

Bases de Données en Génomique



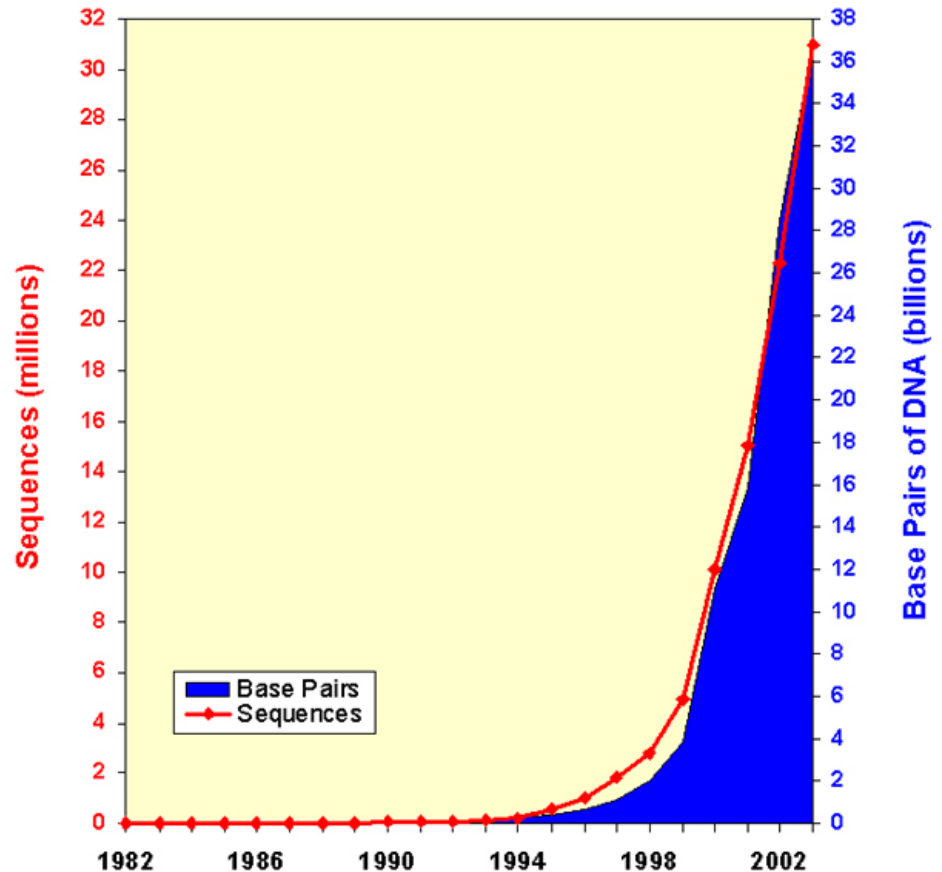
Daniel Gautheret, 2004
ESIL, Université de la Méditerranée

Genbank: La banque d'ADN du NIH

Etat au 2-2004

- ★ 38×10^9 bases
- ★ 32×10^6 séquences
- ★ Genbank double environ tous les 14 mois depuis ses débuts en 1982.
- ★ Nouvelle version tous les 2 mois

Growth of GenBank



Divisions de Genbank

- ★ **ESTs** (Expressed sequence tags):
Principale division ed Genbank. 18
10⁶ sequences, 580 organismes
différents
- ★ **GSS** (Genome Sequence Survey):
résultats de séquençages aléatoire
de BAC, dans le cadre de projets
Génome
- ★ **HTGS** (High Throughput Genomic
Sequences): séquences génomiques
en cours d'assemblage. Une fois
assemblées, les séquences passent
dans les divisions « organisme ».
- ★ Bactéries (**BCT**), virus (**VRL**),
primates (**PRI**), rongeurs (**ROD**) etc:
divisions « organismes ».
- ★ 17 divisions en tout.

