complete cds.

ACCESSION AB000263

ORIGIN

```
1   acaagatgcc   attgtccccc   ggctctctgc   tgctgctgct   ctccggggcc   acggccaccg
61  ctgccctgcc   cctggagggt   ggccccaccg   gccgagacag   cgagcatatg   caggaagcgg
121 caggaataag   gaaaagcagc   ctctgactt   tcctcgcttg   gtggtttgag   tggacctccc
181 aggccagtgc   cgggcccctc   ataggagagg   aagctcggga   ggtggccagg   cggcaggaag
241 gcgcaccccc   ccagcaatcc   gcgcgccggg   acagaatgcc   ctgcaggaac   ttcttctgga
301 agaccttctc   ctctgcaaa   taaaacctca   cccatgaatg   ctcacgcaag   ttaattaca
361 gacctgaa
```

//

- ❑ Peut contenir plusieurs séquences (séparées par « // »)
- ❑ Contient d'abord des mots clés (LOCUS, ACCESSION ...) suivis de leurs valeurs, puis le mot clé ORIGIN suivi de la séquence (60 caractères par ligne)

# Formats de fichiers : IG

```
; comment
; comment
AB000263
ACAAGATGCCATTGTCCCCGGCCTCCTGCTGCTGCTGCTCTCCGGGGCCACGGCCACC
GCTGCCCTGCC
CCTGGAGGGTGGCCCCACCGGCCGAGACAGCGAGCATATGCAGGAAGCGGCAGGAATAA
GGAAAAGCAGC
CTCCTGACTTTCCTCGCTTGGTGGTTTGAGTGGACCTCCCAGGCCAGTGCCGGGGCCCCT
CATAGGAGAGG
AAGCTCGGGAGGTGGCCAGGCGGCAGGAAGGCGCACCCCCCAGCAATCCGCGCGCC
GGGACAGAATGCC
CTGCAGGAACTTCTTCTGGAAGACCTTCTCCTCCTGCAAATAAAACCTCACCCATGAATG
CTCACGCAAGTTTAATTACAGACCTGAAI
```

- Peut contenir plusieurs séquences
- Commence par des lignes de commentaires commençant par « ; », puis une ligne indiquant le nom de la séquence (sans espaces), puis la séquence.
- La séquence se termine par « I » si elle est linéaire, « 2 » si elle est circulaire.

# Formats de fichiers : PDB (1/2)

```
HEADER          OXIDOREDUCTASE                27-OCT-03                IUR5
TITLE          STABILIZATION OF A TETRAMERIC MALATE DEHYDROGENASE BY
TITLE          2 INTRODUCTION OF A DISULFIDE BRIDGE AT THE DIMER/DIMER
TITLE          3 INTERFACE
COMPND         MOL_ID: 1;
COMPND         2 MOLECULE: MALATE DEHYDROGENASE; COMPND 3 CHAIN: A, C;
COMPND         4 EC: 1.1.1.37; COMPND 5 ENGINEERED: YES;
COMPND         6 MUTATION: YES
SOURCE         MOL_ID: 1;
SOURCE         2 ORGANISM_SCIENTIFIC: CHLOROFLEXUS AURANTIACUS;
SOURCE         3 EXPRESSION_SYSTEM: ESCHERICHIA COLI;
SOURCE         4 EXPRESSION_SYSTEM_STRAIN: DH5A
KEYWDS         OXIDOREDUCTASE, TRICARBOXYLIC ACID CYCLE,
KEYWDS         2 MALATE DEHYDROGENASE
EXPDTA        X-RAY DIFFRACTION
AUTHOR        A.BJORK,B.DALHUS,D.MANTZILAS,V.G.H.EIJSINK,R.SIREVAG
```

# Formats de fichiers : PDB (1/2)

```
SEQRES  1 A  309 MET  ARG LYS  LYS  ILE  SER  ILE  ILE  GLY  ALA  GLY  PHE  VAL
SEQRES  2 A  309 GLY  SER  THR  THR  ALA  HIS  TRP  LEU  ALA  ALA  LYS  GLU  LEU
SEQRES  3 A  309 GLY  ASP  ILE  VAL  LEU  LEU  ASP  ILE  VAL  GLU  GLY  VAL  PRO
SEQRES  4 A  309 GLN  GLY  LYS  ALA  LEU  ASP  LEU  TYR  GLU  ALA  SER  PRO  ILE
SEQRES  5 A  309 GLU  GLY  PHE  ASP  VAL  ARG  VAL  THR  GLY  THR  ASN  ASN  TYR
SEQRES  6 A  309 ALA  ASP  THR  ALA  ASN  SER  ASP  VAL  ILE  VAL  VAL  THR  SER
...
SEQRES  1 C  309 MET  ARG LYS  LYS  ILE  SER  ILE  ILE  GLY  ALA  GLY  PHE  VAL
SEQRES  2 C  309 GLY  SER  THR  THR  ALA  HIS  TRP  LEU  ALA  ALA  LYS  GLU  LEU
SEQRES  3 C  309 GLY  ASP  ILE  VAL  LEU  LEU  ASP  ILE  VAL  GLU  GLY  VAL  PRO
SEQRES  4 C  309 GLN  GLY  LYS  ALA  LEU  ASP  LEU  TYR  GLU  ALA  SER  PRO  ILE
SEQRES  5 C  309 GLU  GLY  PHE  ASP  VAL  ARG  VAL  THR  GLY  THR  ASN  ASN  TYR
SEQRES  6 C  309 ALA  ASP  THR  ALA  ASN  SER  ASP  VAL  ILE  VAL  VAL  THR  SER
...
```

- Contient une seule séquence
- Commence par des lignes de descriptions. La séquence suit les lignes débutant par « SEQRES »
- Acides aminés codés par 3 lettres, acides nucléiques par DA, DC, DG, DT, DI