

# BCB 5200 Introduction to Bioinformatics

## **Lecture 03: NCBI Databases**

Bioinformatics and Computational Biology  
Saint Louis University

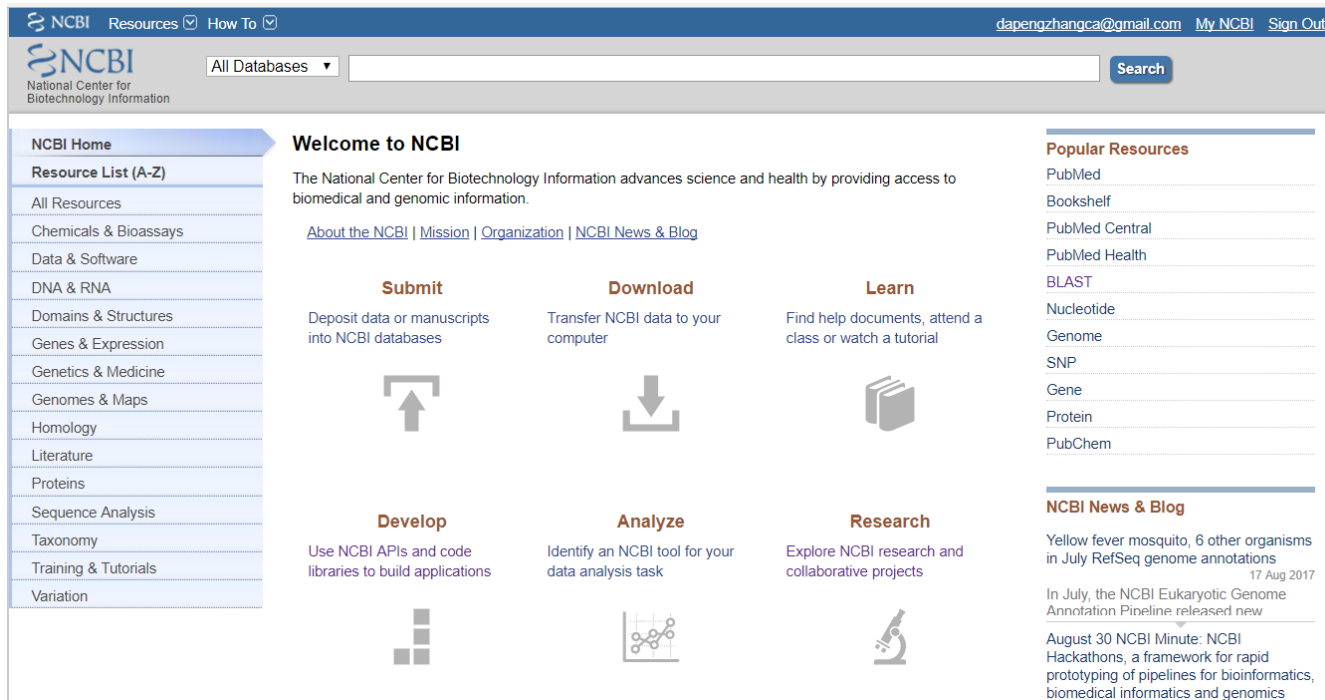
1. How to quickly identify the gene expression patterns of a particular gene?

Maelstrom

C10orf99

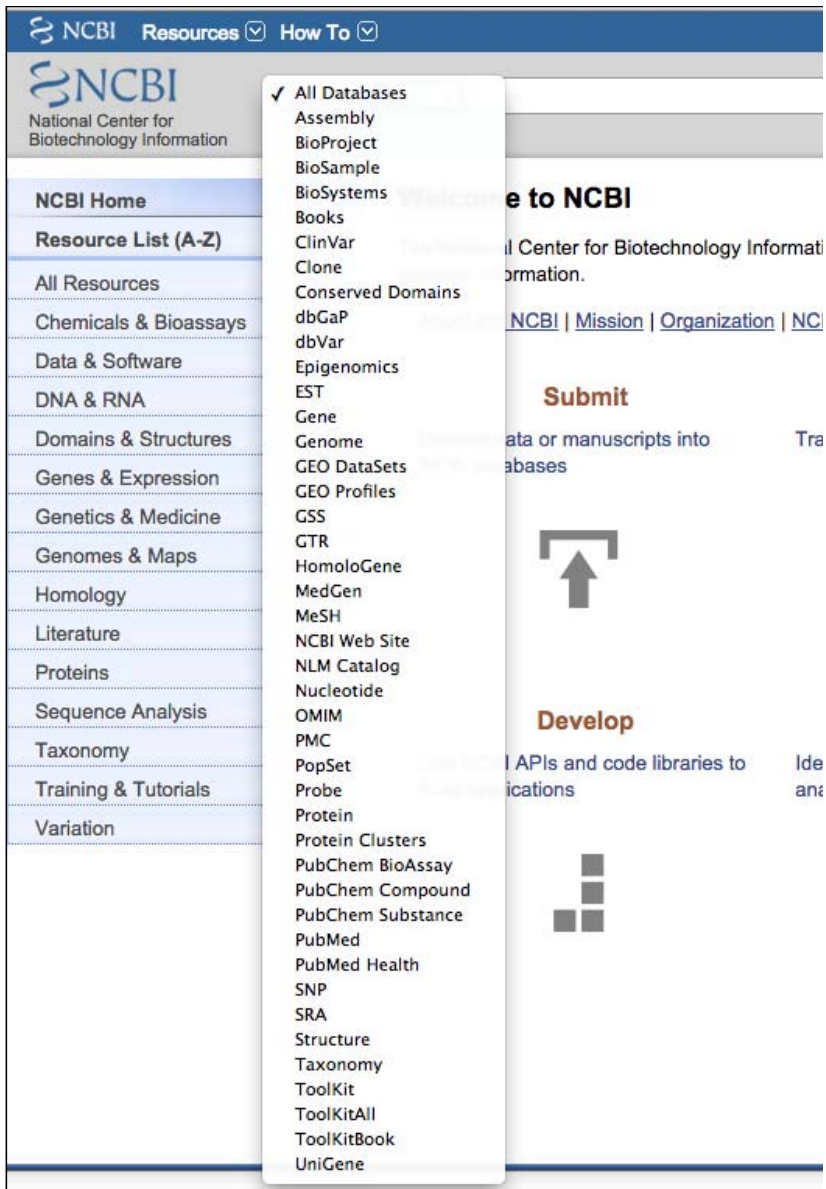
Variations of the data

# NCBI: <https://www.ncbi.nlm.nih.gov/>



The screenshot shows the NCBI homepage with a blue header bar containing the NCBI logo, navigation links (Resources, How To), a user email (dapengzhangca@gmail.com), and links to My NCBI and Sign Out. Below the header is a search bar with a dropdown menu set to 'All Databases' and a 'Search' button. The main content area is divided into several sections. On the left is a 'Resource List (A-Z)' sidebar with links to various databases and tools. The central area features a 'Welcome to NCBI' message, a brief description of the center's mission, and links to 'About the NCBI', 'Mission', 'Organization', and 'NCBI News & Blog'. Below this are six large icons representing different functions: Submit (Deposit data or manuscripts into NCBI databases), Download (Transfer NCBI data to your computer), Learn (Find help documents, attend a class or watch a tutorial), Develop (Use NCBI APIs and code libraries to build applications), Analyze (Identify an NCBI tool for your data analysis task), and Research (Explore NCBI research and collaborative projects). On the right side, there are two sections: 'Popular Resources' listing PubMed, Bookshelf, PubMed Central, PubMed Health, BLAST, Nucleotide, Genome, SNP, Gene, Protein, and PubChem; and 'NCBI News & Blog' featuring a recent article about yellow fever mosquito genome annotations and a link to the August 30 NCBI Minute.

