

HMSN204

Info. Biologique et Outils bioinformatiques

Banques de données biologiques
(3h de Cours/TD + 4.5h de TP/TD)

Anne-Muriel Chifolleau

<http://www.lirmm.fr/~arigon/enseignement/HMSN204/>



LIRMM - UM



-
- Introduction
 - Banques de données de séquences nucléiques
 - Banques de données de séquences protéiques
 - Banques de données spécialisées
 - Interrogation des banques de données
 - Formats de fichiers en bioinformatique
 - Références

- Comprendre ce qui constitue le “livre de la vie” , comprendre comment une collection de simple nucléotides est à la base de la vie, ...
 - collecter d’immense quantité de données de séquences et les stocker de manière à pouvoir les analyser et les retrouver facilement.
- Volume de données généré trop important pour être géré par les techniques de publication traditionnelle.

Création et maintenance des sources de données biologiques pour l’archivage, le stockage, la diffusion et l’exploitation des données biologiques

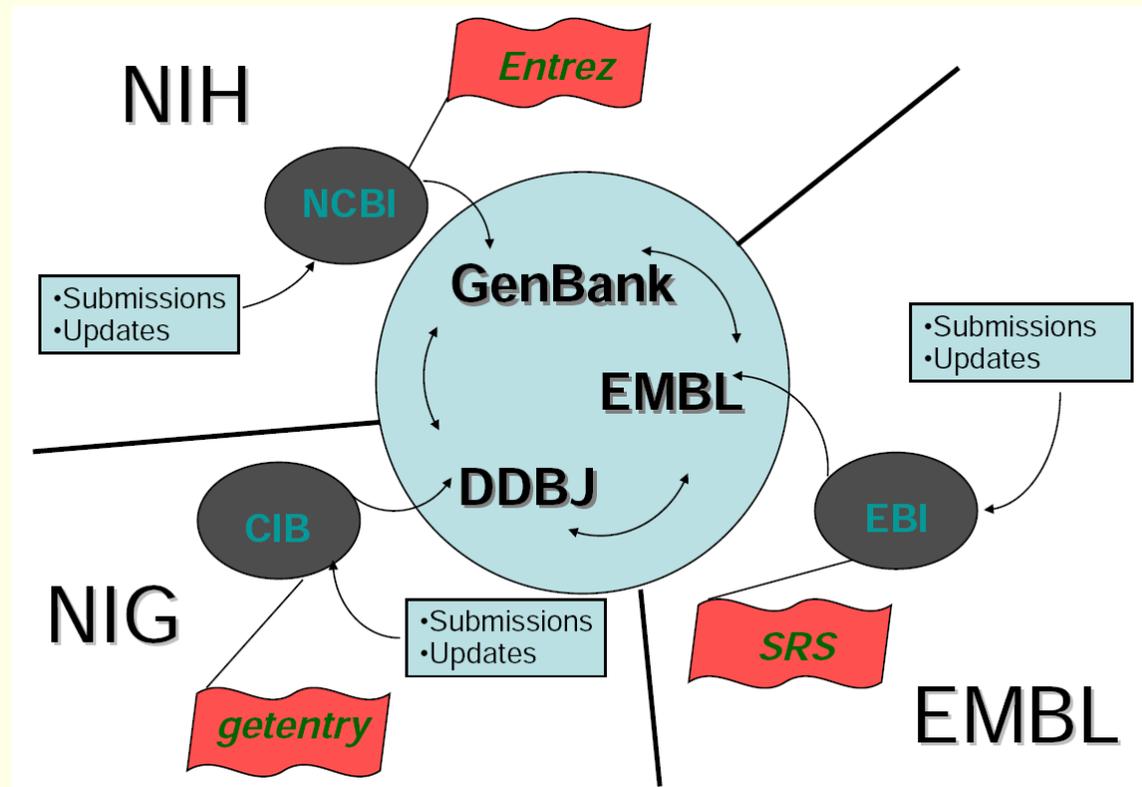
- Source de données en biologie = collection organisée d'informations mémorisées sur un support permanent et pouvant être restituées au biologiste de manière opportune
 - ⇒ Construire des **banques de données** ou des **bases de données** pour collecter les informations, stocker et organiser les données, distribuer l'information, faciliter l'exploitation des données
 - **Banque de données** : collection de fichiers structurés
 - **Base de données** : collection de données structurées

Différence = élément de base considéré Fichier de données ou Donnée
--

- Séquence = élément central autour duquel les banques et bases de données se sont constituées

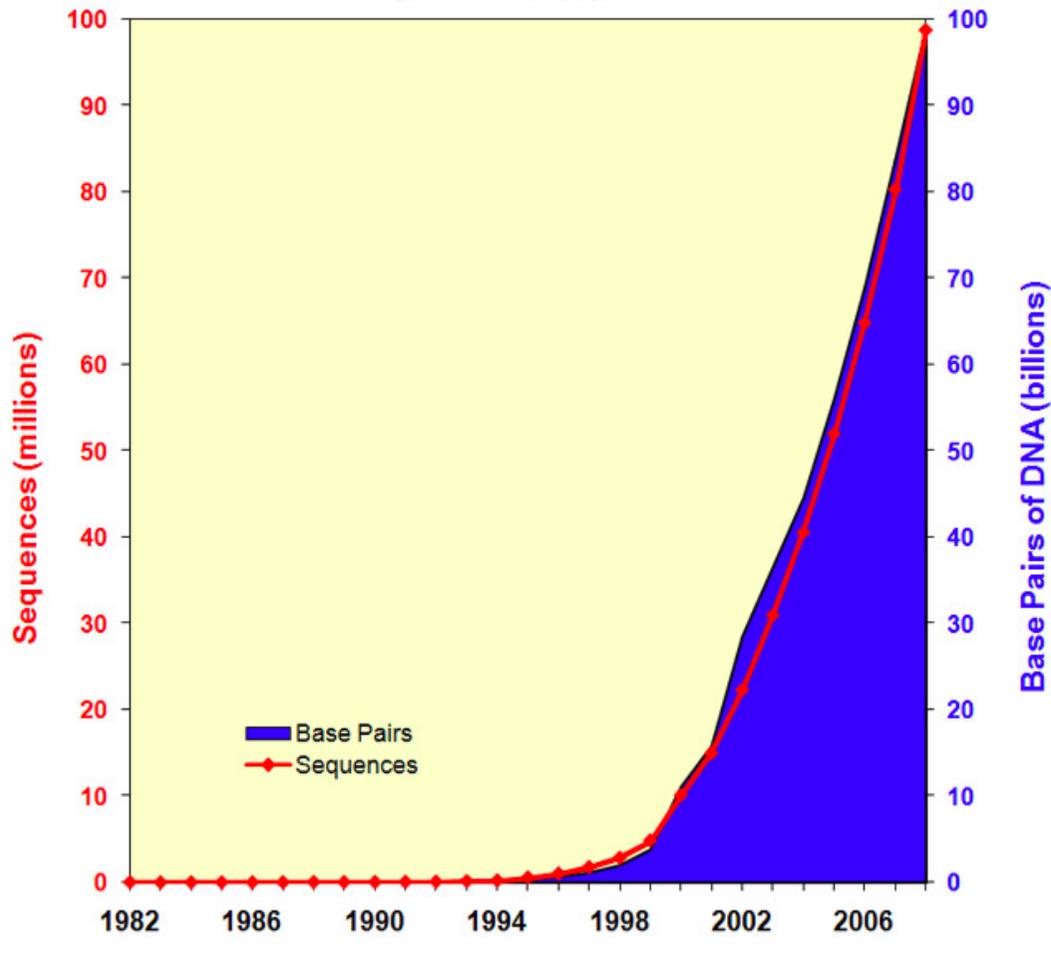
- ~ 1960s Margaret Dayhoff et ses collaborateurs collectent toutes les séquences de protéines connues et les regroupent dans l'Atlas of protein
 - 50 entrées
 - Version papier jusqu'en 78, puis version électronique
 - En 1984, l'Atlas devient PIR (Protein Information Resource)
- A partir de 1982 - EMBL (European Molecular Biology Laboratory) et NCBI (National Center for Biotechnology Information) : transcription et interprétation des séquences publiées dans les journaux au format électronique
 - Véritable explosion de la quantité de séquences disponibles
 - Rapidement, la DDBJ (DNA DataBank of Japan) se joint à la collaboration pour la collecte de données

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



- **1986 - SWISSPROT** (3900 prot.) : conversion du PIR en un format similaire à celui d'EMBL + ajout d'information pour chaque protéine
⇒ Notoriété de SWISSPROT comme banque de haute qualité
- Collaboration entre SWISSPROT et EMBL
But : mettre à disposition rapidement les séquences de protéines non encore annotées par SWISSPROT
⇒ **TrEMBL** (Translation of **EMBL** nucleotide sequences)
- **Fin 2003 - UNIPROT** (**UNI**versal **PROT**ein Ressource) : jonction des informations de SWISSPROT, TrEMBL et PIR
 - UniProtKB (*UniProt KnowledgeBase*) = SWISSPROT + TrEMBL
 - UniRef (*UniProt Reference Clusters*) : regroupement en 1 entrée des séquences similaires pour accélérer les recherches
 - UniParc (*UniProt Archive*) : utilisé pour garder une trace des séquences et de leur identifiants

Growth of GenBank
(1982 - 2008)



- En 1982 :
606 séquences
2 427 bases
- En décembre 2008 :
98 868 465 séquences
99 116 431 942 bases
- En décembre 2017 :
206 293 625 séquences
249 722 163 594 bases

[<https://www.ncbi.nlm.nih.gov/genbank/genbankstats-2008/>]

- Recherche publique donc données publiques
 - Partage international des données
 - Évite qu'une expérience soit refaite par différents laboratoires
 - Permet la comparaison des résultats
- Beaucoup de données générées par l'expérimentation
 - Outils adaptés à de grands volumes de données
 - Gestion des grandes banques par des organismes spécialisés
- Accès fréquent à ces données
 - Interrogation *via* internet
- Une source d'informations pour l'interprétation
 - Ma séquence est-elle déjà présente dans les banques ?
 - Quelles sont les séquences similaires à la mienne ?
 - Quelles sont les protéines similaires à la mienne (structure, fonction) ?
 - ...