Nb entrées	Nb. bases	Espèce
1355113	854232260	Homo sapiens
378892	179249409	Mus musculus
76471	139699685	Caenorhabditis elegans
66177	69663817	Arabidopsis thaliana
48963	53428355	Drosophila melanogaster
10571	28658828	Saccharomyces cerevisiae
39568	25816686	Rattus norvegicus
4923	17859484	Escherichia coli
32221	16490243	Fugu rubripes
31480	13072925	Oryza sativa
28406	11746328	Rattus sp.
9540	10912762	Schizosaccharomyces pombe
24125	10712174	Human immunodeficiency virus type 1
1086	9893044	Bacillus subtilis
15370	5794059	Brugia malayi
661	5701954	Mycobacterium tuberculosis
4852	5585160	Gallus gallus
4680	5400457	Plasmodium falciparum
5063	4559072	Bos taurus
10845	4409926	Toxoplasma gondii

Organismes dans Genbank (en 2002)

Enregistrement Genbank

- ★ Chaque enregistrement se voit attribuer un numéro d'accession, stable et unique, et chaque séquence un numéro GI.
- ★ Quand un changement est effectué dans un enregistrement Genbank, le num. d'accession reste, le GI change.

Enregistrement Genbank avec annotation

```
LOCUS       L10986               47233 bp    DNA     linear   INV 21-SEP-2004
DEFINITION  Caenorhabditis elegans cosmid F10E9, complete sequence.
ACCESSION   L10986
VERSION     L10986.2  GI:38638818
KEYWORDS    HTG.
SOURCE      Caenorhabditis elegans
  ORGANISM  Caenorhabditis elegans
            Eukaryota; Metazoa; Nematoda; Chromadorea; Rhabditida;
            Rhabditoidea; Rhabditidae; Peloderinae; Caenorhabditis.
REFERENCE   1  (bases 1 to 47233)
AUTHORS     .
CONSRTM     WormBase Consortium
TITLE       Genome sequence of the nematode C. elegans: a platform for
            investigating biology. The C. elegans Sequencing Consortium
JOURNAL     Science 282 (5396), 2012-2018 (1998)
MEDLINE     99069613
PUBMED      9851916
FEATURES             Location/Qualifiers
     source            1..47233
                     /organism="Caenorhabditis elegans"
                     /mol_type="genomic DNA"
                     /strain="Bristol N2"
                     /db_xref="taxon:6239"
                     /chromosome="III"
                     /clone="F10E9"
     gene              265..26728
                     /gene="mig-10"
                     /locus_tag="F10E9.6"
     CDS               join(265..338,3266..3515,15194..15317,21507..21
                     21727..21887,23171..23335,24302..24472,24524..24608,
                     25012..25827,26284..26430,26478..26728)
                     /gene="mig-10"
     /translation="MDSCEEECDLEVDSDEEDQLFGEKICISLLSSLLPLSSSTLLSNA
            INLELDEVERPPPLLNVLVEEQQFPKVCANIEEENELEADTEEDIAETADDEESKDPVE
            KTFNFEPVSTMDTYDFDPYPVQIRARPVQPPKPPIDTVRYSMNNIKESADWQLDELL
            EELEALETQLNSSNGGDQLLLGVSGIPASSSRENVKSISTLPPPPALS YHQTPQQPQ
            . . .
            QVYTIGIWEEKYKSPTPWCISIKLTALQMKRSQFIKYICAEDEMTFKKWLVALRIAKN
            GAELLENYERACQIRRETLGPASSMSAASSSTAISEVPHLSLHHQRTPSVASSIQLSS
            HMMNNPHTPLSVNVRNQSPASFVNSCQQSHPSRTSAKLEIQYDEQPTGTIKRAPLDV
            LRRVSRASTSSPTIPQEESSDSEDFPAPPPVASVMRMPVTPPKPCTPLTSKKAPPP
            PPKRSDTTKLQSASPMAPAKNDLEAALARRREKMATMEC"
     . . .
```

```
BASE COUNT      2598 a    2024 c    1888 g    2449 t
ORIGIN
1  ttctaaaagt cgaaaaacga gcaatTTTTg atgctagatt ttttgattg acgaatTTTT
61 tcagttTTTT ttcttTaaaa aaggTTTTtg acccctTaaa gttttccttt cccttccaat
121 tttttccttc ttctttatac gacttctcaa gtttcaactc taaaacaag ctacatgtac
181 atttccggta aactttgtgt ctcagaagat ccattttctt tttgttacct ttattcaaga
241 ttgaattcca aaatttcagc caatatggac agttgcgaag aggaatgcga tctggaagtt
301 gacagtgcgc aagaagatca acctttttgt gaaaagtgtg gagttcttat tgtggtaacc
361 aaagaaatgt cagtggtcgc taaacacttg actcccaaat ggtttctcgt aattacctta
421 tgcacacttt tcaagtgttt gccgtttgat cttagccaat ttgaacgctt tagatgttaa
481 atggaaaatg ggtaaaagtt tttattttat agaaaaaagg tttggaaaaa aatcgagtca
541 ctgaatagtt tgaagaacgc aaaaaataaa ctttccaaaa atcataaaac atttagtgtt
601 tcgaaaatta tagtgTTTTt tttgttggtg tgttttgaca aaagctaaac catctttatt
661 gtagttttgt aaaatgttca caaagatgcg tttttttttc aaatttggca ggctatcttt
721 acattcacat ttggataaatt caaatTTTTc ttatcgctaa caaatTTTcc tatttttcca
781 attattcgtt ttataaaagc tttggtagta tgttgtgtct atcttttagt gtcacagtt
. . .
//
```