- Develop and maintain molecular and bibliographic databases.
- Develop software for searching, and analysis of these data.
- Provide Web access point for data and software.



## Molecular Databases in NCBI

- Sequences: Nucleotide, Protein
- Gene annotation: Gene, RefSeq
- Genomes: Assembly, Genome
- Expression: GEO, EST
- Protein Domains: CDD
- Homologous Genes
- Genetic Variation: SNP, ClinVar, dbVar
- Taxonomy
- 3D Structures
- Pathways
- Literature: PubMed, PMC
- Small molecules: PubChem



## NCBI Search Services and Tools

- Entrez integrated literature and molecular databases
- BLAST sequence similarity search service
- Graphical Sequence Viewer annotation viewer and analysis tool
- Genome Workbench standalone sequence analysis annotation platform
- SRA Utilities
  - SRA Run Browser: web access for viewing, searching and downloading next generation reads
  - SRA toolkit: standalone SRA manipulator and client

# Entrez is NCBI's primary text search and retrieval system

**Search NCBI databases**

About 1,353,092,513 search results for "all[sb]"

**Literature**

<b>Books</b>	334,634	books and reports
<b>MeSH</b>	253,991	ontology used for PubMed indexing
<b>NLM Catalog</b>	1,511,375	books, journals and more in the NLM Collections
<b>PubMed</b>	24,324,879	scientific & medical abstracts/citations
<b>PubMed Central</b>	3,264,562	full-text journal articles

**Health**

<b>ClinVar</b>	125,404	human variations of clinical significance
<b>dbGaP</b>	166,753	genotype/phenotype interaction studies
<b>GTR</b>	35,977	genetic testing registry
<b>MedGen</b>	260,910	medical genetics literature and links
<b>OMIM</b>	23,795	online mendelian inheritance in man
<b>PubMed Health</b>	51,412	clinical effectiveness, disease and drug reports

**Genomes**

<b>Assembly</b>	33,680	genomic assembly information
<b>BioProject</b>	138,999	biological projects providing data to NCBI
<b>BioSample</b>	2,865,750	descriptions of biological source materials
<b>Clone</b>	37,024,042	genomic and cDNA clones
<b>dbVar</b>	4,305,990	genome structural variation studies
<b>Epigenomics</b>	6,635	epigenomic studies and display tools
<b>Genome</b>	10,777	genome sequencing projects by organism
<b>GSS</b>	37,646,655	genome survey sequences
<b>Nucleotide</b>	155,458,773	DNA and RNA sequences
<b>Probe</b>	31,890,151	sequence-based probes and primers
<b>SNP</b>	444,458,769	short genetic variations
<b>SRA</b>	1,064,488	high-throughput DNA and RNA sequence read archive
<b>Taxonomy</b>	1,310,804	taxonomic classification and nomenclature catalog

**Genes**

<b>EST</b>	75,736,497	expressed sequence tag sequences
<b>Gene</b>	18,350,280	collected information about gene loci
<b>GEO DataSets</b>	1,332,018	functional genomics studies
<b>GEO Profiles</b>	108,708,851	gene expression and molecular abundance profiles
<b>HomoloGene</b>	141,268	homologous gene sets for selected organisms
<b>PopSet</b>	212,126	sequence sets from phylogenetic and population studies
<b>UniGene</b>	6,473,284	clusters of expressed transcripts

**Proteins**

<b>Conserved Domains</b>	49,955	conserved protein domains
<b>Protein</b>	153,454,195	protein sequences
<b>Protein Clusters</b>	820,546	sequence similarity-based protein clusters
<b>Structure</b>	103,644	experimentally-determined biomolecular structures

**Chemicals**

<b>BioSystems</b>	653,443	molecular pathways with links to genes, proteins and chemicals
<b>PubChem BioAssay</b>	1,112,105	bioactivity screening studies
<b>PubChem Compound</b>	62,041,347	chemical information with structures, information and links
<b>PubChem Substance</b>	177,330,151	deposited substance and chemical information

The Entrez system: 39 (and counting) integrated databases

# The Syntax ...

- **Key words search**
- **Boolean operators:** AND, OR, NOT must be entered in **UPPERCASE** (e.g., promoters OR response elements). The default is AND.
- **Parentheses:** Entrez processes all Boolean operators in a **left-to-right** sequence. You can change the order by enclosing individual concepts in **parentheses**. The terms inside the parentheses are processed first.
  - For example, the search statement: g1p3 OR (response AND element AND promoter).
- **Quotation marks:** The term inside the quotation marks is read as one phrase (e.g. “public health” is different than public health, which will also include articles on public latrines and their effect on health workers).
- **Asterisk:** Extends the search to all terms that start with the letters before the asterisk.
  - cilio\* will include such terms as ciliopathy, ciliopathies, and ciliogenesis.

# The easiest Entrez search in Gene database

## Specific gene:

XXX[Symbol] AND YYY[Organism]

`MAEL[symbol] AND human[organism]`

`apt[sym] AND Escherichia coli[orgn]`

## All genes:

YYY[organism]

`zebrafish[orgn]`

`Zea mays[orgn]`



# NCBI-Nucleotide/Protein

- MAEL[All Fields] AND "Homo sapiens"[Organism]

NCBI Resources How To

dapengzhangcai@gmail.com My NCBI Sign Out

Protein Protein MAEL AND human[organism] Search

Create alert Advanced Help

Species  
Animals (23)  
Customize ...

Source databases  
RefSeq (8)  
UniProtKB / Swiss-Prot (7)  
Customize ...

Sequence length  
Custom range...

Molecular weight  
Custom range...

Release date  
Custom range...

Revision date  
Custom range...

Clear all  
Show additional filters

Summary 20 per page Sort by Default order Send to: Filters: Manage Filters

See the [results of this search \(2 items\)](#) in our new [Identical Protein Groups](#) database.

Items: 1 to 20 of 23

<< First < Prev Page 1 of 2 Next > Last >>

- ☐ [MAEL protein \[Homo sapiens\]](#)  
1. 156 aa protein  
Accession: AAH34310.1 GI: 71296887  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [protein maelstrom homolog isoform 2 \[Homo sapiens\]](#)  
2. 403 aa protein  
Accession: NP\_001273306.1 GI: 556503424  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [protein maelstrom homolog isoform 3 \[Homo sapiens\]](#)  
3. 378 aa protein  
Accession: NP\_001273307.1 GI: 556503335  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [protein maelstrom homolog isoform 1 \[Homo sapiens\]](#)  
4. 434 aa protein  
Accession: NP\_116247.1 GI: 14249590  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)
- ☐ [RecName: Full=Protein maelstrom homolog](#)  
5. 434 aa protein  
Accession: Q96JY0.1 GI: 74760900  
[GenPept](#) [Identical Proteins](#) [FASTA](#) [Graphics](#)

Find related data  
Database: Select

Find items

Search details  
MAEL[All Fields] AND "Homo sapiens"[Organism]

Search See more...