- Quelle est la qualité des données ?
 - Erreurs de frappe
 - Erreurs dans les séquences
 - Erreurs d'annotation
 - Redondance des données
- Où sont les données ?
 - Dans les laboratoires
 - Dans des banques ou bases de données spécialisées : souvent maintenues par un laboratoire
 - Dans des banques de données généralistes : maintenues par des consortiums
- Différentiations possibles des banques de données
 - Banque primaire ou généraliste
 - Banque secondaire ou spécialisée

- Banques primaires/généralistes ~ archives
- Banques secondaires/spécialisées ~ données vérifiées (corrigées et annotées)
- La plus grande contribution des banques de données à la communauté des biologistes est de rendre les séquences **accessibles**
- Les banques primaires contiennent majoritairement des résultats expérimentaux (avec qqs interprétations), mais qui ne sont **pas vérifiés, ni analysés**
 - *les séquences nucléiques d'EMBL/GenBank/DDBJ sont issues du séquençage d'une molécule biologique qui existe dans un tube à essai, qqpart dans un labo ; elles ne représentent pas les séquences qui sont un consensus d'une population*
- Ceci a des conséquences sur l'interprétation de l'analyse des séquences : **informations parfois très utiles ou trompeuses**

- Ces banques contiennent des données hétérogènes (collecte la plus exhaustive possible)
- Banques de séquences nucléiques (GenBank, EMBL, DDBJ)
- Banques de séquences protéiques (PIR, SWISSPROT, UNIPROT)
- Banques d'articles scientifiques (PubMed)

Avantage : tout est consultable en une fois

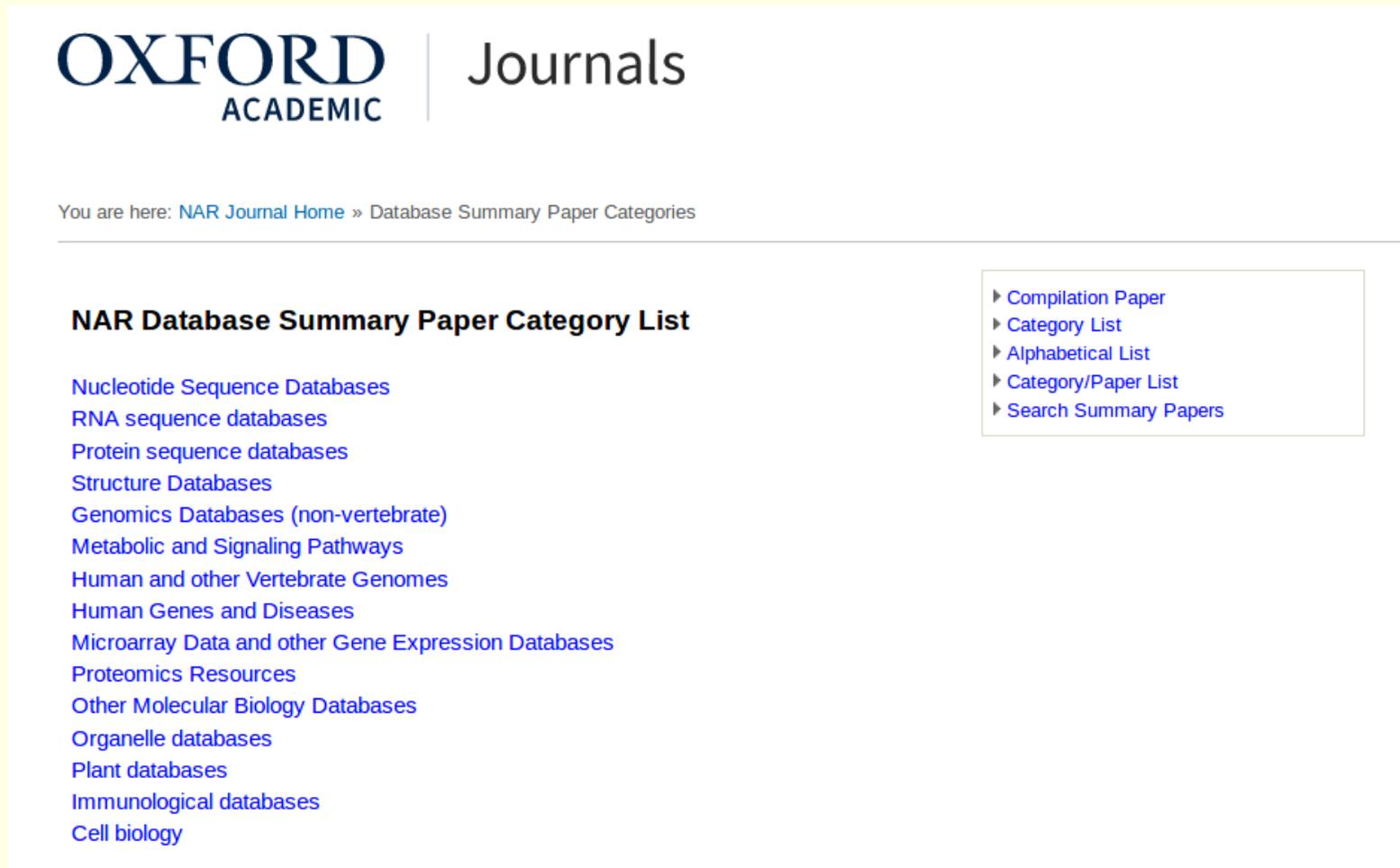
Inconvénient : difficiles à maintenir, difficiles à interroger

- Ces banques contiennent des données homogènes
- Collecte établie autour d'une thématique particulière telle que :
 - Banques de localisation (*mapping databases* GeneLoc)
 - Banques de structures 3D de macromolécules (PDB)
 - Banques génomiques (UCSC, Ensembl)
 - Banque spécialisée pour un génome spécifique, banques de séquences immunologiques, de voies métaboliques, de cartes génétiques, de motifs protéiques, d'expression de gènes, de structures, ...

Avantage : facilité pour mettre à jour les données, vérifier leur intégrité, offrir une interface adaptée, ...

Inconvénient : ne cible pas toujours exactement ce que l'on veut, toutes les banques possibles n'existent pas

- La revue Nucleic Acids Research consacre un numéro spécial "database issue" chaque année
(<http://www.oxfordjournals.org/nar/database/c/>).



OXFORD
ACADEMIC | Journals

You are here: [NAR Journal Home](#) » Database Summary Paper Categories

NAR Database Summary Paper Category List

- [Nucleotide Sequence Databases](#)
- [RNA sequence databases](#)
- [Protein sequence databases](#)
- [Structure Databases](#)
- [Genomics Databases \(non-vertebrate\)](#)
- [Metabolic and Signaling Pathways](#)
- [Human and other Vertebrate Genomes](#)
- [Human Genes and Diseases](#)
- [Microarray Data and other Gene Expression Databases](#)
- [Proteomics Resources](#)
- [Other Molecular Biology Databases](#)
- [Organelle databases](#)
- [Plant databases](#)
- [Immunological databases](#)
- [Cell biology](#)

- ▶ [Compilation Paper](#)
- ▶ [Category List](#)
- ▶ [Alphabetical List](#)
- ▶ [Category/Paper List](#)
- ▶ [Search Summary Papers](#)

-
- Introduction
 - **Banques de données de séquences nucléiques**
 - Banques de données de séquences protéiques
 - Banques de données spécialisées
 - Interrogation des banques de données
 - Formats de fichiers en bioinformatique
 - Références

- Les trois principales banques :
 - EMBL (Europe, 82)
 - GenBank (É-U, 82)
 - DDBJ (Japon, 86)
- Échange quotidien des données entre ces banques depuis 88
- Répartition de la collecte des données : chaque banque collecte les données de son continent
- Mise à jour : nouvelles version disponibles plsr fois par an (date et num de version), mise à disposition des “Updates”
- Format de stockage similaire : les fichiers (*flatfiles*) représentent l’unité d’information élémentaire

- Chaque *entrée* (séquence + informations) est stockée dans un fichier (*flatfile*)
- Ces fichiers se composent de trois parties :
 - Entête (*header*) : description générale de l'entrée
 - Les caractéristiques (*features*) : objets biologiques présents sur la séquence
 - La séquence elle-même
- 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 //

```

ID   U49845; SV 1; linear; genomic DNA; STD; FUN; 5028 BP.
XX
AC   U49845;
XX
DT   07-MAY-1996 (Rel. 47, Created)
DT   17-APR-2005 (Rel. 83, Last updated, Version 4)
XX
DE   Saccharomyces cerevisiae TCP1-beta gene, partial cds; and Axl2p (AXL2) and
DE   Rev7p (REV7) genes, complete cds.
XX
KW   .
XX
OS   Saccharomyces cerevisiae (baker's yeast)
OC   Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
OC   Saccharomycetales; Saccharomycetaceae; Saccharomyces.
XX
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/DBJ databases.
RL   Terry Roemer, Biology, Yale University, New Haven, CT, USA
XX
FH   Key          Location/Qualifiers
FH
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

```

```

LOCUS      SCU49845                5028 bp    DNA     linear   PLN 21-JUN-1999
DEFINITION Saccharomyces cerevisiae TCP1-beta gene, partial cds; and Axl2p
            (AXL2) and Rev7p (REV7) genes, complete cds.
ACCESSION  U49845
VERSION   U49845.1  GI:1293613
KEYWORDS   .
SOURCE     Saccharomyces cerevisiae (baker's yeast)
  ORGANISM Saccharomyces cerevisiae
            Eukaryota; Fungi; Ascomycota; Saccharomycotina; Saccharomycetes;
            Saccharomycetales; Saccharomycetaceae; Saccharomyces.
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
FEATURES   Location/Qualifiers
     source          1..5028
                    /organism="Saccharomyces cerevisiae"
                    /mol_type="genomic DNA"
                    /db_xref="taxon:4932"
                    /chromosome="IX"
                    /map="9"
     CDS             <1..206
                    /codon_start=3
                    /product="TCP1-beta"
                    /protein_id="AAA98665.1"
                    /db_xref="GI:1293614"
                    /translation="SSIYNGISTSGLDLNNGTIADMRQLGIVESYKLRVSSASEA
                    AEVLLRVDNIIRARPRTANRQHM"
     gene           687..3158
                    /gene="AXL2"
     CDS             687..3158
                    /gene="AXL2"

```

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)
	and
DE	Rev7p (REV7) genes, complete cds.

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

Utilisée dans les 3 principales banques de séquences nucléiques

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

- 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 ; Saccharomycetes ;
OC	Saccharomycetales ; Saccharomycetaceae ; Saccharomyces.

SOURCE	Saccharomyces cerevisiae (baker's yeast)
ORGANISM	Saccharomyces cerevisiae Eukaryota ; Fungi ; Dikarya ; Ascomycota ; Saccharomycotina ; Saccharomycetes ; Saccharomycetales ; Saccharomycetaceae ; Saccharomyces.

- 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
```

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 <i>Saccharomyces cerevisiae</i>
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

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 supplémentaire 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

```
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="SSIYNGISTSGLDLNNGTIADMRQLGIVESYKCLKRAVVSSASEAA
FT EVLLRVDNIIRARPRTANRQHM"
[... ]
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 DQIITETE VFDEF RSSLNSLIMHLEKLPKVND D TITFEAVINAI EELGHKLD RNR RVD
FT SLEEKAEIERDSN WVKCQEDENLPD NNGFQPPKIKL TSLVGSDV GPLIIHQFSEK LISG
FT DDKILNGVYSQYEEGESIFGSLF"
```

FEATURES	Location/Qualifiers
source	1..5028 /organism="Saccharomyces cerevisiae" /mol_type="genomic DNA" /db_xref="taxon :4932" /chromosome="IX" /map="9"
CDS	<1..206 /codon_start=3 /product="TCP1-beta" /protein_id="AAA98665.1" /db_xref="GI :1293614" /translation="SSYNGISTSGLDLNGTIADMRQLGIVESYKLGKRAVVSSASEA AEVLLRVDNIIRARPRTANRQHM"
[...]	
gene	complement(3300..4037) /gene="REV7"
CDS	complement(3300..4037) /gene="REV7" /codon_start=1 /product="Rev7p" /protein_id="AAA98667.1" /db_xref="GI :1293616" /translation="MNRWVEKWL RVYLKCYINLILFYRNVYPPQSFDYTTYQSFNLPQ FVPINRHPALIDYIEELILDVLSKLT HVYRFSICIINKKNDLCIEKYVLD FSELQHVD KDDQIITETEVFDEF RSSLNSLIMHLEKLPKVNDTITFEAVINAI EELGHKLDNR RVDSLEEKAEIERDSN WVKCQEDENLPDNNGFQPPKIKLTSLVGSDVGPLIIHQFSEK LISGDDKILNGVYSQYEEGESIFGSLF"

3. La séquence de nucléotides

→ EMBL

```
SQ Sequence 5028 BP; 1510 A; 1074 C; 835 G; 1609 T; 0 other;
gatcctccat atacaacggt atctccacct caggtttaga tctcaacaac ggaaccattg      60
ccgacatgag acagttagggt atcgtcgaga gttacaagct aaaacgagca gtagtcagct     120
ctgcatctga agccgctgaa gttctactaa ggggtgataa catcatccgt gcaagaccaa     180
gaaccgcaa tagacaacat atgtaacata tttaggatat acctcgaaaa taataaacg      240
ccacactgtc attattataa ttagaaacag aacgcaaaaa ttatccacta tataattcaa     300
agacgcgaaa aaaaaagaac aacgcgcat agaacttttg gcaattcgcg tcacaaataa     360
atthtgcaa cttatgtttc ctcttcgagc agtactcgag ccctgtctca agaatgtaat     420
aataccatc gtaggtatgg ttaaagatag catctccaca acctcaaagc tccttgccga     480
gagtcgccct ctttgctcga gtaattttca cttttcatat gagaacttat tttcttattc     540
```

→ GenBank