Les banques de gènes

Nom	Type	Organisme	Description	Nb. Enrgst.
Refseq	gene+ mRNA+ prot	H, M, R	Itération à partir de Blast seed mRNA/EST vs contigs de Genbank. +Annotation manuelle: publications, UTR prolongées, etc.	11 405(H) 5 749 (M)
HGI (Human Gene Index)	mRNA	H, M, R, D, Arabid., Fugu, riz, etc.	Genbank mRNA+EST contigués. Toutes les solutions alternatives sont conservées.	388 000
Ensembl	gènes + transcrits	H + eucaryotes?	Genscan sur contig, puis Blast vs prot, mRNA, EST, PFAM	35 500 genes, 44 860 transcrits
Unigene	clusters	H, M, R, bovin, zebrafish, blé, riz, mais, orge	Clustering itératif à partir de mRNA+CDS génomique+ESTs. Pas de contigage.	89 371 cl

★ [Unigene](#): banque d'ESTs classifiés ("clusterisés"). Dans chaque cluster Unigene sont regroupés des EST ayant une similarité de séquence significative. On peut donc trouver des transcripts différents et des artefacts (chimères, etc.). Unigene ne propose pas de mRNA reconstruits (contigs) à partir des séquences d'un cluster.

★ TIGR Human Gene Index ([HGI](#)). Ici encore on a clusterisé les EST, mais HGI est une banque de "contigs", c.a.d. de séquences de mRNA reconstruites à partir des EST d'un même cluster. Les clusters étant souvent hétérogènes, ils produisent souvent plusieurs contigs. Ces contigs doivent théoriquement correspondre à des mRNA alternatifs.

Autres banques de séquences

Nucléotidiques

- ★ gbEST / dbEST: Division EST dans Genbank
- ★ EMBL: Equivalent européen de Genbank. Format différent, contenu presque identique.
- ★ Banques spécialisées Certaines collections de séquences, bien que généralement présentes dans Genbank, sont beaucoup plus utiles lorsqu'elles sont rassemblées dans des banques spécialisées, par ex:
 - Récepteurs des lymphocytes T (Réarrangements de l'ADN)
 - Génomes HIV, etc.
- ★ NR nucléique (Non-redundant). Banque combinée: Genbank+refseq (20x10e9 nt / oct. 2002)