Recent activity  
Turn Off Clear

- MAEL AND human[organism] (23) Protein
- protein maelstrom homolog isoform 1 [Homo sapiens] Protein
- MAEL[All Fields] AND human[organism] (23) Protein
- MAEL maelstrom spermatogenic transposon silencer [Homo sapiens] Gene
- MAEL[symbol] AND human[organism] AND (alive[prop]) (1) Gene

# NCBI-Nucleotide/Protein

[https://www.ncbi.nlm.nih.gov/protein/NP\\_116247.1](https://www.ncbi.nlm.nih.gov/protein/NP_116247.1)

## protein maelstrom homolog isoform 1 [Homo sapiens]

NCBI Reference Sequence: NP\_116247.1

[Identical Proteins](#) [FASTA](#) [Graphics](#)

GenPept format

Go to: ☒

LOCUS NP\_116247 434 aa linear PRI 03-JUN-2017  
DEFINITION protein maelstrom homolog isoform 1 [Homo sapiens].  
ACCESSION NP\_116247  
VERSION NP\_116247.1  
DBSOURCE REFSEQ: accession [NM\\_032858.2](#)  
KEYWORDS RefSeq.  
SOURCE Homo sapiens (human)  
ORGANISM [Homo sapiens](#)  
Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;  
Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;  
Catarrhini; Hominidae; Homo.  
REFERENCE 1 (residues 1 to 434)  
AUTHORS Li Q, Wei P, Huang B, Xu Y, Li X, Li Y, Cai S and Li D.  
TITLE MAEL expression links epithelial-mesenchymal transition and stem  
cell properties in colorectal cancer  
JOURNAL Int. J. Cancer 139 (11), 2502-2511 (2016)  
PUBMED [27537253](#)  
REMARK GeneRIF: Study demonstrated that MAEL interacts with Snail and  
inhibit E-cadherin promoter activity. MAEL is an oncogene that  
plays an important role in the development and progression of colon  
cancer.  
REFERENCE 2 (residues 1 to 434)  
AUTHORS Li Q, Wojciechowski R, Simpson CL, Hysi PG, Verhoeven VJ, Ikram MK,  
Hohn R, Vitart V, Hewitt AW, Oexle K, Makela KM, MacGregor S,  
Pirastu M, Fan Q, Cheng CY, St Pourcain B, McMahon G, Kemp JP,  
Northstone K, Rahi JS, Cumberland PM, Martin NG, Sanfilippo PG, Lu  
Y, Wang YX, Hayward C, Polasek O, Campbell H, Bencic G, Wright AF,  
Wedenoja J, Zeller T, Schillert A, Mirshahi A, Lackner K, Yip SP,

FASTA ▾

Fasta format

## protein maelstrom homolog isoform 1 [Homo sapiens]

NCBI Reference Sequence: NP\_116247.1

[GenPept](#) [Identical Proteins](#) [Graphics](#)

>NP\_116247.1 protein maelstrom homolog isoform 1 [Homo sapiens]  
MPNRKASRNAYYFFVQEKIPELRRRGLPVARVADAIPYCSDWALLREEEKEKYAEMAREWRAAQGKDPG  
PSEKQKPVFTPLRRPMLVPKQNVSPDMSALSLKGDQALLGGIFYFLNIFSHGELPPHCEQRFLPCEIG  
CVKYSLQEGIMADFHSEINPGEIPRGFRFHCQAASDSSHKIPISNFERGHNQATVLQNLRYFIHPNPGNW  
PPIYCKSDDRTRVNWCLKHKMAKASEIRQDLQLLTVEDLVVGIYQKFLKEPSKTWIRSLLDVAMWDYSSN  
TRCKWHEENDILFCALAVCKKIAYCISNSLATLFGIQLTEAHVPLQDYEASNSVTPKMMVLDAGRYQKLR  
VGSSGFSHFNSNNEEQRSNTPIGDYPRAKISGQNSSVRGRGITRLLESISNSSNIHKFSNCDTSLSPY  
MSQKDGYSFSSLS

# NCBI-Genome (www.ncbi.nlm.nih.gov/genome/)

[NCBI](#) [Resources](#) [How To](#) dapengzhangca@gmail.com [My NCBI](#) [Sign Out](#)

Genome    [Create alert](#) [Limits](#) [Advanced](#) [Help](#)

**Homo sapiens (human)**  
**Reference genome: Homo sapiens (assembly GRCh38.p11)**  
Download sequences in FASTA format for **genome, transcript, protein**  
Download genome annotation in GFF, GenBank or tabular format  
BLAST against Homo sapiens **genome**  
**All 61 genomes for species:**  
Browse the **list**  
Download sequence and annotation from **RefSeq** or **GenBank**

**NCBI Resources**  
[Genome Data Viewer](#)  
[FTP Human annotation \(GFF\)](#)  
[FTP Human chromosomes](#)  
[Map Viewer](#)

**Tools**  
[BLAST Genome](#)

**Related information**  
[Assembly](#)  
[BioProject](#)  
[Gene](#)  
[Components](#)  
[Protein](#)  
[PubMed](#)  
[Taxonomy](#)

**Search details**  
  
 [See more...](#)

Display Settings: [Overview](#) [Send to:](#)

**Organism Overview** ; [Genome Assembly and Annotation report \[61\]](#) ; [Organelle Annotation Report \[19\]](#)  
**Homo sapiens (human)**  
Human genome projects have generated an unprecedented amount of knowledge about human genetics and health.  
Lineage: [Eukaryota\[2448\]](#); [Metazoa\[818\]](#); [Chordata\[348\]](#); [Craniata\[340\]](#); [Vertebrata\[340\]](#); [Euteleostomi\[335\]](#); [Mammalia\[143\]](#); [Eutheria\[138\]](#); [Euarchontoglires\[63\]](#); [Primates\[28\]](#); [Haplorrhini\[22\]](#); [Catarrhini\[17\]](#); [Hominidae\[5\]](#); [Homo\[1\]](#); [Homo sapiens\[1\]](#)  
Study of the human condition such as genetic and infectious disease, the intersection between genetics and the environment, and population variation is supported by a wealth of genome-scale data. These data sets include: a) numerous sequenced genomes including several which have been assembled; b) studies that examine transcript and protein existence, [More...](#)

**Summary**  
**Sequence data:** genome assemblies: 61; sequence reads: 387 (See [Genome Assembly and Annotation report](#))  
**Statistics:** median total length (Mb): 2996.43  
median protein count: 79257  
median GC%: 40.9  
**NCBI Annotation Release:** 108

**Publications**

1. Long-read sequencing and de novo assembly of a Chinese genome. Shi L, et al. Nat Commun 2016 Jun 30
2. De novo assembly and phasing of a Korean human genome. Seo JS, et al. Nature 2016 Oct 13
3. Interchromosomal core duplicons drive both evolutionary instability and disease susceptibility of the Chromosome 8p23.1 region. Mohajeri K, et al. Genome Res 2016 Nov

[More...](#)

# NCBI-Genome ([www.ncbi.nlm.nih.gov/genome/](http://www.ncbi.nlm.nih.gov/genome/))

 **Representative** (genome information for reference and representative genomes)

**Reference genome:** [\[see all organisms\]](#)

◦  [Homo sapiens GRCh38.p11](#)