```
ORIGIN
1 gatcctccat atacaacggt atctccacct caggtttaga tctcaacaac ggaaccattg
61 ccgacatgag acagttagggt atcgtcgaga gttacaagct aaaacgagca gtagtcagct
121 ctgcatctga agccgctgaa gttctactaa ggggtgataa catcatccgt gcaagaccaa
181 gaaccgcaa tagacaacat atgtaacata tttaggatat acctcgaaaa taataaacg
241 ccacactgtc attattataa ttagaaacag aacgcaaaaa ttatccacta tataattcaa
301 agacgcgaaa aaaaaagaac aacgcgcat agaacttttg gcaattcgcg tcacaaataa
361 atthtgcaa cttatgtttc ctcttcgagc agtactcgag ccctgtctca agaatgtaat
421 aataccatc gtaggtatgg ttaaagatag catctccaca acctcaaagc tccttgccga
481 gagtcgccct ctttgctcga gtaattttca cttttcatat gagaacttat tttcttattc
541 ttactctca catctgtag tgattgacac tgcaacagcc accatcacta gaagaacaga
```

- Répondez aux questions suivantes en vous référant aux 2 fiches distribuées
- Une des fiches vient d'EMBL, l'autre de GenBank, identifiez-les
- Trouvez le nom des organismes d'où proviennent les séquences
- Que sont ces séquences ?
- À quelle division d'organisme appartiennent-elles ?
- Quelle séquence a été entrée le + récemment ? Ont-elles été modifiées ?
- Qui a soumis ces séquences ?
- De quel type de molécule s'agit-il ?
- Peut-on localiser ces séquences sur un chromosome ?

-
- Introduction
 - Banques de données de séquences nucléiques
 - Banques de données de séquences protéiques
 - Banques de données spécialisées
 - Interrogation des banques de données
 - Formats de fichiers en bioinformatique
 - Références

- Origine des séquences

- Traduction automatique de séquences d'ADN (majoritairement)
- Séquençage de protéines (rare car long et coûteux)
- Protéines dont la structure 3D est connue

- Origine des annotations

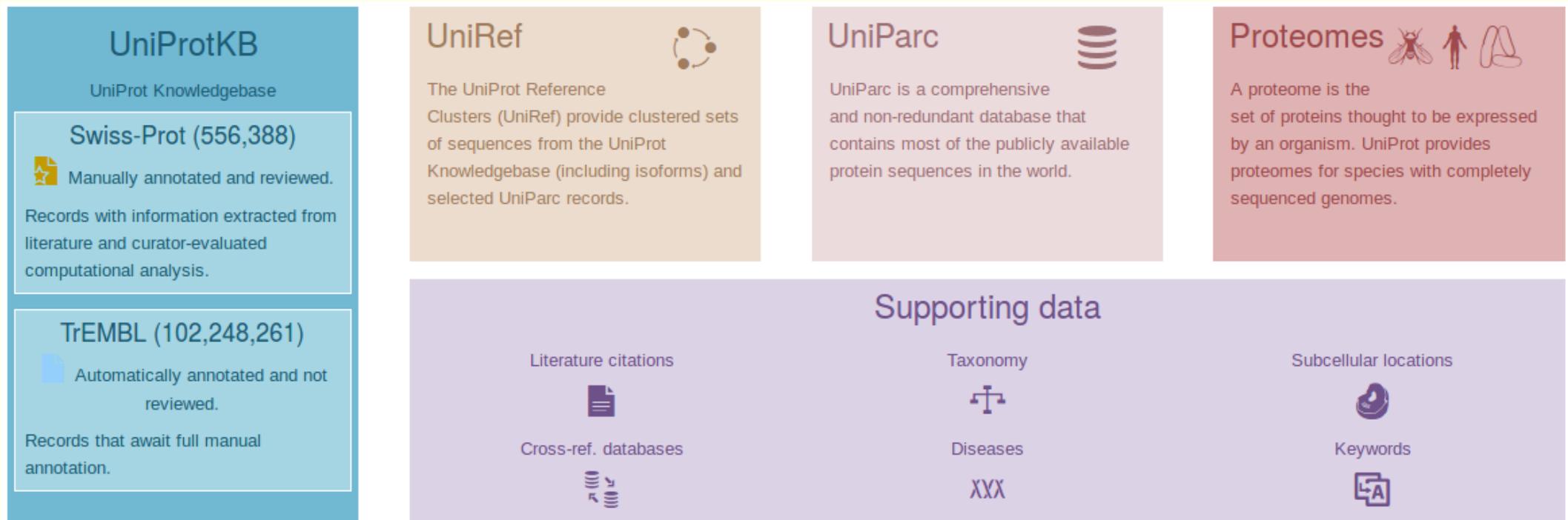
- Spectrométrie de masse : régulation, rythme et localisation de l'expression des protéines ; mais aussi identification et modification post-transcriptionnelle
- Études d'interactions : comment les protéines s'assemblent entre elles ou avec d'autres molécules pour former des complexes moléculaires
- Cristallographie et résonance magnétique nucléaire : pour déterminer la forme 3D finale de la protéine

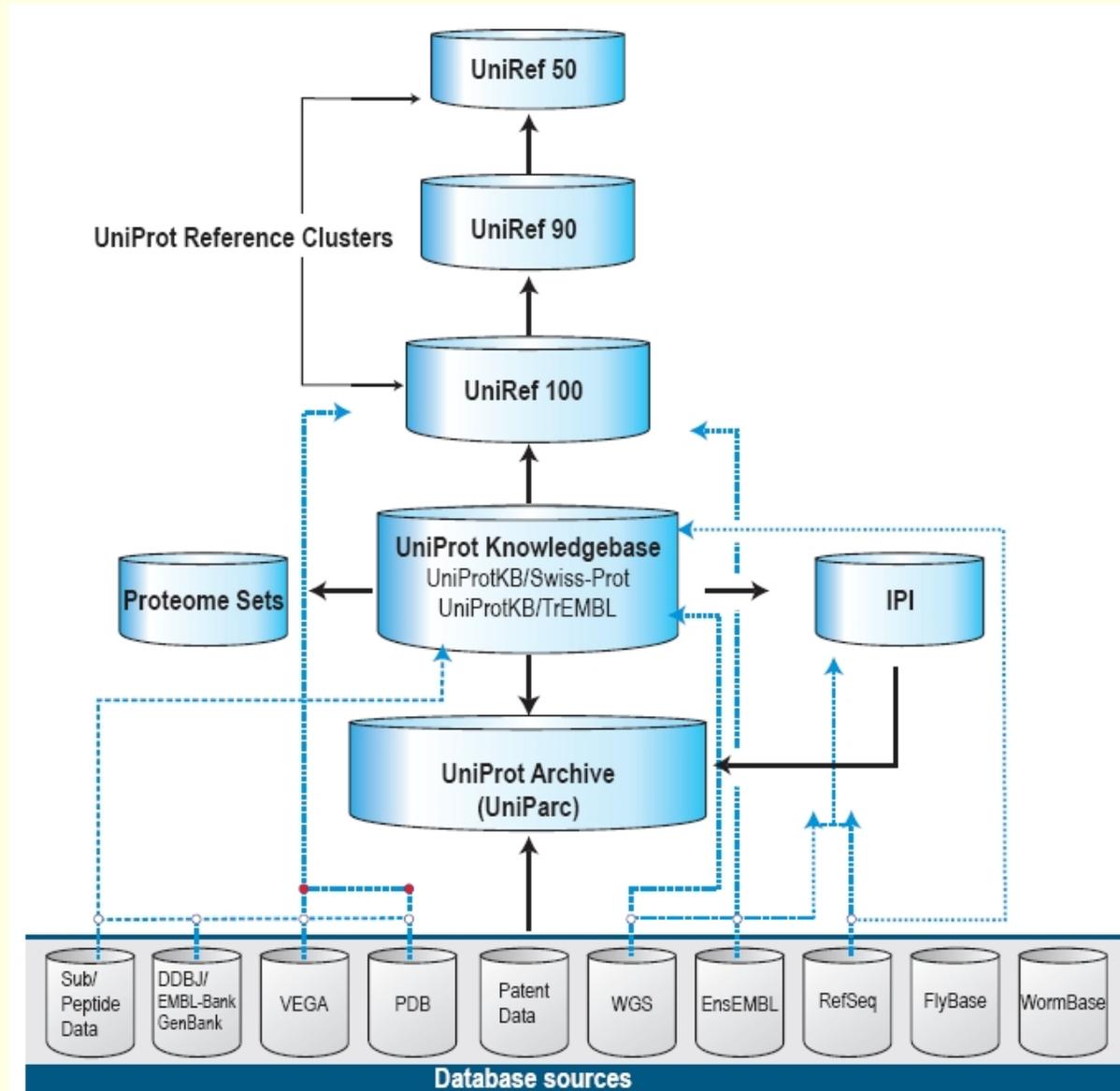
- Les **données stockées** : séquences + annotations (protéines entières ou fragment de protéines)
 - Banques généralistes : protéines de toutes les espèces
 - Banques spécialisées : familles de protéines particulières, groupes de protéines ou d'un organisme particulier
- **GenPept** : entrées = traductions des séquences de DDBJ/EMBL/GenBank (champ CDS) ; les annotations sont les mêmes ; la banque n'est pas vérifiée ~ *archive basique*
- **RefSeq** : projet NCBI ; but : fournir une vue d'ensemble, intégrée et non-redondante de séquences d'ADN, d'ARN et de protéines ; lien explicite entre les séq. nucléiques et protéiques ; numéro d'accèsion particulier format **2 + 6**, ex : NT_123456
- **UniProt** : assemblage de SWISSPROT, TrEMBL et PIR

- **Objectif** : proposer les séquences les plus représentatives et décrites par les experts
- **Avantages**
 - Entrées non redondante
 - Liens explicites entre les séquences nucléiques et protéiques
 - Mise à jour régulière par le personnel du NCBI (statut indiqué)
 - Validation des données et consistance des formats
 - Synthèse des informations issues de plusieurs entrées nucl. ou prot.
- Différents **niveaux de correction** des données (champ **COMMENT**)
 - Reviewed : revu par un membre du NCBI (ajout d'informations)
 - Validated : une 1^{ère} révision a été effectuée par un membre du NCBI (annotation en cours)
 - Provisional : entrée non lue par un annotateur (contient sûrement un vrai transcrit ou une vrai protéine)
 - Predicted : transcrit ou protéine issu d'une prédiction (prog info)

- Un numéro d'accèsion au format spécifique : un préfixe à 2 lettres + un numéro à 6 chiffres (convention de nommage)
- Liste des **préfixes** des numéros d'accèsion de RefSeq
 - NC_ : Complete genomic molecules
 - NG_ : Incomplete genomic region
 - NM_ : mRNA
 - NR_ : ncRNA
 - NP_ : Protein
 - XM_ : predicted mRNA model
 - XR_ : predicted ncRNA model
 - XP_ : predicted Protein model (eukaryotic sequences)
 - WP_ : predicted Protein model (prokaryotic sequences)

”The mission of UniProt is to provide the scientific community with a comprehensive, high-quality and freely accessible resource of protein sequence and functional information”





Sources and flow of data for UniProt's component databases

- **SWISSPROT**

- Données corrigées et validées par des experts
- Haut niveau d'annotation
- Redondance minimale
- Nombreux liens vers d'autres banques (~ 60)

- **TrEMBL**

- Entrées supplémentaires à SWISSPROT (pas encore annotées)