Ensembl (www.ensembl.org)

★ Plusieurs banques en une:

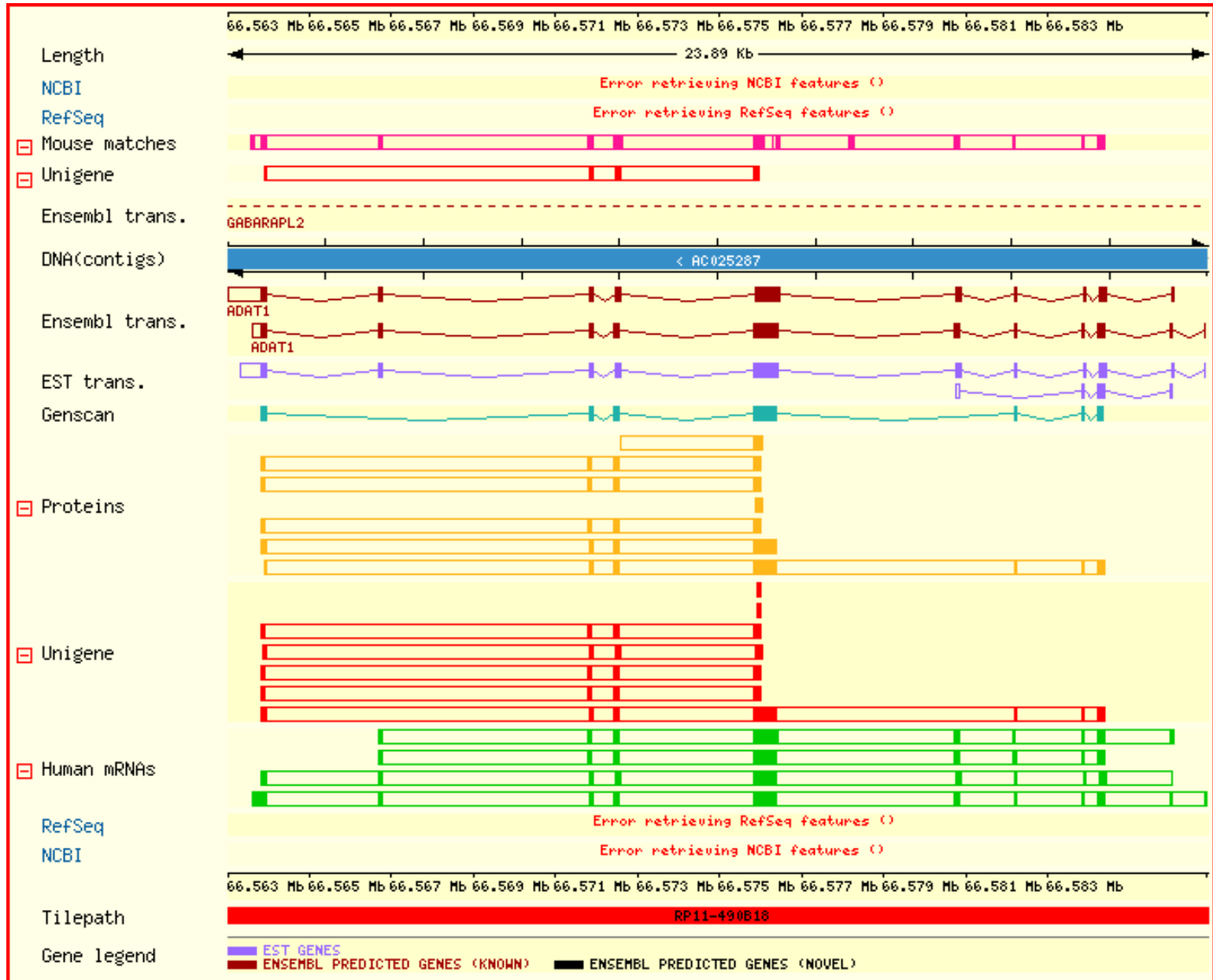
- Peptides confirmés
- Transcrits confirmés
- peptides prédits
- Transcrits prédits
- Génome assemblé (golden path)

★ Méthode de prédiction (système Genewise): Genscan sur contig, puis Blast contre: protéines, mRNA, EST

★ Version Juillet 2001, humain: Confirmed genes: 21921; Predicted genes: 24636; Confirmed exons: 143479; Predicted exons: 770562; Transcripts: 23931; Contigs: 329154; Sequences: 29080; base pairs: 4318661441.