Submitter: [Genome Reference Consortium](#)

Loc	Type	Name	RefSeq	INSDC	Size (Mb)	GC%	Protein	rRNA	tRNA	Other RNA	Gene	Pseudogene
Nuc	Chr	22	<a href="#">NC_000022.11</a>	<a href="#">CM000684.2</a>	50.82	47.7	2,493	-	-	965	1,172	348
Nuc	Chr	1	<a href="#">NC_000001.11</a>	<a href="#">CM000663.2</a>	248.96	42.3	11,046	17	90	4,350	5,078	1,372
Nuc	Chr	2	<a href="#">NC_000002.12</a>	<a href="#">CM000664.2</a>	242.19	40.3	8,054	-	8	3,638	3,862	1,166
Nuc	Chr	3	<a href="#">NC_000003.12</a>	<a href="#">CM000665.2</a>	198.3	39.7	6,790	-	4	2,723	2,971	887
Nuc	Chr	4	<a href="#">NC_000004.12</a>	<a href="#">CM000666.2</a>	190.22	38.3	4,374	-	1	2,209	2,441	799
Nuc	Chr	5	<a href="#">NC_000005.10</a>	<a href="#">CM000667.2</a>	181.54	39.5	4,590	-	17	2,225	2,578	766

**Homo sapiens (human)**

**Reference genome:** [Homo sapiens \(assembly GRCh38.p11\)](#)

Download sequences in FASTA format for [genome](#), [transcript](#), [protein](#)

Download genome annotation in [GFF](#), [GenBank](#) or [tabular](#) format

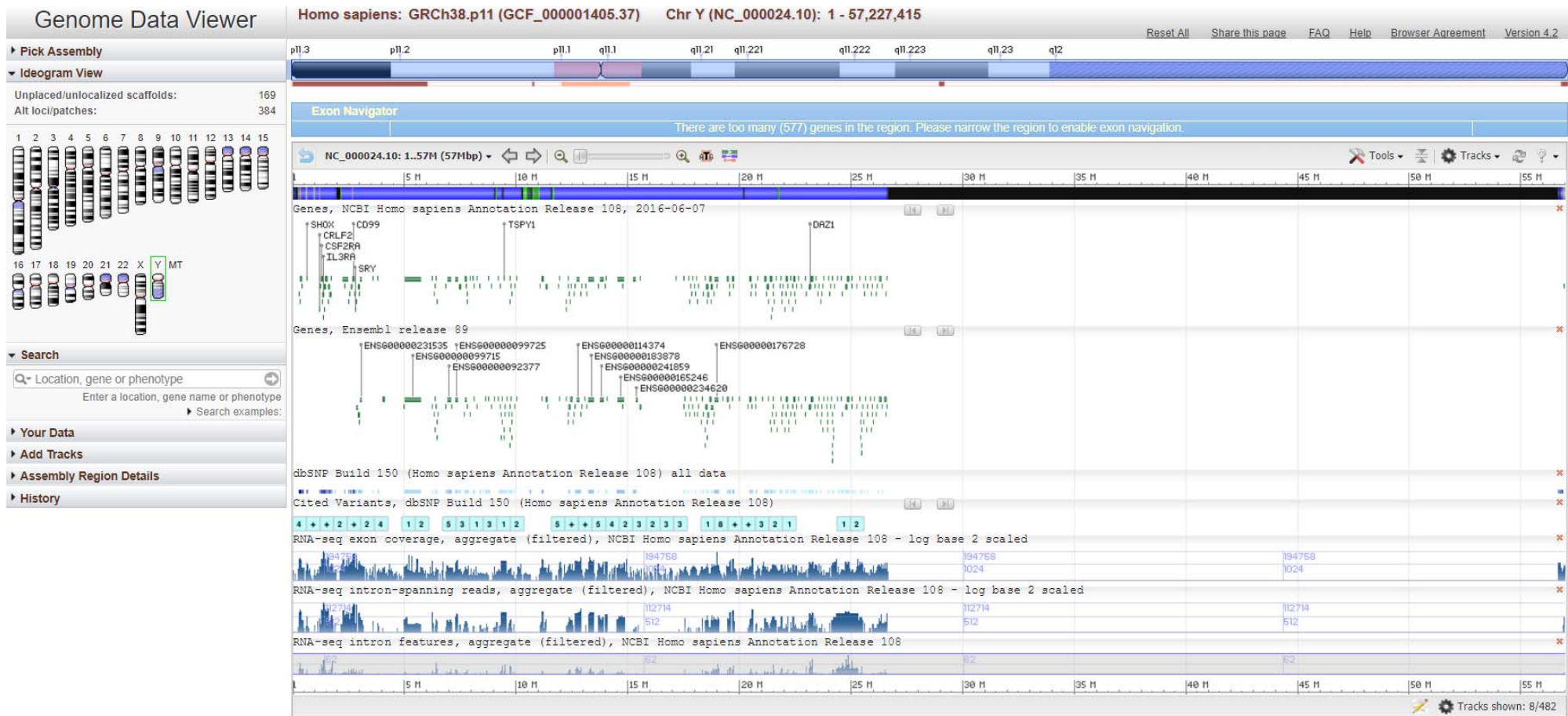
BLAST against [Homo sapiens genome](#)

**All 61 genomes for species:**


Browse the [list](#)

Download sequence and annotation from [RefSeq](#) or [GenBank](#)

# Genome Data Viewer



# NCBI-SRA (sequence read archive)


 NCBI Resources ▾ How To ▾

dapengzhangca@gmail.com My NCBI Sign Out

SRA SRA ▾

Search

Advanced Help



## SRA

Sequence Read Archive (SRA) makes biological sequence data available to the research community to enhance reproducibility and allow for new discoveries by comparing data sets. The SRA stores raw sequencing data and alignment information from high-throughput sequencing platforms, including Roche 454 GS System®, Illumina Genome Analyzer®, Applied Biosystems SOLiD System®, Helicos Heliscope®, Complete Genomics®, and Pacific Biosciences SMRT®.

### Getting Started

- [How to Submit](#)
- [Log in to SRA \(for updating and troubleshooting submissions\)](#)
- [Log in to Submission Portal \(for submitting sequence data\)](#)
- [SRA Documentation](#)
- [Download Guide](#)
- [SRA Fact Sheet \(.pdf\)](#)

### Tools and Software

- [Download SRA Toolkit](#)
- [SRA Toolkit Documentation](#)
- [SRA-BLAST](#)
- [SRA Run Browser](#)
- [SRA Run Selector](#)

### Related Resources

- [Submission Portal](#)
- [Trace Archive](#)
- [dbGaP Home](#)
- [BioProject](#)
- [BioSample](#)



To look for all human gut microbiome data:  
(gut microbiome[All Fields]) AND "Homo sapiens"[Organism]

SRA  (gut microbiome[All Fields]) AND "Homo sapiens"[Organism]  [Create alert](#) [Advanced](#) [Help](#)

Access  
Controlled (470)  
Public (1,403)

Source  
DNA (1,779)  
RNA (36)

Type  
genome (23)

[Clear all](#)  
[Show additional filters](#)