- Traduction automatique des CDS d'EMBL et soumissions spontanées
- Annotation automatique des protéines

UniProtKB

Advanced Search

BLAST Align Retrieve/ID mapping Peptide search Help Contact

UniProtKB - P04585 (POL_HV1H2)

Basket

Display

Entry Publications Feature viewer Feature table

BLAST Align Format Add to basket History

Feedback Help video Other tutorials and videos

Protein | **Gag-Pol polyprotein**

Gene | **gag-pol**

Organism | *Human immunodeficiency virus type 1 group M subtype B (isolate HXB2) (HIV-1)*

Status | Reviewed - Annotation score: ●●●●● - Experimental evidence at protein levelⁱ

Functionⁱ

Gag-Pol polyprotein: Mediates, with Gag polyprotein, the essential events in virion assembly, including binding the plasma membrane, making the protein-protein interactions necessary to create spherical particles, recruiting the viral Env proteins, and packaging the genomic RNA via direct interactions with the RNA packaging sequence (Psi). Gag-Pol polyprotein may regulate its own translation, by the binding genomic RNA in the 5'-UTR. At low concentration, the polyprotein would promote translation, whereas at high concentration, the polyprotein would encapsidate genomic RNA and then shut off translation.

Matrix protein p17: Targets the polyprotein to the plasma membrane via a multipartite membrane-binding signal, that includes its myristoylated N-terminus (By similarity). Matrix

Indication de la fiabilité des annotations - "Annotation score"
(SWISSPROT / TrEMBL)

-
- Introduction
 - Banques de données de séquences nucléiques
 - Banques de données de séquences protéiques
 - **Banques de données spécialisées**
 - Interrogation des banques de données
 - Formats de fichiers en bioinformatique
 - Références

- Localisation génomique :
 1. “Cartographie” : carte physique, carte cytogénétique, liaison génétique (*genetic linkage*), ...
 2. “Identification” : trouver où sont les objets d'intérêt biologiques comme les gènes, les variations génétiques ou les locus de prédisposition à des maladies \Rightarrow association d'une signature moléculaire à un résultat biologique
- Ex : localiser des nouveaux gènes ou affiner des régions d'intérêt
- Genome Database (GDB), eGenome, LBD2000, GeneCards et GeneLoc, GeneLynx, EuGenes, AceView, ...
- Cartes comparatives entre plusieurs génomes (Mouse Genome Informatics Database (MGI), ...)

- **1996** 1^{re} séquence complète d'un génome eucaryote
Saccharomyces cerevisiae, chromosomes entre 270 et 1500 Kb ;
limite d'une entrée dans GenBank à cette époque : 350 Kb
⇒ Mise en place d'une **section spécifique pour les génomes**,
1^{re} vue graphique des séquences génomiques
- **2001** première ébauche du génome humain ; chromosomes entre
46 et 246 **Mb**
- “Navigateur” de génomes : **NCBI Map Viewer**, **UCSC Genome Browser**, **Ensembl** (EBI et Sanger Institute)
- Attention au version d'assemblage qui peuvent changer sans
avertissement

1. Register
Register your project information and Metadata in the Genomes Online Database
[Register](#)

2. Annotate
Annotate your microbial genome or metagenome with IMG/ER or IMG/MER
[Annotate](#)

3. Publish
Publish your genome or metagenome in open access standards-supportive journal.
[Publish](#)

Studies
Metagenomic [1,375](#)
Non-Metagenomic [30,008](#)

Biosamples
Classification
Ecosystems
Host-associated [21,399](#)
Engineered [4,271](#)
Environmental [18,093](#)

Sequencing Projects
Complete Projects [12,795](#)
Permanent Drafts [104,867](#)
Incomplete Projects [66,515](#)
Targeted Projects [1,237](#)

Analysis Projects
Genome Analysis [95,065](#)
Metagenome Analysis [27,920](#)
Metagenome - Cell Enrichment [935](#)
Metagenome - Single Particle Sort [2,959](#)
Metagenome - Assembled Genome (MAG) [5,066](#)
Metatranscriptome Analysis [3,375](#)
Combined Assembly [139](#)
Single Cell - Screened (SAG) [2,140](#)
Single Cell - Unscreened (SAG) [1,154](#)
Transcriptome Analysis [421](#)

Special Projects
Type Strain Projects [8,690](#)
Strains at Genbank [7,450](#)
GEBA Projects [3,203](#)
HMP Projects [2,912](#)

Projects with Genbank Data
Seq. Projects [90,731](#)
Archaeal Projects [812](#)
Bacterial Projects [77,144](#)
Eukaryal Projects [4,268](#)
Viral Projects [8,507](#)

JGI Projects
JGI Studies [1,266](#)
JGI Biosamples [14,087](#)
JGI Sequencing Projects [72,365](#)
JGI Analysis Projects [34,237](#)

Organisms
Organisms [289,784](#)
Archaea [2,439](#)
Bacteria [256,794](#)
Eukarya [21,646](#)
Viruses [8,876](#)
Bacterial Type Strains [15,497](#)
Archaeal Type Strains [560](#)

Prokaryotic Type Strains in GOLD
At Genbank (7 450) | With Projects (8 690) | Total in GOLD (16 057)