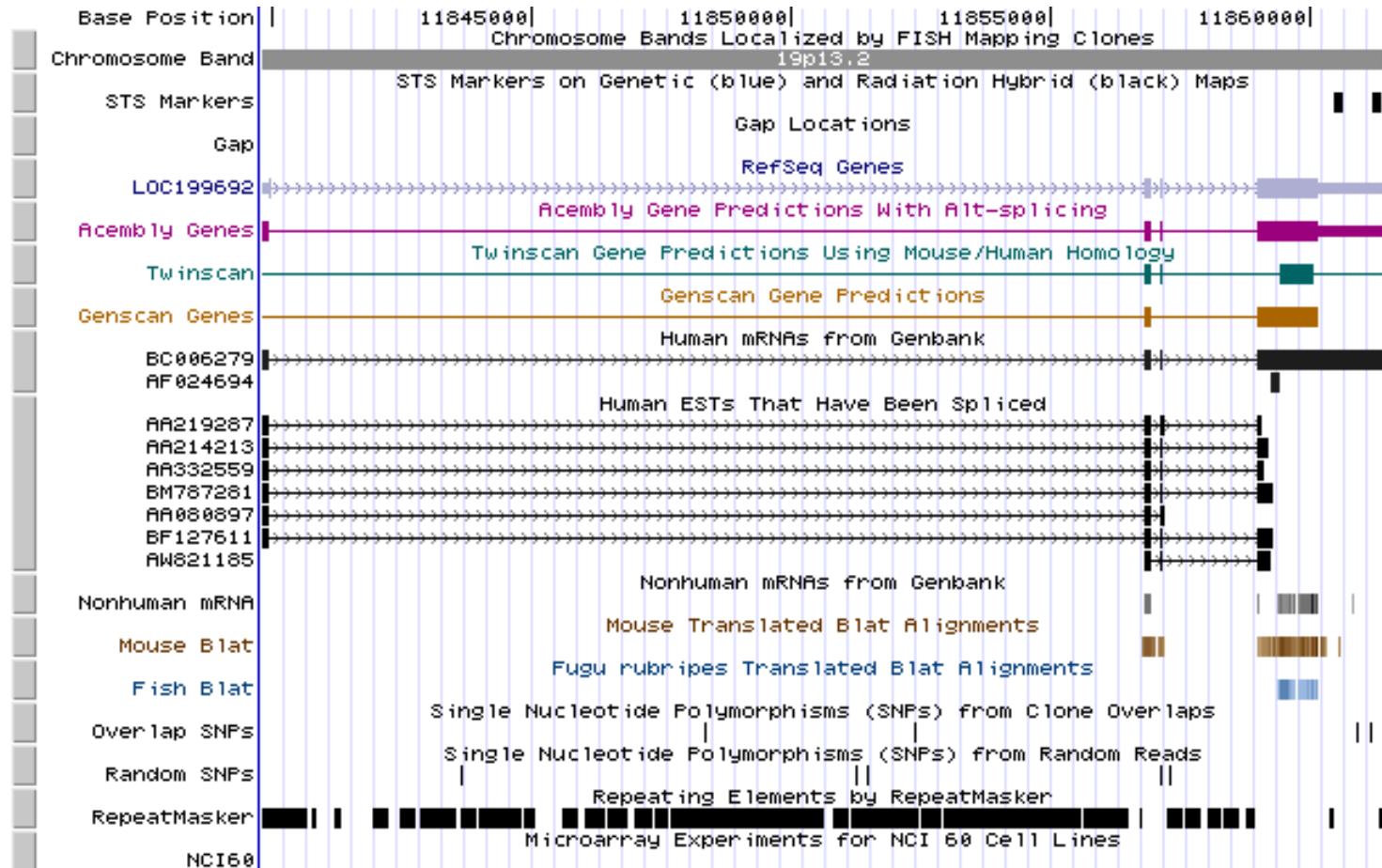
Species - Ensembl v24		
Human	<i>pre!</i>	NCBI 34 Jul 04
Mouse		NCBI m33 Jul 04
Zebrafish		WTSI Zv4 Sep 04
Rat		RGSC 3.1 Jul 04
Chicken		WASHUC1 Jul 04
Mosquito		MOZ 2 Apr 04
Fugu		Fugu v2.0 May 04
Fruitfly		BDGP 3.1 Jul 03
Chimp		CHIMP1 May 04
Honeybee		Amel1.1 Sep 04
Tetraodon		TETRAODON7 Sep 04
Dog	<i>pre!</i>	BROADD1
<i>C. elegans</i>		WS 116 Apr 04
<i>C. briggsae</i>		cb25_aqp8 Jul 03