Summary ▾ 20 per page ▾ Send to: ▾ Filters: [Manage Filters](#)

View results as an expanded interactive table using the RunSelector. [Send results to Run selector](#)

**Search results**  
Items: 1 to 20 of 1873

<< First < Prev Page 1 of 94 Next > Last >>

☐ [454 GS FLX Titanium sequencing](#)  
1. 1 LS454 (454 GS FLX Titanium) run: 3,136 spots, 1.3M bases, 2.8Mb downloads  
Accession: ERX458806

☐ [454 GS FLX Titanium sequencing](#)  
2. 1 LS454 (454 GS FLX Titanium) run: 8,572 spots, 3.4M bases, 7.7Mb downloads  
Accession: ERX458801

☐ [454 GS FLX Titanium sequencing](#)  
3. 1 LS454 (454 GS FLX Titanium) run: 4,605 spots, 1.8M bases, 4.1Mb downloads  
Accession: ERX458800

☐ [454 GS FLX Titanium sequencing](#)  
4. 1 LS454 (454 GS FLX Titanium) run: 15,212 spots, 6.1M bases, 13.6Mb downloads  
Accession: ERX458799

**Search in related databases**

Database	Access		all
	public	controlled	
BioSample	<a href="#">652</a>	<a href="#">312</a>	<a href="#">964</a>
BioProject	<a href="#">22</a>	<a href="#">2</a>	<a href="#">24</a>
dbGaP		<a href="#">1</a>	<a href="#">1</a>
GEO Datasets			

**Find related data**

Database:

**Search details**

("gut metagenome"[Organism] OR gut microbiome[All Fields]) AND "Homo sapiens"[Organism]

**RNA/transcriptome**

**metagenome**

# NCBI-Taxonomy



The image shows the NCBI Taxonomy Browser interface. At the top, there is a navigation bar with links to Entrez, PubMed, Nucleotide, Protein, Genome, Structure, and PMC. Below this is a search bar with the text "Search for" followed by a text input field, a dropdown menu set to "as complete name", a checkbox for "lock", and a "Go" button. Below the search bar is a "Display" section with a dropdown menu set to "3" and a text input field for "levels using filter:" followed by a dropdown menu set to "none".

The "Token set" option returns longer names that include the search terms, e.g., hybrid taxa. See what happens if you query "Bos taurus" using the "Complete match" option versus the "Set of tokens" when you are not sure about the exact spelling of a organism name. It tries to find the phonetically closest strings (try "Drozofila" as an example).

**This is the top level of the taxonomy database maintained by NCBI/GenBank. You can explore any of the taxa listed below by clicking it.**

- [Archaea](#)
- [Bacteria](#)
- [Eukaryota](#)
- [Viroids](#)
- [Viruses](#)
- [Other](#)
- [Unclassified](#)

**These are direct links to some of the organisms commonly used in molecular research projects:**

[Arabidopsis thaliana](#)

[Bos taurus](#)

[Caenorhabditis elegans](#)

[Chlamydomonas reinhardtii](#)

[Danio rerio \(zebrafish\)](#)

[Dictyostelium discoideum](#)

[Drosophila melanogaster](#)

[Escherichia coli](#)

[Hepatitis C virus](#)

[Homo sapiens](#)

[Mus musculus](#)

[Mycoplasma pneumoniae](#)

[Oryza sativa](#)

[Plasmodium falciparum](#)

[Pneumocystis carinii](#)

[Rattus norvegicus](#)

[Saccharomyces cerevisiae](#)

[Schizosaccharomyces pombe](#)

[Takifugu rubripes](#)

[Xenopus laevis](#)

[Zea mays](#)



# Homo sapiens

Search for  as  complete name ☐ lock   
Display  3 levels using filter:  none

## Homo sapiens

Taxonomy ID: 9606

Genbank common name: **human**

Inherited blast name: **primates**

Rank: species

Genetic code: [Translation table 1 \(Standard\)](#)

Mitochondrial genetic code: [Translation table 2 \(Vertebrate Mitochondrial\)](#)

Other names:

common name: **man**

authority: **Homo sapiens Linnaeus, 1758**

### [Lineage](#)(full)

[cellular organisms](#); [Eukaryota](#); [Opisthokonta](#); [Metazoa](#); [Eumetazoa](#); [Bilateria](#); [Deuterostomia](#); [Chordata](#); [Craniata](#); [Vertebrata](#); [Gnathostomata](#); [Teleostomi](#); [Euteleostomi](#); [Sarcopterygii](#); [Dipnotetrapodomorpha](#); [Tetrapoda](#); [Amniota](#); [Mammalia](#); [Theria](#); [Eutheria](#); [Boreoeutheria](#); [Euarchontoglires](#); [Primates](#); [Haplorrhini](#); [Simiiformes](#); [Catarrhini](#); [Hominoidea](#); [Hominidae](#); [Homininae](#); [Homo](#)

Entrez records		
Database name	Subtree links	Direct links
Nucleotide	<a href="#">14,908,727</a>	<a href="#">14,908,681</a>
Nucleotide EST	<a href="#">8,705,106</a>	<a href="#">8,705,106</a>
Nucleotide GSS	<a href="#">1,783,249</a>	<a href="#">1,781,923</a>
Protein	<a href="#">1,175,370</a>	<a href="#">1,175,052</a>
Structure	<a href="#">35,858</a>	<a href="#">35,858</a>
Genome	<a href="#">1</a>	<a href="#">1</a>
Popset	<a href="#">23,768</a>	<a href="#">23,767</a>
SNP	<a href="#">336,828,091</a>	<a href="#">336,828,091</a>
Domains	<a href="#">43</a>	<a href="#">43</a>
GEO Datasets	<a href="#">1,256,551</a>	<a href="#">1,256,551</a>
UniGene	<a href="#">130,056</a>	<a href="#">130,056</a>
PubMed Central	<a href="#">24,983</a>	<a href="#">24,918</a>
Gene	<a href="#">221,647</a>	<a href="#">221,574</a>
HomoloGene	<a href="#">18,713</a>	<a href="#">18,713</a>
SRA Experiments	<a href="#">833,162</a>	<a href="#">832,633</a>
Probe	<a href="#">27,382,489</a>	<a href="#">27,382,489</a>
Assembly	<a href="#">108</a>	<a href="#">108</a>
Bio Project	<a href="#">36,789</a>	<a href="#">36,776</a>
Bio Sample	<a href="#">3,071,592</a>	<a href="#">3,071,440</a>
Bio Systems	<a href="#">3,077</a>	<a href="#">3,077</a>
Clone DB	<a href="#">17,630,270</a>	<a href="#">17,630,270</a>
dbVar	<a href="#">5,169,896</a>	<a href="#">5,157,957</a>
GEO Profiles	<a href="#">61,958,910</a>	<a href="#">61,958,910</a>
PubChem BioAssay	<a href="#">311,518</a>	<a href="#">311,510</a>
Protein Clusters	<a href="#">13</a>	<a href="#">13</a>
Taxonomy	<a href="#">3</a>	<a href="#">1</a>