GOLD : Genome Online Database (<http://gold.jgi.doe.gov/>)
Informations concernant les projets de séquençage de génomes et métagénomiques (et les métadonnées associées)

-
- Regroupement des protéines ayant des fonctions identiques ou proches
 - Construites souvent par comparaison de toutes les protéines entre elles puis constitution de groupes
 - Nombreuses banques :
 - HomoloGene, COG (NCBI)
 - KEGG SSDB (Sequence Similarity DataBase)
 - Clustr (EBI)
 - HOGENOM, OrthoMaM (LBBE, ISEM)
 - Autres banques basées sur la structure 3D des protéines
 - SSF, SCOP, CATH, ...

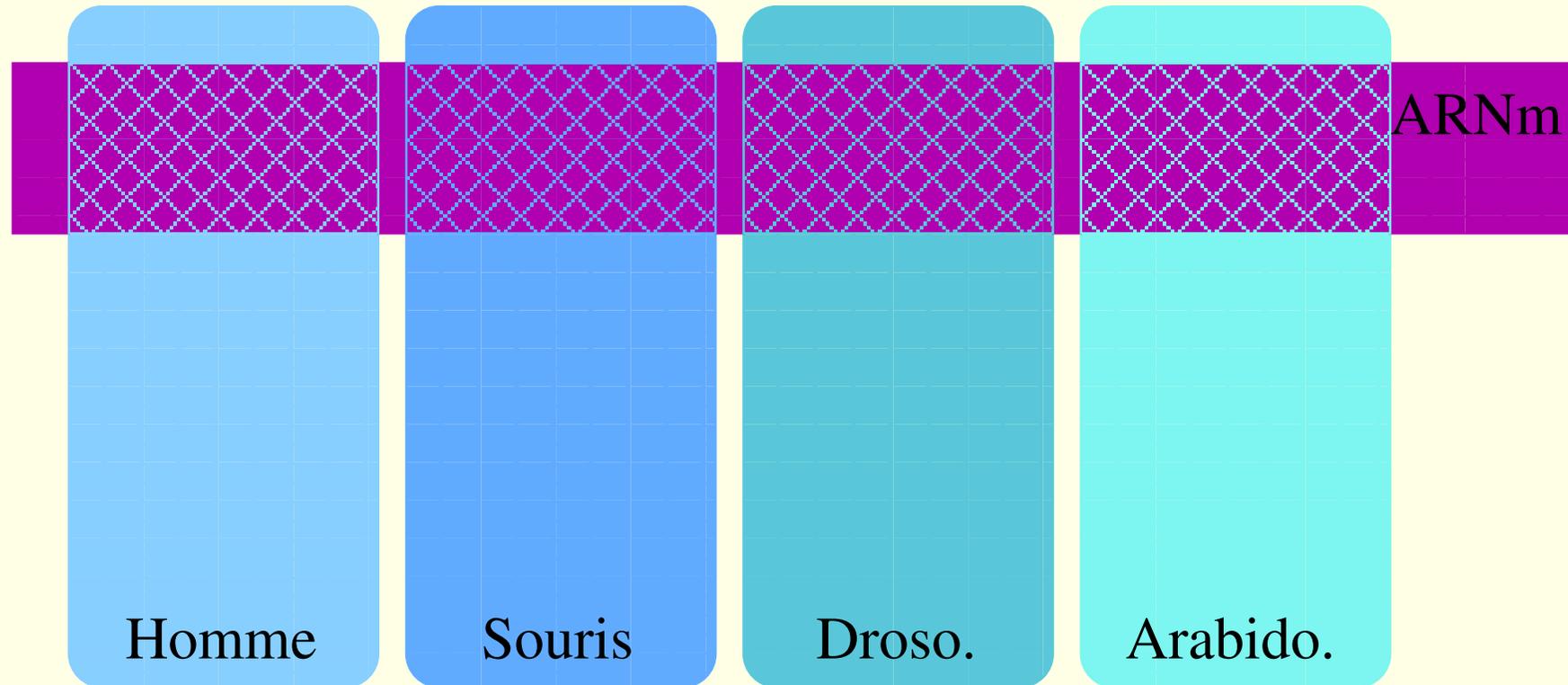
- Une famille de protéines peut-être caractérisée par un motif ou un domaine protéique
 - Séquence plus ou moins conservée importante pour la fonction des protéines de la famille
 - Déterminée à partir d'un alignement multiple
 - Plusieurs représentations possibles : consensus, expression régulière, alignement, matrices, HMM, ...
- Nombreuses banques
 - Prosite, PFAM, BLOCKS, Prodom, CDD, ...

-
- Introduction
 - Banques de données de séquences nucléiques
 - Banques de données de séquences protéiques
 - Banques de données spécialisées
 - Interrogation des banques de données
 - Formats de fichiers en bioinformatique
 - Références

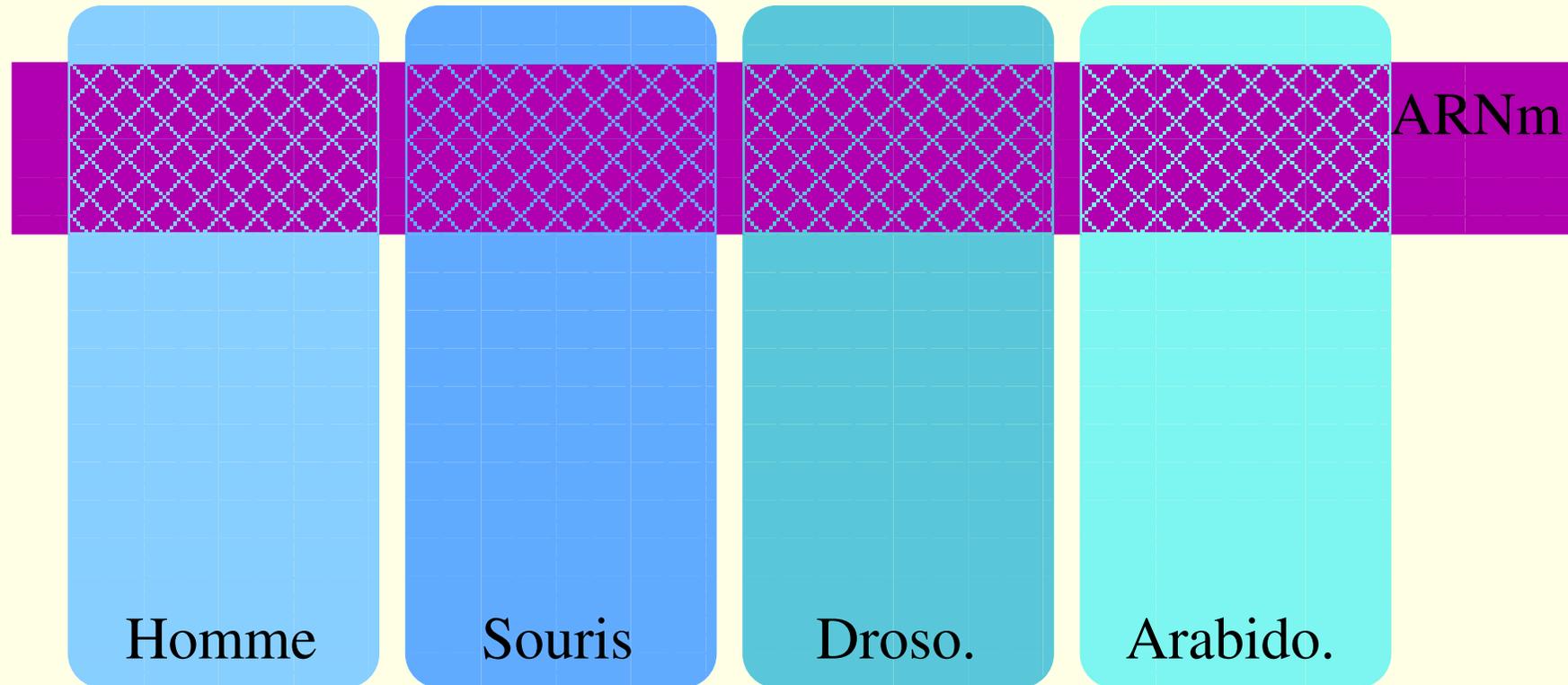
- Le réseau **internet** : interface web, API, programme client, ...
- Des **centres de ressources** : **NCBI** (<http://www.ncbi.nlm.nih.gov/>), **EBI** (<http://www.ebi.ac.uk/>), ...
- Des **catalogues d'outils** : **EMBOSS** (<http://emboss.sourceforge.net>), **Institut Pasteur** (<http://galaxy.pasteur.fr/>), ...
- Des **systèmes d'interrogation**
 - **GQuery** (anciennement Entrez) : développé et utilisé par le NCBI
 - **EB-Eye** ou **EBI Search** : développé et utilisé par l'EBI, utilisé par d'autres projets comme EnsemblGenome, Emboss, Interpro,...
 - **ACNUC** développé et utilisé par le PRABI-Doua : un des 1^{er} système d'interrogation pour les données biologiques
 - Systèmes spécifique comme **UniProt search engine** (~ EB-Eye)
 - **Base de données** relationnelles / NoSQL
 - ...

- **But**
 - Obtenir des information nouvelles et pertinentes
 - Aide à la mise au point d'expérience
 - Validation des résultats d'une expérience

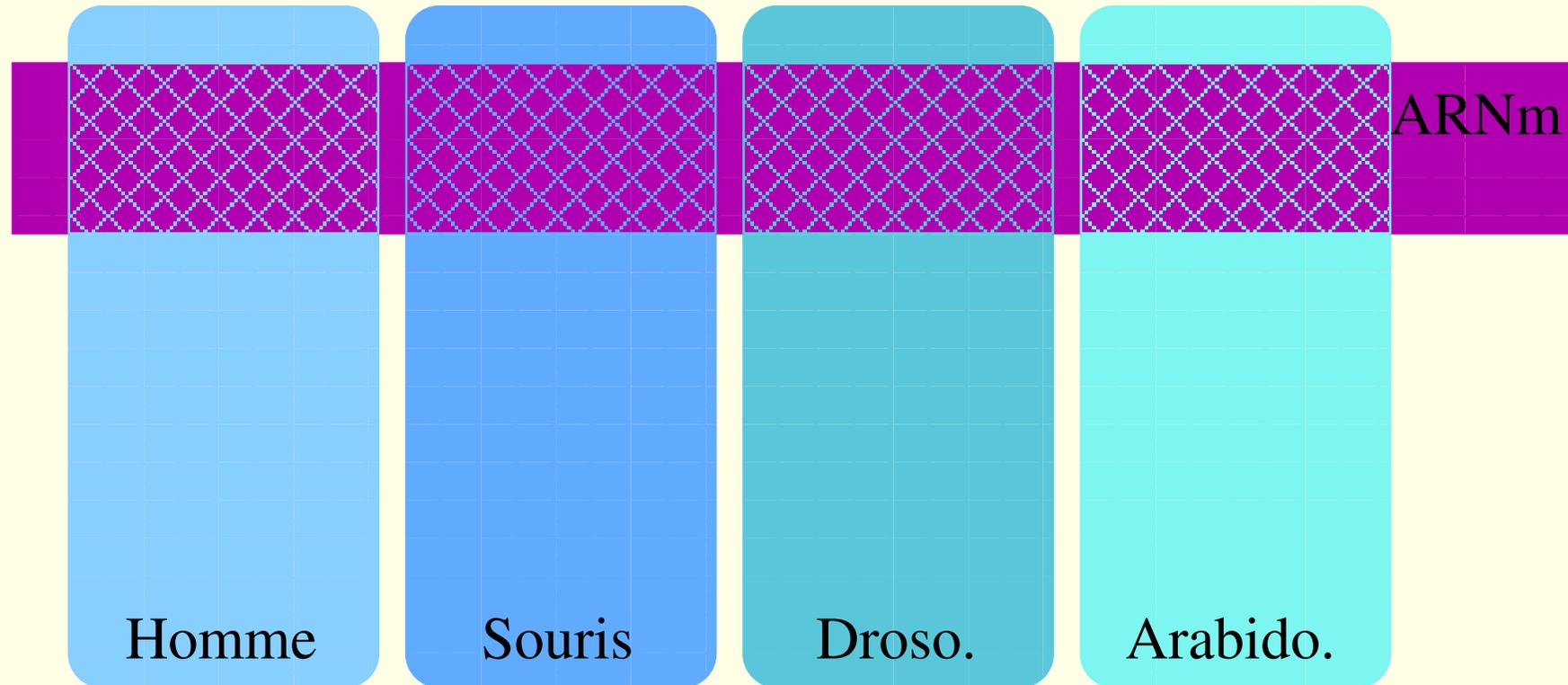
- **Contraintes** pour un système d'interrogation
 - Obtention de données pertinentes (pas trop de résultats mais tous ceux relatifs à notre problématique)
 - Simplicité d'utilisation (syntaxe d'interrogation intuitive)
 - Réponse rapide
 - Possibilité d'analyse des résultats (couplage à des outils)



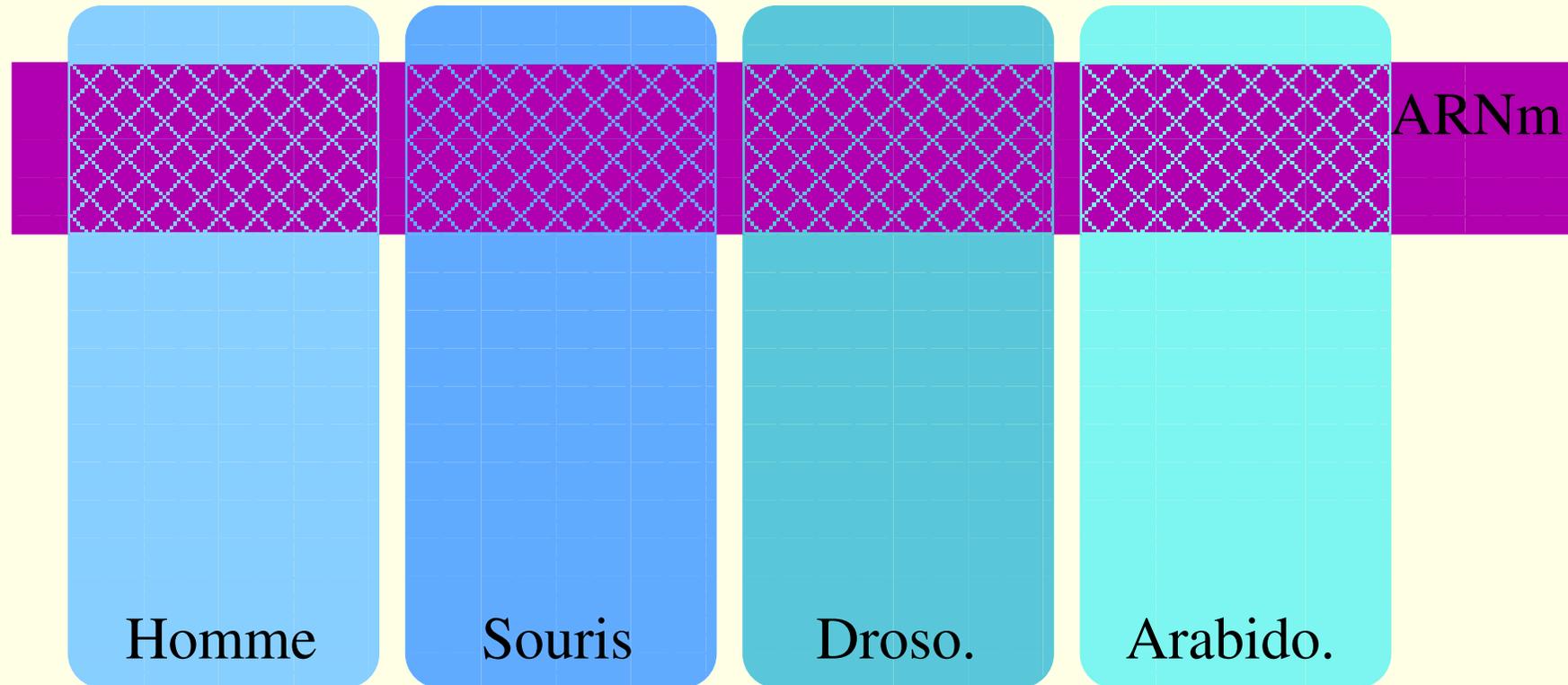
- Toutes séquences contenant 'Homme' :
- Séquences contenant 'Homme' ou 'Souris' :
- Tous les ARNm de la drosophile :
- Tous les ARNm sauf ceux d'arabidopsis :



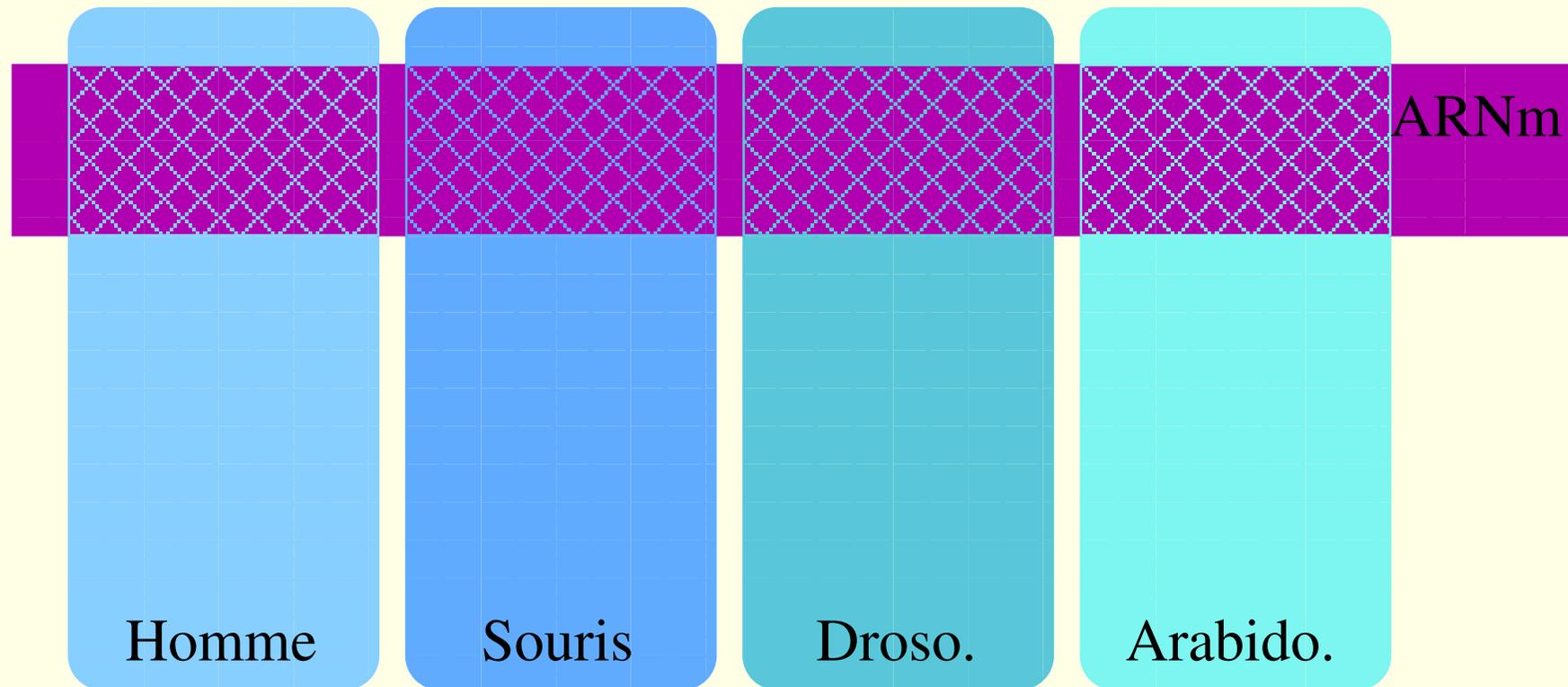
- Toutes séquences contenant 'Homme' : Homme
- Séquences contenant 'Homme' ou 'Souris' :
- Tous les ARNm de la drosophile :
- Tous les ARNm sauf ceux d'arabidopsis :



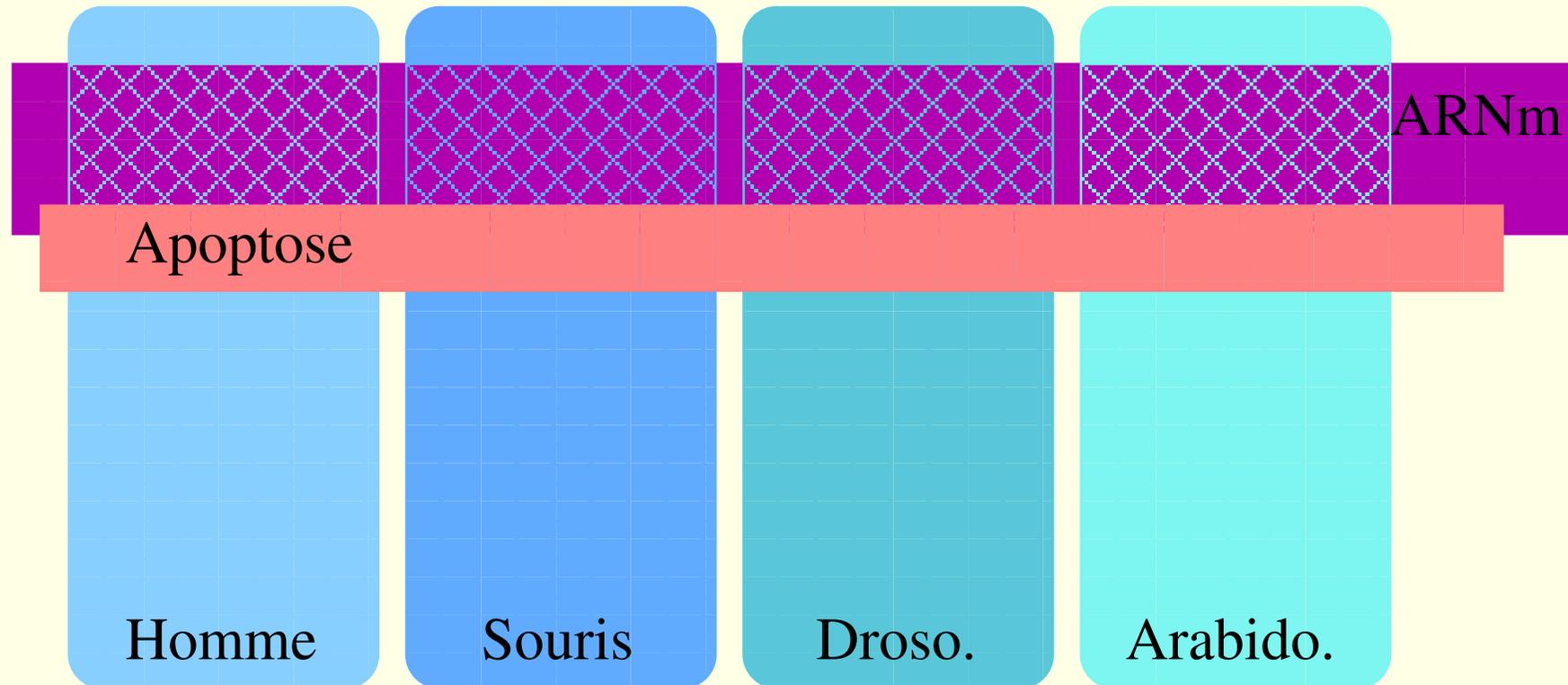
- Toutes séquences contenant 'Homme' : Homme
- Séquences contenant 'Homme' ou 'Souris' : Homme OU Souris
- Tous les ARNm de la drosophile :
- Tous les ARNm sauf ceux d'arabidopsis :



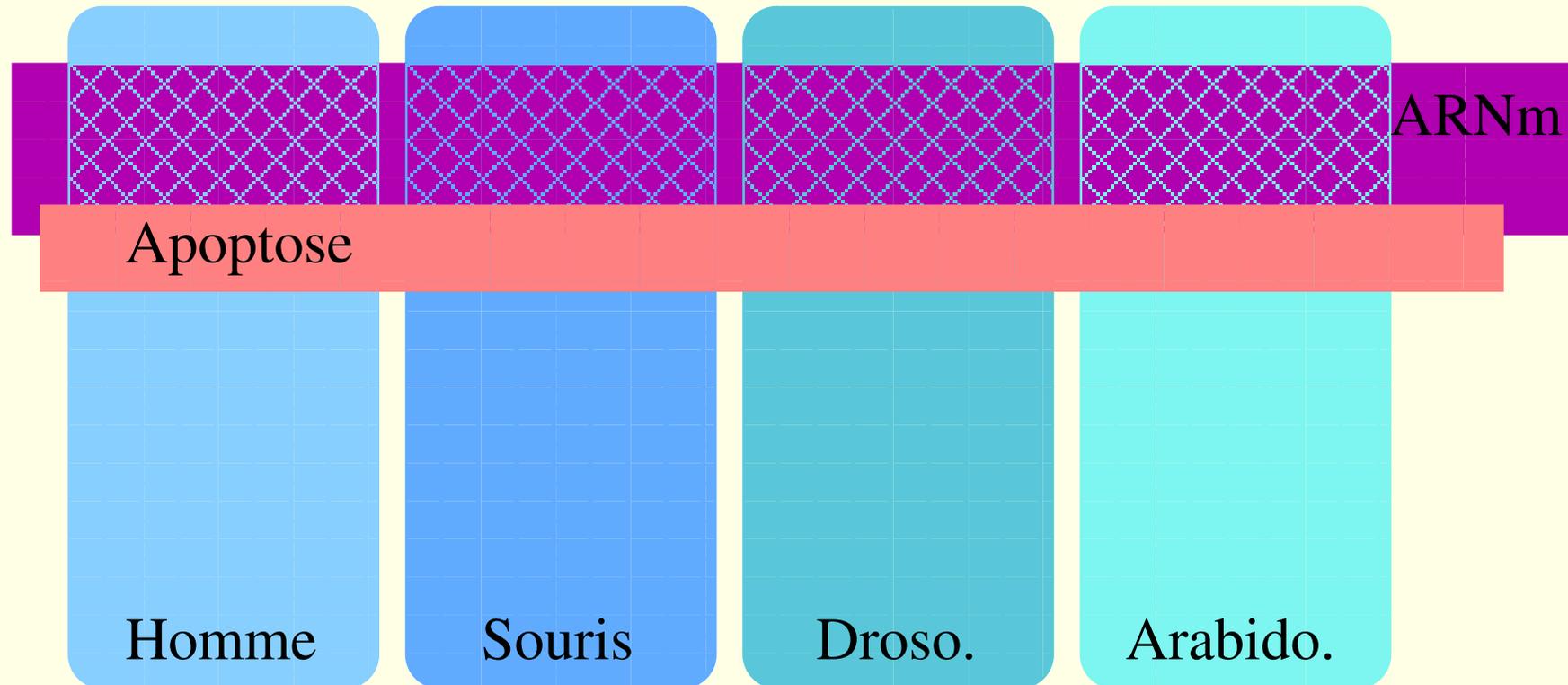
- Toutes séquences contenant 'Homme' : Homme
- Séquences contenant 'Homme' ou 'Souris' : Homme OU Souris
- Tous les ARNm de la drosophile : Droso. ET ARNm
- Tous les ARNm sauf ceux d'arabidopsis :



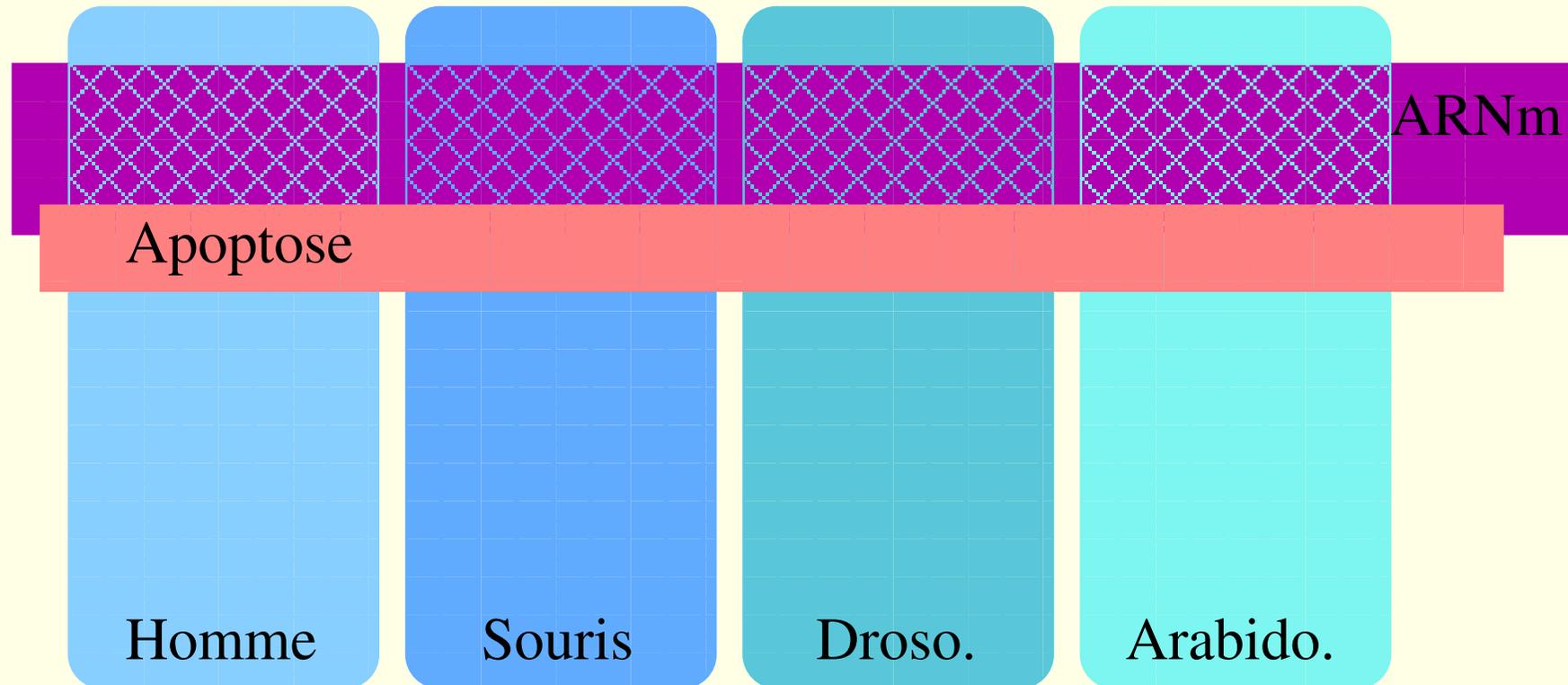
- Toutes séquences contenant 'Homme' : Homme
- Séquences contenant 'Homme' ou 'Souris' : Homme OU Souris
- Tous les ARNm de la drosophile : Druso. ET ARNm
- Tous les ARNm sauf ceux d'arabidopsis : ARNm NON Arabido.



- ARNm impliqués dans l'apoptose chez l'homme :
- Séquences de la souris ou de la droso. impliquées dans l'apoptose mais qui ne sont pas des ARNm :

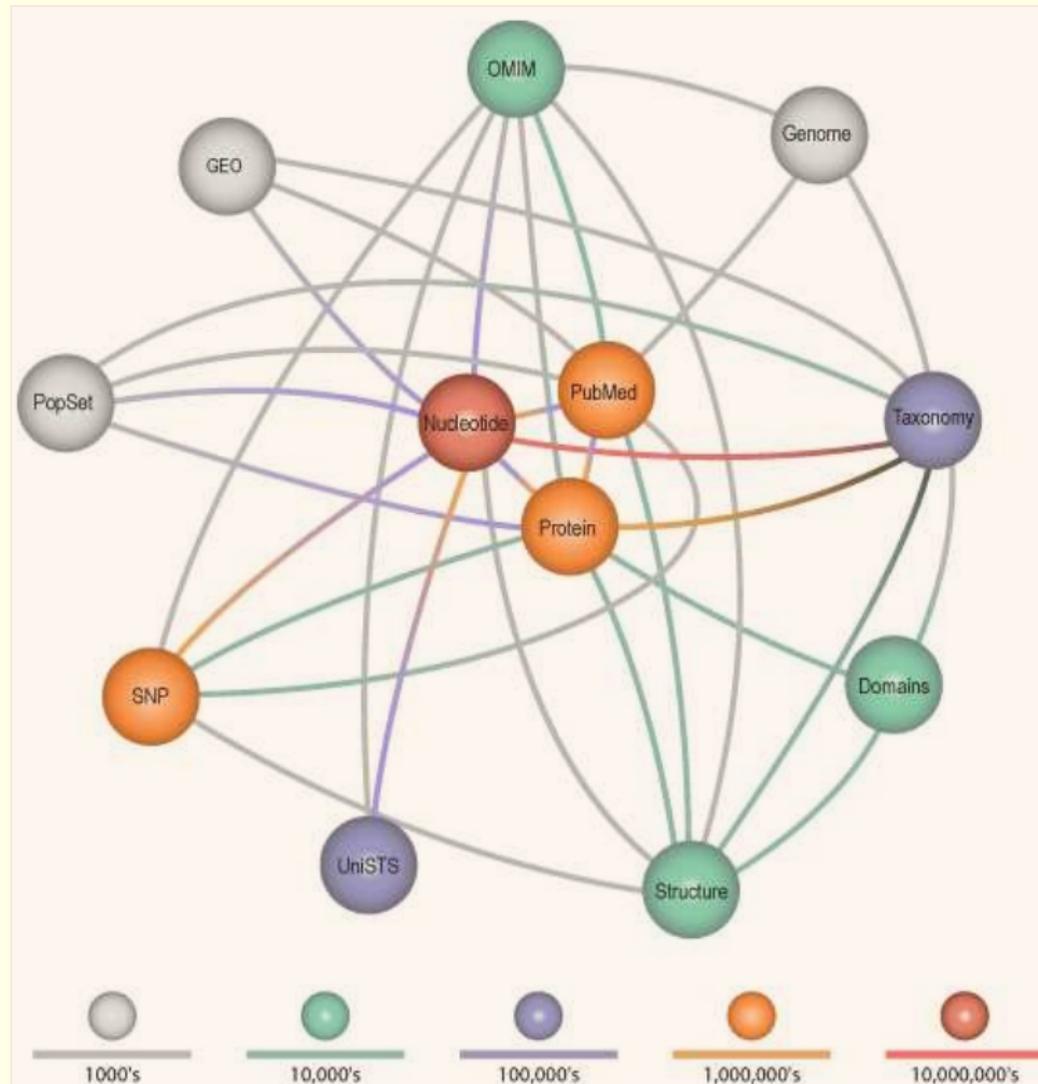


- ARNm impliqués dans l'apoptose chez l'homme :
ARNm ET Apoptose ET Homme
- Séquences de la souris ou de la droso. impliquées dans l'apoptose mais qui ne sont pas des ARNm :



- ARNm impliqués dans l'apoptose chez l'homme :
ARNm ET Apoptose ET Homme
- Séquences de la souris ou de la droso. impliquées dans l'apoptose mais qui ne sont pas des ARNm :
((Souris OU Droso.) ET Apoptose) NON ARNm

- Intégration "légère" de diverses banques de données au travers de la notion de référence croisée
→ choix de rubriques



- Interface **propriétaire** (ne peut être installée par autrui)
- Requête simple/manuelle (possibilité de mettre le nom du champs entre crochets)
Exemple : homo sapiens [organism]
- Requête avancée (choix des champs dans une liste déroulante)