Ensembl: « contig view »



Banque génomique UCSC

★ <http://www.genome.ucsc.edu/>



Banques protéiques

- ★ Swissprot. La mieux annotée des banques protéiques.
Release 39 (2001): 101247 entrées, 37 135 523 aa.
Attention: toutes les protéines connues n'y sont pas! Visiter le serveur
- ★ PIR (Protein Identification Resource), EMBL.
- ★ NR Protéique (Non-redundent): Banque protéique du NCBI
= Traduction de tous les CDS de GenBank + PDB + SwissProt + PIR + PRF - redondances.
- ★ Banques spécialisées
 - Cazy (Carbohydrate Active Enzymes)
 - Etc.

SRS

Sequence Retrieval System

Database selection page

The screenshot shows the SRS@EMBL-EBI website interface in a Netscape browser window. The browser's address bar displays the URL: <http://srs.ebi.ac.uk/srsbin/cgi-bin/wgetz?page+top+id+ufno10a1eD>. The website features a navigation menu with options: Quick Search, Library Page, Query Form, Tools, Results, Projects, Views, Databanks, and HELP. A search bar with a "Quick Search" button is located at the top. Below the search bar, there are sections for "Search Options" and "Available Databanks".

Search Options:

1. Select the **databanks** you want to search
2. Enter your **search terms** in the **Quick Search** box, or choose a **query form** from below

Buttons: [Standard Query Form](#), [Extended Query Form](#), [Browse Entries](#)

Tips:

- ▶ bookmark this [link](#) to return to your project
- ▶ [Linking to SRS?](#)
- Please read this [document](#) for important information regarding linking to our SRS server.

BookMarkLets

Available Databanks:

Expand all Collapse all Show databanks tooltips:

- Literature, Bibliography and Reference Databases**
 - MEDLINE OMIM TAXONOMY GENETICCODE
 - Patent Abstracts KarynsGenomes
 - Literature, Bibliography and Reference Databases - subsections*
 - MEDLINE (Main Release) MEDLINE (Updates) OLDMEDLINE MED2PUB
- Nucleotide sequence databases**
 - EMBL EMBL (Contig) EMBL (Contigs expanded) Genome Reviews
 - IMGTHLA IMG/LIGM-DB LiveLists PATENT_DNA
 - RefSeq Genome DB
 - Nucleotide sequence databases - subsections*
 - EMBL (Release) EMBL (Updates) EMBL (Third Party Annotation) EMBL (Coding Sequences)
 - RefSeq Genome Release RefSeq Genome Updates
- UniProt Universal Protein Resource**
 - UniProt UniParc UniRef100 UniRef90 UniRef50
 - UniProt/Swiss-Prot UniProt/TrEMBL
- Other protein sequence databases**
 - RefSeq Proteome DB IPI EPO Proteins JPO Proteins
 - USPTO Proteins MHCBN BCIPEP SWISSCHANGE
 - Protein sequence databases - subsections*
 - Refseq Proteome Release Refseq Proteome Updates
- Deprecated Protein Databases**
- Nucleotide related databases**
- Protein function databases**
- Protein structure databases**
- Protein expression and metabolite databases**