# Capsaspora

Taxonomy Browser

EntrezPubMedNucleotideProteinGenomeStructurePMCTaxonomyBooks

Search for  as  complete name ☐ lock

Display  2 levels using filter:  has genome sequences

### Capsaspora owczarzaki ATCC 30864

*Taxonomy ID:* 595528  
*Inherited blast name:* **eukaryotes**  
*Rank:* no rank  
*Genetic code:* [Translation table 1 \(Standard\)](#)  
*Mitochondrial genetic code:* [Translation table 4 \(Mold Mitochondrial; Protozoan Mitochondrial; Coelenterate Mitochondrial; Mycoplasma; Spiroplasma\)](#)

Lineage( full )  
[cellular organisms](#); [Eukaryota](#); [Opisthokonta](#); [Opisthokonta incertae sedis](#); [Ichthyosporea](#); [Capsaspora](#); [Capsaspora owczarzaki](#)

Entrez records	
Database name	Direct links
Nucleotide	<a href="#">8,993</a>
Protein	<a href="#">18,915</a>
Genome	<a href="#">1</a>
PubMed Central	<a href="#">3</a>
Gene	<a href="#">9,504</a>
SRA Experiments	<a href="#">11</a>
Assembly	<a href="#">2</a>
Bio Project	<a href="#">2</a>
Bio Sample	<a href="#">12</a>
Taxonomy	<a href="#">1</a>

## External Information Resources (NCBI LinkOut)

LinkOut	Subject	LinkOut Provider
<a href="#">GOLD: Go0003455</a>	organism-specific	<a href="#">Genomes On Line Database</a>
<a href="#">WebScipio: Capsaspora owczarzaki ATCC 30864</a>	organism-specific	<a href="#">WebScipio - eukaryotic gene identification</a>
<a href="#">diArk: Capsaspora owczarzaki ATCC 30864</a>	organism-specific	<a href="#">diArk a resource for eukaryotic genome research</a>

# Capsaspora

NCBI Taxonomy Browser

Entrez PubMed Nucleotide Protein

Search for  as  complete name ☐ lock

Display  2 levels using filter:  has genome sequences

☐ Nucleotide ☐ Nucleotide EST ☐ Nucleotide GSS ☐ Protein ☐ Structure ☐ Genome ☐ Popset  
☐ Domains ☐ GEO Datasets ☐ UniGene ☐ PubMed Central ☐ Gene ☐ HomoloGena ☐ SRA Expe  
☐ Assembly ☐ LinkOut ☐ BLAST ☐ TRACE ☐ Host ☐ Viral Host ☐ Bio Project  
☐ Bio Systems ☐ Clone DB ☐ dbVar ☐ GEO Profiles ☐ PubChem BioAssay ☐ Protein Clusters

Lineage (full): [root](#); [cellular organisms](#); [Eukaryota](#)

- [Opisthokonta](#) *Click on organism name to get more information.*
  - [Choanoflagellida](#)
    - [Craspedida](#)
  - [Fungi](#) (fungi)
    - [Blastocladiomycota](#)
    - [Chytridiomycota](#)
    - [Cryptomycota](#)
    - [Dikarya](#)
    - [Microsporidia](#)
    - [Mucoromycota](#)
    - [Neocallimastigomycota](#)
    - [Zoopagomycota](#)
    - [unclassified Fungi](#)
  - [Metazoa](#) (metazoans)
    - [Eumetazoa](#)
    - [Mesozoa](#)
    - [Placozoa](#) (placozoans)
    - [Porifera](#) (sponges)
  - [Nucleariidae and Fonticula group](#)
    - [Fonticula](#)
      - [Fonticula-like sp. SCN 57-25](#)
  - [Opisthokonta incertae sedis](#)
    - [Ichthyosporea](#)

# Entrez: Integrated Molecular and Sequence Databases

Search NCBI databases

all[**sb**]

Search

About 1,353,092,513 search results for "all[**sb**]"

Literature

Books	334,634	
MeSH	253,991	
NLM Catalog	1,511,375	
PubMed	24,324,879	
PubMed Central	3,264,562	

Health

ClinVar	125,404	
dbGaP	166,753	
GTR	35,977	
MedGen	260,910	
OMIM	23,795	online mendelian inheritance in man
PubMed Health	51,412	clinical effectiveness, disease and drug reports

Genomes

Assembly	33,680	genomic assembly information
BioProject	138,999	biological projects providing data to NCBI
BioSample	2,865,750	descriptions of biological source materials
Clone	37,024,042	genomic and cDNA clones
dbVar	4,305,990	genome structural variation studies
Epigenomics	6,635	epigenomic studies and display tools
Genome	10,777	genome sequencing projects by organism
GSS	37,646,655	genome survey sequences
Nucleotide	155,458,773	DNA and RNA sequences
Probe	31,890,151	sequence-based probes and primers
SNP	444,458,769	short genetic variations
SRA	1,064,488	high-throughput DNA and RNA sequence read archive
Taxonomy	1,310,804	taxonomic classification and nomenclature catalog

Central Resources / Databases

- Taxonomy
- BioProject
- Assembly
- Gene

Follow links to others when needed

Nucleotide, Protein, SRA

Domains

Protein	153,454,195	expressed sequence tag sequences
Protein Clusters	820,546	collected information about gene loci
Structure	103,644	functional genomics studies

Chemicals

BioSystems	653,443	gene expression and molecular abundance profiles
PubChem BioAssay	1,112,105	homologous gene sets for selected organisms
PubChem Compound	62,041,347	sequence sets from phylogenetic and population studies
PubChem Substance	177,330,151	clusters of expressed transcripts




The Entrez system: 39 (and counting) integrated databases


# Start in High Level Resources

If your question is about data for ...

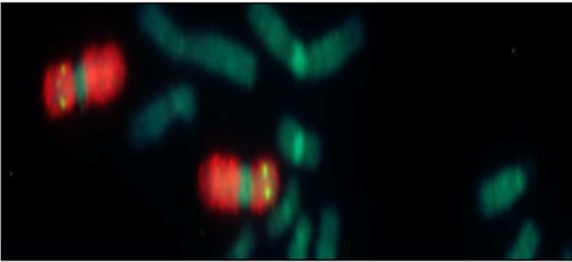
- an organism -> Taxonomy
- a gene name -> Gene (common organisms)
- a large-scale project -> BioProject
- a bacterial genome -> Genome
- a genome sequence -> Assembly

# NCBI-Gene ([www.ncbi.nlm.nih.gov/gene](http://www.ncbi.nlm.nih.gov/gene))

 [Resources](#)  [How To](#)  [Sign in to NCBI](#)

Gene    [Help](#)


[Advanced](#)



## Gene

Gene integrates information from a wide range of species. A record may include nomenclature, Reference Sequences (RefSeqs), maps, pathways, variations, phenotypes, and links to genome-, phenotype-, and locus-specific resources worldwide.