- Différentes options
 - Aide dans **“Preview/Index”**
 - Historique (lien **“History”**)
 - Ajout de limites (lien **“Limits”**)
 - Sauvegarde, format : **“Display”**, **“send to”**, menus déroulants associés

- Terme
- Locution (groupe de terme) : " mot1 mot2 ... "
- Connecteurs logiques : ET (AND), OU (OR), NON (NOT)
- Caractère joker : *
- Expression complexe : parenthèse () (ET prioritaire)
- Couple descripteur (champ) / valeur en fonction des rubriques (ex : organism, gene, slen, author, title, ...) : valeur [descripteur]
- Exemples de requêtes
 - Human hemoglobin
 - Human [organism] hemoglobin* [protein name]
 - Human [organism] OR Mouse [organism]
 - 1000 :2000 [SLEN] (Human [organism] OR Mouse [organism])
 - 1000 :2000 [SLEN] NOT Human [organism]

- Syntaxe : valeur [descripteur]
- Descripteur Properties [prop]
 - type de molecule : biomol_genomic [prop], biomol_mrna [prop]
 - rubrique : srcdb_refseq [prop], srcdb_ddbj/embl/genbank [prop], srcdb_genbank [prop]
 - groupe d'organismes : gbdiv_pri [prop], gbdiv_rod [prop], gbdiv_pln [prop]
- Descripteur annotation [fkey] : exon [fkey] cds [fkey] 5'UTR [fkey] mRNA [fkey]
- Autres descripteurs en fonction de la rubrique : [organism], [gene], ...

- **Nucleotide** : séquences nucléiques
- **Gene** : synthèse sur les gènes
- **Protein** : séquences protéiques (Genpept, Swissprot)
- **CDD** : domaines protéiques fonctionnels (Conserv Domain protein database)
- **PubMed** : publications dans les journaux, revues, conférences (Medline)
- **Taxonomy** : classification du vivant
- **OMIM** : maladies génétiques humaines
- **RefSeq** : séquences de la banque RefSeq
- ...

Exemple de champs de la rubrique Gene

Nom du champ	Abréviation	Définition	Exemples
Symbol	SYM	nom normalisé du gène	SRY [SYM]
Chromosome	CHR	numéro chromosome portant le gène	9 [CHR]
Filter	Filter	lien avec une autre banque	gene_omim [FILTER]
Gene Ontology	GO	descripteur normalisé	cell adhesion [GO]
Title	TI	Terme dans la définition	binding [TI]
Properties	PROP	différentes propriétés	gene_in_mitochondrion[PROP]
Organism	ORGN	noms pour l'ensemble des taxons	human[ORGN]
All Fields	[ALL]	tous les champs	SRY [ALL]

GQuery, formulaire de requête et résultats d'une recherche sur toutes les rubriques

NCBI Resources How To Sign in to NCBI

Search NCBI databases [Help](#)

Homo sapiens

Results found in 35 databases for "Homo sapiens"

Literature

Books	63,968	books and reports
MeSH	5	ontology used for PubMed indexing
NLM Catalog	35,564	books, journals and more in the NLM Collections
PubMed	16,291,694	scientific & medical abstracts/citations
PubMed Central	2,007,759	full-text journal articles

Health

ClinVar	12	human variations of clinical significance
dbGaP	915	genotype/phenotype interaction studies
GTR	0	genetic testing registry
MedGen	5	medical genetics literature and links
OMIM	67	online mendelian inheritance in man
PubMed Health	12,135	clinical effectiveness, disease and drug reports

Genomes

Assembly	90	genome assembly information
BioProject	32,758	biological projects providing data to NCBI
BioSample	2,720,671	descriptions of biological source materials
Clone	17,567,281	genomic and cDNA clones
dbVar	5,041,963	genome structural variation studies
Genome	1	genome sequencing projects by organism
GSS	1,763,652	genome survey sequences
Nucleotide	22,846,232	DNA and RNA sequences
Probe	27,386,680	sequence-based probes and primers
SNP	165,297,523	short genetic variations
SRA	896,085	high-throughput DNA and RNA sequence read archive
Taxonomy	1	taxonomic classification and nomenclature catalog

Genes

EST	8,864,000	expressed sequence tag sequences
Gene	935,188	collected information about gene loci
GEO DataSets	1,277,435	functional genomics studies
GEO Profiles	61,958,910	gene expression and molecular abundance profiles
HomoloGene	18,732	homologous gene sets for selected organisms
PopSet	36,331	sequence sets from phylogenetic and population studies
UniGene	66,572	clusters of expressed transcripts

Proteins

Conserved Domains	83	conserved protein domains
Protein	52,158,336	protein sequences
Protein Clusters	14	sequence similarity-based protein clusters
Structure	34,497	experimentally-determined biomolecular structures

Chemicals

BioSystems	26,037	molecular pathways with links to genes, proteins and chemicals
PubChem BioAssay	377,185	bioactivity screening studies
PubChem Compound	0	chemical information with structures, information and links
PubChem Substance	0	deposited substance and chemical information

GQuery, formulaire de requête avancé sur une rubrique (ex : Nucleotide)

55

NCBI Resources How To Sign in to NCBI

Nucleotide Home Help