SRS

Query:

Standard Query Form - Netscape

File Edit View Go Bookmarks Tools Window Help

http://srs.embl.ac.uk/srsbin/cgi-bin/wgetz

Home Local Institutions Journaux Mot/Annu Cours/Guides MolBio 1 RNA FP6-ATD trad

SRS@EMBL-EBI Quick Search Library Page Query Form Tools Results Projects Views Database

Reset search UniProt/Swiss-Prot

Search Options

Combine search terms with:

Use wildcards

Get results of type:

Result Display Options

View results using:

or

Create a view

Show results per page

Fields you can search

Your search terms

In a single field, you can separate multiple values by &, |, !

<input type="checkbox"/> AllText	<input type="text" value="argonaute"/>
<input type="checkbox"/> Organism Name	<input type="text" value="drosophila"/>
<input type="checkbox"/> AllText	<input type="text"/>
<input type="checkbox"/> AllText	<input type="text"/>

Create a view

Select the fields you want displayed in your view and choose the format

Choose 1 or more fields:

- ID
- EntryName
- AccessionNumber
- Creation Date
- Seq Mod Date
- Annot Mod Date
- Description

Display As: Table List

Sequence Format:

Entrez

The screenshot shows the Entrez cross-database search page in a Netscape browser window. The browser title is "Entrez cross-database search - Netscape" and the address bar shows "http://www.ncbi.nlm.nih.gov/Entrez/". The page features the NCBI logo and the Entrez logo with the tagline "Entrez, The Life Sciences Search Engine". A navigation bar includes links for HOME, SEARCH, SITE MAP, PubMed, Entrez, Human Genome, GenBank, Map Viewer, and BLAST. A search bar is present with "GO" and "CLEAR" buttons and a "Help" link. The main content area is titled "Welcome to the new Entrez cross-database search page" and lists various databases in a grid format, each with an icon and a help link.

NCBI **Entrez, The Life Sciences Search Engine**

HOME SEARCH SITE MAP PubMed Entrez Human Genome GenBank Map Viewer BLAST

Search across databases Help

Welcome to the new Entrez cross-database search page

PubMed: biomedical literature citations and abstracts	Books: online books
PubMed Central: free, full text journal articles	OMIM: online Mendelian Inheritance in Man
Nucleotide: sequence database (GenBank)	UniGene: gene-oriented clusters of transcript sequences
Protein: sequence database	CDD: conserved protein domain database
Genome: whole genome sequences	3D Domains: domains from Entrez Structure
Structure: three-dimensional macromolecular structures	UniSTS: markers and mapping data
Taxonomy: organisms in GenBank	PopSet: population study data sets
SNP: single nucleotide polymorphism	GEO Profiles: expression and molecular abundance profiles
Gene: gene-centered information	GEO DataSets: experimental sets of GEO data
HomoloGene: eukaryotic homology groups	Cancer Chromosomes: cytogenetic databases
PubChem Compound: small molecule chemical structures	PubChem BioAssay: bioactivity screens of chemical substances
PubChem Substance: chemical substances screened for bioactivity	
Journals: detailed information about the journals indexed in PubMed and other Entrez databases	MeSH: detailed information about NLM's controlled vocabulary
NLM Catalog: catalog of books, journals, and audiovisuals in the NLM collections	

Glossaire de génomique...

[http://www.sciencemag.org/cgi/content/full/291/
5507/1197](http://www.sciencemag.org/cgi/content/full/291/5507/1197)

(*Science*, Vol 291, 1197)