### Using Gene

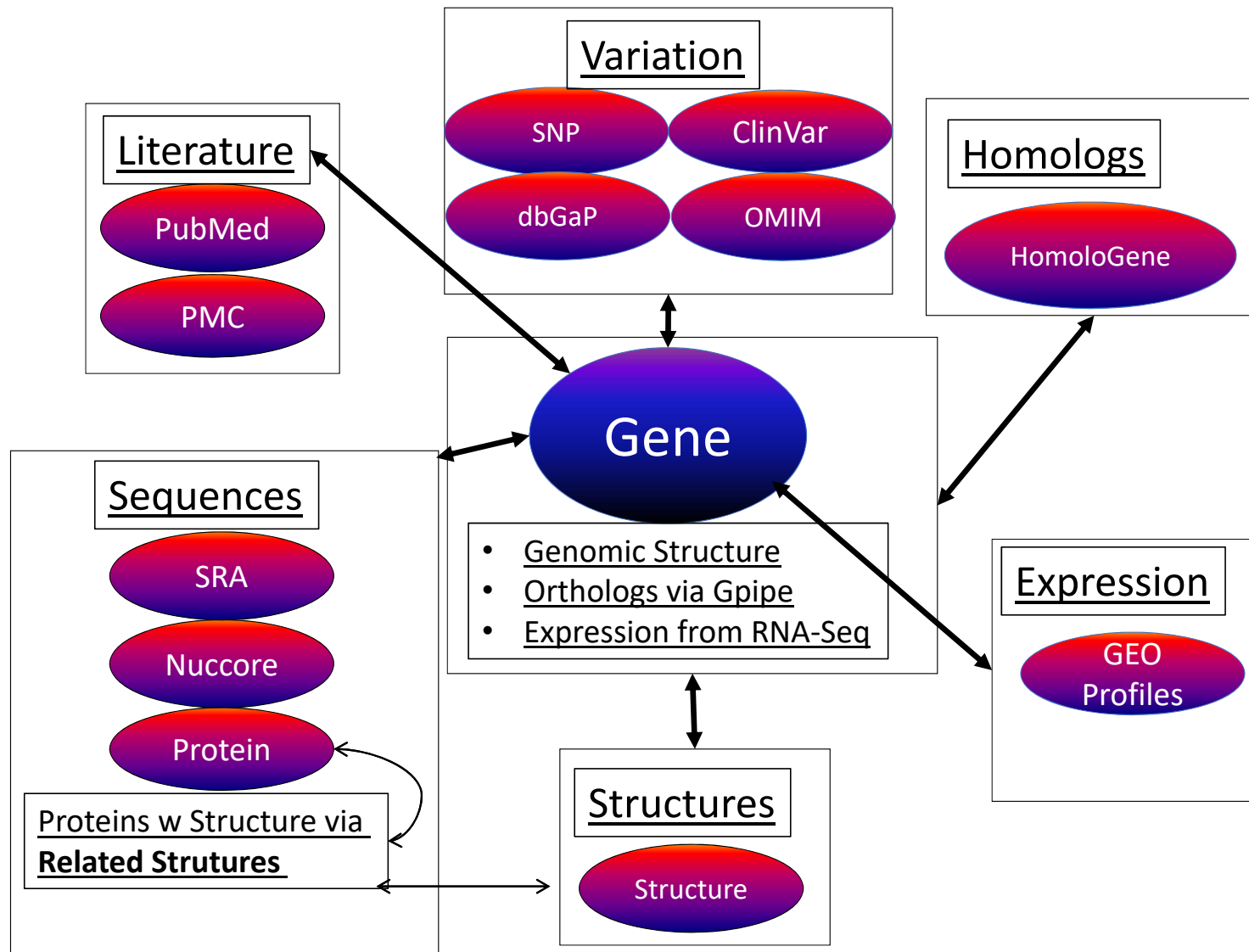
- [Gene Quick Start](#)
- [FAQ](#)
- [Download/FTP](#)
- [RefSeq Mailing List](#)
- [Gene News](#) 
- [Factsheet](#)

### Gene Tools

- [Submit GeneRIFs](#)
- [Submit Correction](#)
- [Statistics](#)
- [BLAST](#)
- [Genome Workbench](#)
- [Splign](#)

### Other Resources

- [HomoloGene](#)
- [OMIM](#)
- [RefSeq](#)
- [RefSeqGene](#)
- [UniGene](#)
- [Protein Clusters](#)



# A gene record: SORT1

(<https://www.ncbi.nlm.nih.gov/gene/6272>)

The image shows a screenshot of the NCBI Gene database record for SORT1 (sortilin 1) in Homo sapiens. The record is titled "SORT1 sortilin 1 [Homo sapiens (human)]" with Gene ID 6272, updated on 20-Feb-2017. The interface includes a search bar at the top, a "Full Report" dropdown, and a "Send to" dropdown. The main content area is divided into two columns. The left column contains a list of sections: Summary, Genomic context, Genomic regions, transcripts, and products, Expression, Bibliography, Phenotypes, Variation, Pathways from BioSystems, Interactions, General gene information, General protein information, NCBI Reference Sequences (RefSeq), Related sequences, and Additional links. The right column contains a list of links: Table of contents, Genome Browsers, Related information, Links to other resources, General information, Related sites, Feedback, Subscription, and Recent activity. Three blue arrows point from the text "Genomic Structure" to the "Genomic context" and "Genomic regions, transcripts, and products" sections. Another set of blue arrows points from the text "Functional information" to the "Expression", "Bibliography", "Phenotypes", "Variation", "Pathways from BioSystems", "Interactions", "General gene information", and "General protein information" sections. A third set of blue arrows points from the text "Sequence" to the "NCBI Reference Sequences (RefSeq)" and "Related sequences" sections.

Gene

Gene

Advanced

Search

Help

Full Report

Send to

Hide sidebar >>

**SORT1** sortilin 1 [ *Homo sapiens* (human) ]

Gene ID: 6272, updated on 20-Feb-2017

Summary

Genomic context

Genomic regions, transcripts, and products

Expression

Bibliography

Phenotypes

Variation

Pathways from BioSystems

Interactions

General gene information

General protein information

NCBI Reference Sequences (RefSeq)