Nucleotide Advanced Search Builder

homo sapiens[Organism] [Edit](#) [Clear](#)

Builder

Organism [-](#) [Show index list](#)

AND [-](#) [+](#) [Show index list](#)

[Search](#) or [Add to history](#)

History [Download history](#) [Clear history](#)

Search	Add to builder	Query	Items found	Time
#4	Add	Search Homo sapiens	13456380	11:53:12
#1	Add	Search homo sapiens [ORGN]	10408382	11:29:52

GQuery, résultats d'une recherche sur une rubrique (ex : Nucleotide)

The screenshot shows the NCBI Nucleotide search interface. At the top, the search criteria are 'Nucleotide' and 'homo sapiens [ORGN]'. The search results indicate that 208,423 nucleotide sequences were found. The results are displayed in a list format, with the first 8 items shown. Each item includes a checkbox, a link to the sequence, and its accession number and GI number. The right sidebar contains several sections: 'Filter your results' with a dropdown for 'All (10408382)', 'Top Organisms' with a tree view, 'Find related data' with a database selector, 'Search details' with a search box, and 'Recent activity' with a list of recent searches.

NCBI Resources How To Sign in to NCBI

Nucleotide Nucleotide homo sapiens [ORGN] Search Save search Limits Advanced Help

Display Settings: Summary, 20 per page, Sorted by Default order Send to: Filter your results:

Found 20842423 nucleotide sequences. Nucleotide (10408382) EST (8704845) GSS (1729196)

Results: 1 to 20 of 10408382 << First < Prev Page 1 of 520420 Next > Last >>

[Homo sapiens uncharacterized LOC101927285 \(AC007131.1\), long non-coding RNA](#)

1. 2,036 bp linear ncRNA, lncRNA
Accession: NR_110219.1 GI: 572153117
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927482 \(AC020571.3\), transcript variant 2, long non-coding RNA](#)

2. 1,413 bp linear ncRNA, lncRNA
Accession: NR_110226.1 GI: 572153116
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927217 \(AC097500.2\), long non-coding RNA](#)

3. 605 bp linear ncRNA, lncRNA
Accession: NR_110218.1 GI: 572153115
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927196 \(AC007966.1\), transcript variant 3, long non-coding RNA](#)

4. 1,196 bp linear ncRNA, lncRNA
Accession: NR_110216.1 GI: 572153114
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927196 \(AC007966.1\), transcript variant 4, long non-coding RNA](#)

5. 1,069 bp linear ncRNA, lncRNA
Accession: NR_110217.1 GI: 572153113
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927142 \(AC016738.3\), long non-coding RNA](#)

6. 771 bp linear ncRNA, lncRNA
Accession: NR_110213.1 GI: 572153112
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927554 \(AC019118.2\), long non-coding RNA](#)

7. 2,344 bp linear ncRNA, lncRNA
Accession: NR_110228.1 GI: 572153111
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927289 \(AC131097.3\), long non-coding RNA](#)

8. 576 bp linear ncRNA, lncRNA
Accession: NR_110220.1 GI: 572153110
[GenBank](#) [FASTA](#) [Graphics](#)

[Homo sapiens uncharacterized LOC101927196 \(AC007966.1\), transcript variant 1, long non-coding RNA](#)

Filter your results: All (10408382) Bacteria (0) INSDC (GenBank) (10124092) mRNA (374596) RefSeq (112968) Manage Filters

Top Organisms [Tree] Homo sapiens (10369653) synthetic construct (38398) Mus musculus (157) Human papillomavirus type 16 (78) Homo sapiens neanderthalensis (27) All other taxa (99) More...

Find related data Database: Select Find items

Search details "Homo sapiens"[Organism] Search See more...

Recent activity Turn Off Clear homo sapiens[ORGN] (10408382) Nucleotide Entrez Programming Utilities Help Bookshelf Entrez Help - Entrez Help Bookshelf See more...

- Introduction
- Banques de données de séquences nucléiques
- Banques de données de séquences protéiques
- Banques de données spécialisées
- Interrogation des banques de données
- **Formats de fichiers en bioinformatique**
- Références

1. Les fichiers des banques de séquences (*flatfiles*)
→ Ceux qu'on a vu mais aussi ceux d'autres banques
2. Le format **FASTA** : une 1re 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 TCP1-beta gene, partial cds; and Axl2p (AXL2) and Rev7p (REV7) genes, complet  
gatcctccatatacaacggtatctccacctcaggtttagatctcaacaacggaaccattg  
ccgacatgagacagttaggtatcgctcgagagttacaagctaaaacgagcagtagtcagct  
ctgcatctgaagccgctgaagttctactaagggtggataaacatcatccgtgcaagaccaa  
gaaccgccaatagacaacatatgtaacatatttaggatatacctcgaaaataataaacg  
ccacactgtcattattataattagaaacagaacgcaaaaattatccactatataattcaa  
agacgcgaaaaaaaaagaacaacgcgtcatagaacttttggcaattcgcgtcacaaataa  
atthttggcaacttatgthttcctcttcgagcagtactcgagccctgtctcaagaatgtaat  
aataccatcgtaggtatggttaaagatagcatctccacaacctcaaagctccttgccga
```

3. Les formats des outils phylogénétiques

→ Le format **NEXUS** [Maddison *et al.*, 97] (utilisé par PAUP)

```
#nexus

BEGIN Taxa;
DIMENSIONS ntax=6;
TAXLABELS
[1] 'Europe'
[2] 'Albania'
[3] 'Andorra'
[4] 'Belarus'
[5] 'Belgium'
[6] 'BosniaHerzeg'
;
END; [Taxa]
BEGIN Distances;
DIMENSIONS ntax=6;
FORMAT labels=left diagonal triangle=both;
MATRIX
[1] 'Europe'          0 64 48 37 41 57
[2] 'Albania'        64 0 73 81 51 75
[3] 'Andorra'        48 73 0 50 70 68
[4] 'Belarus'        37 81 50 0 46 68
[5] 'Belgium'        41 51 70 46 0 70
[6] 'BosniaHerzeg'  57 75 68 68 70 0
;
END; [Distances]
```

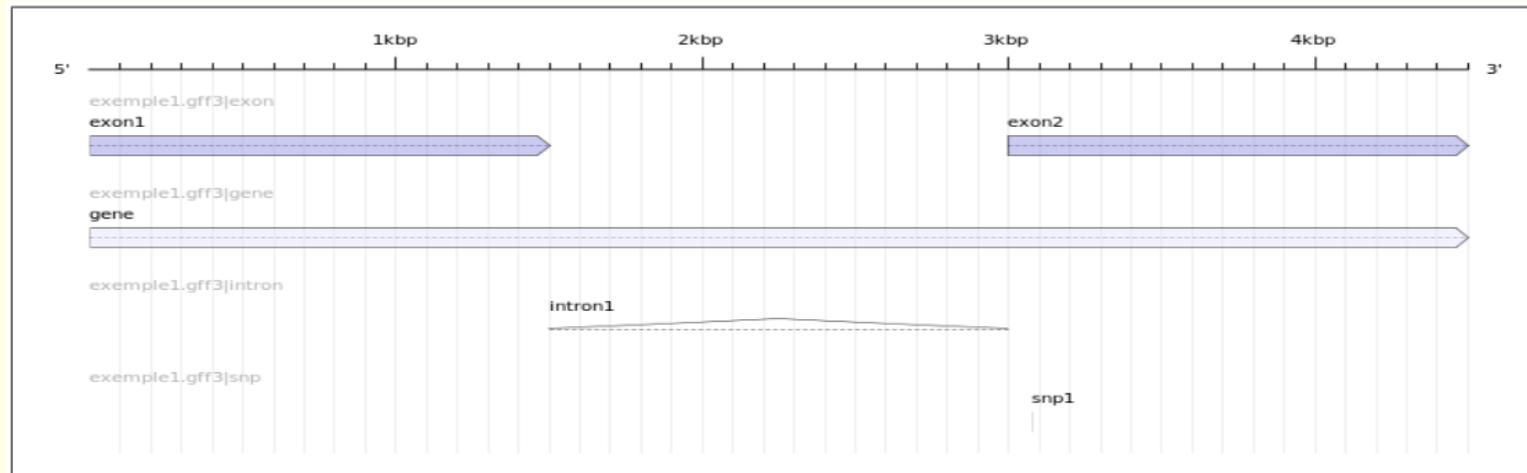
→ Le format **Phylip** (suite EMBOSS) : la 1re ligne contient le nb de séquences suivi de leur taille (les séq. ayant été alignées, elles ont nécessairement toutes la même taille). Viennent ensuite les informations concernant chaque espèce ou gène

```
7          50
C1          GCCAACCCCA CGGTCACTCT GTTCCCGCCC TCCTGGAGCT CCAAGACAAG
C2          GCTGCCCCCT CGGTCACTCT GTTCCCGCCC TCCTGGAGCT TCAAGACAAG
C3          GCTGCCCCCT CGGTCACTCT GTTCCCACCC TCCTGGAGCT TCAAGACAAG
C4          GAGACACCTT CATCTCCTCT GACCCAGAG GCAGGGAGCT CCAAGACAAG
C5          GCCACCCCTT TGGTCACTCT GTTC---CCC TCCTGGAGCT CCAAGACAAG
C6          GCTGCCCCAT CGGTCACTCT GTTCCCGCCC TCCTGGAGCT TCAAGACAAG
C7          GCTGCCCCCT CGGTCACTCT GTTCCCACCC TCCTGGAGCT TCAAGACAAG
```

4. Les formats d'échange

- **GFF** (*Generic Feature Format*) : exploiter la richesse informationnelle des annotations au travers d'une **structure tabulaire à 9 colonnes**
- Rendre compte de la présence de **sous-régions fonctionnelles imbriquées**
- Outils de **validation**

```
##gff-version 3
##gff-version 3
ctg123 GenBank exon 1 1500 . + . ID=exon1;note=chromobox homolog 8
ctg123 GenBank gene 1 4500 . + . ID=gene;Dbxref=GeneID:779897
ctg123 GenBank exon 3000 4500 . + . ID=exon2
ctg123 GenBank intron 1501 2999 . + . ID=intron1
ctg123 GenBank snp 3080 3080 . + . ID=snp1
```



→ **ASN.1** (*Abstract Syntax Notation number one*) : norme qui définit un formalisme de description de types de données abstraits

Standard international ; format semi-structuré ; format de base pour les données du NCBI

```
Seq-entry ::= set {
  level 1 ,
  class nuc-prot ,
  descr {
    title "Mus musculus Brca1 mRNA, and translated products" ,
    source {
      org {
        taxname "Mus musculus" ,
        db {
          {
            db "taxon" ,
            tag
              id 10090 } } ,
        orgname {
          name
            binomial {
              genus "Mus" ,
              species "musculus" } , ...
```

→ XML (eXtensible Markup Language) : fournit un moyen pour implémenter des ontologies (vocabulaire structuré)

Langage à balise comme HTML ; standard international ; format semi-structuré

```
<?xml version="1.0"?>
<!DOCTYPE GBSeq PUBLIC "-//NCBI//NCBI GBSeq/EN" http://www.ncbi.nlm.nih.gov/dtd/NCBI_GBSeq.dtd">
<GBSet>
<GBSeq>
  <GBSeq_locus>MMU35641</GBSeq_locus>
  <GBSeq_length>5538</GBSeq_length>
  <GBSeq_strandedness value="not-set">0</GBSeq_strandedness>
  <GBSeq_moltype value="mrna">5</GBSeq_moltype>
  <GBSeq_topology value="linear">1</GBSeq_topology>
  <GBSeq_division>ROD</GBSeq_division>
  <GBSeq_update-date>18-OCT-1996</GBSeq_update-date>
  <GBSeq_create-date>25-OCT-1995</GBSeq_create-date>
  <GBSeq_definition>Mus musculus Brca1 mRNA, complete cds</GBSeq_definition>
  <GBSeq_primary-accession>U35641</GBSeq_primary-accession>
  <GBSeq_accession-version>U35641.1</GBSeq_accession-version>
  ...
```

5. Formats correspondants aux alignements multiples

→ **MSF** *Multiple Sequence Format*

→ **MAF** *Multiple Alignment Format*

→ **ALN** (CLUSTALW)

→ ...

6. ...

- Cette liste de formats classiques en bioinformatique n'est bien sûr pas exhaustive
- Pour certains formats, représentation entrelacée (interleaved) ou séquentielle (sequential) possible

Certains formats tendent à devenir des **standards** ou sont **plus génériques** que d'autres, il faut encourager leur utilisation

- Les logiciels peuvent n'accepter qu'un type de format particulier
- Des outils existent pour transformer un format en un autre :
 - **ReadSeq** (<http://iubio.bio.indiana.edu/cgi-bin/readseq.cgi>)
Formats supportés : IG/Stanford, GenBank/GB, NBRF, EMBL, GCG, DNASTrider, Fitch, Pearson/Fasta, Phylip3.2, Phylip, PIR/CODATA, MSF, ASN.1 et PAUP/NEXUS
 - **Outils de conversion de l'EBI** (<http://www.ebi.ac.uk/Tools/sfc/>)
 - **Fonctions BioPython, BioPerl, ...**
 - ...
- Des outils de visualisation de données existent et acceptent, en général, plusieurs formats (AnnotationSketch, Seaview, MEGA, ...)

-
- Introduction
 - Banques de données de séquences nucléiques
 - Banques de données de séquences protéiques
 - Banques de données spécialisées
 - Interrogation des banques de données
 - Formats de fichiers en bioinformatique
 - **Références**

- Chap 1-4 dans "Bioinformatics : A Practical Guide to the Analysis of Genes and Proteins, 3rd Edition"
[Baxevanis & Ouellette,2005]
- "Bioinformatique - Cours et cas pratique"
[Deléage Gouy, 2013]
- Cours en ligne de Maude Pupin et Jean-Stéphane Varré
(université de Lille)