Related sequences

Additional links

Table of contents

Genome Browsers

Related information

Links to other resources

General information

Related sites

Feedback

Subscription

Recent activity

Genomic Structure

Functional information

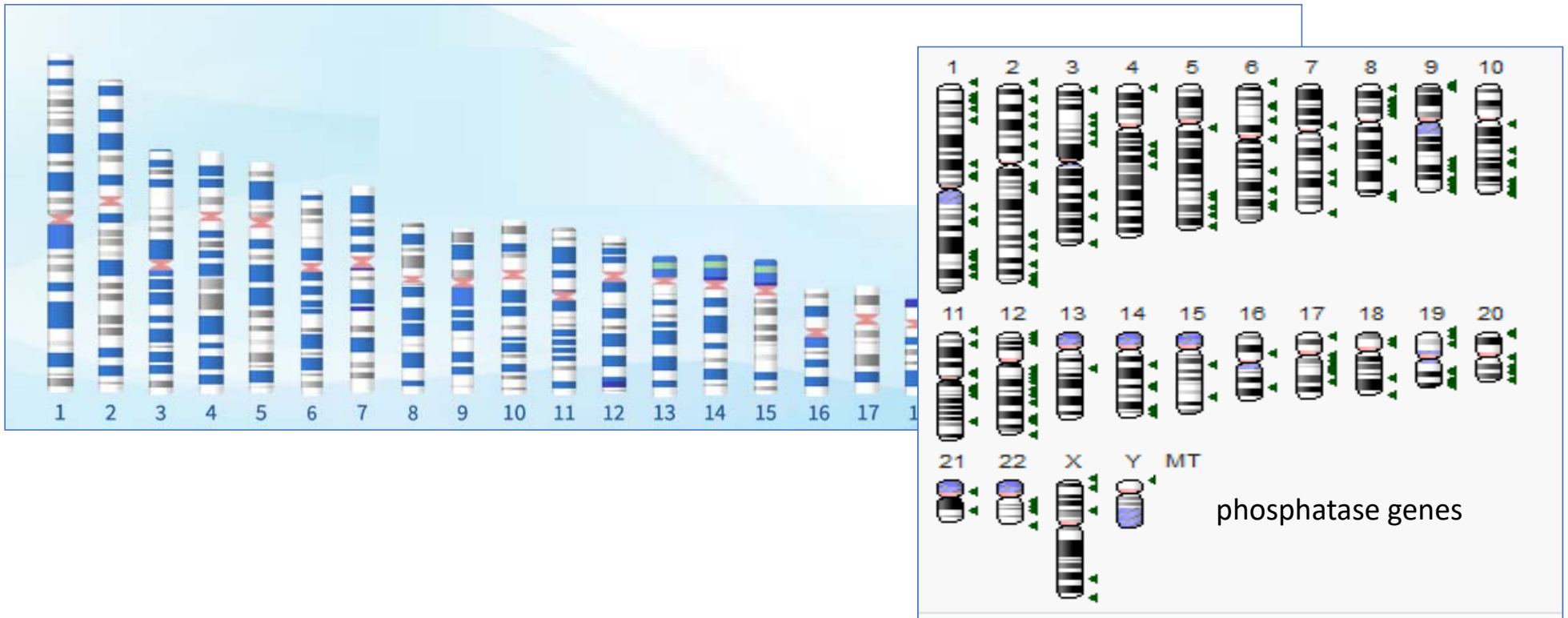
Sequence



# Some Examples

1. Where is the gene located (chromosome and position) in the genome?
2. What are the Reference genomic, transcript and protein sequences for the gene?
3. What variations are present in the gene and are they associated with disease?
4. What are the equivalent genes (homologs) in other species?

# 1. Where is the gene located in the genome?



# Genome context and structure:

DSG2[All Fields] AND "Homo sapiens"[Organism]

**DSG2 desmoglein 2 [ *Homo sapiens* (human) ]**  
Gene ID: 1829, updated on 20-Feb-2017

**Summary**

Official Symbol DSG2 provided by HGNC

**Table of contents**

- Summary
- Genomic context
- Genomic regions, transcripts, and products
- Expression

**Genomic context**

Location: 18q12.1

Exon count: 15

Annotation release	Status	Assembly
<a href="#">108</a>	current	GRCh38.p7 ( <a href="#">GCF</a> )
<a href="#">105</a>	previous assembly	GRCh37.p13 ( <a href="#">GCF</a> )

**Chromosome**

**Genomic regions, transcripts, and products**

Go to [reference sequence details](#)

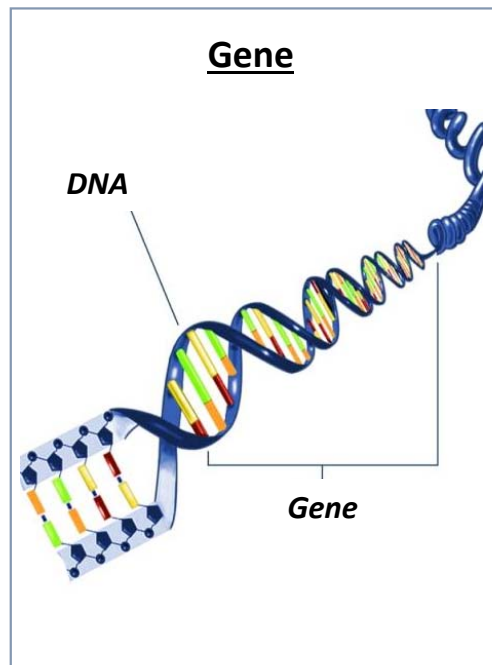
Genomic Sequence: [NC\\_000018.10](#) Chromosome 18 Reference GRCh38.p7 Primary Assembly

Go to nucleotide: [Graphics](#) [FASTA](#) [GenBank](#)

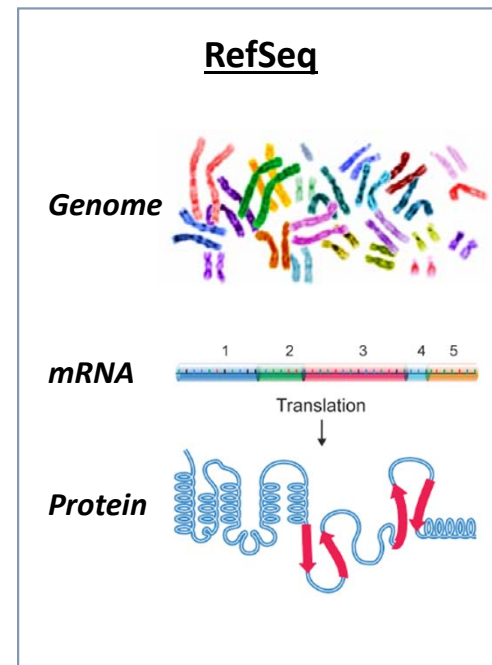
A lab report has shown that a patient has a deletion of chromosome 5p which corresponds to nts 204,700-5,500,000. How do I find out what genes are involved?

A recent research article has shown that a genomic region is associated with a disease. How can I find what genes are there?

2. What are the Reference genomic, transcript and protein sequences for the gene?



*Gene level*



*Transcript and protein level*

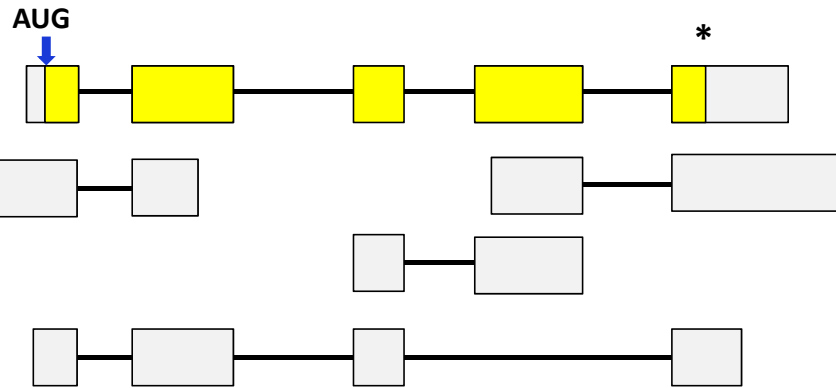
# NCBI Reference Sequences (RefSeq) Project

Reference genome sequence

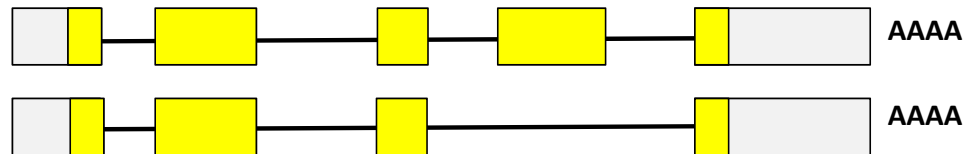


INSDC transcripts  
(Primary data)

INSDC: The International  
Nucleotide Sequence  
Database Collaboration  
shared by NCBI, EMBL-  
EBI and DDBJ.



RefSeq variants



# RefSeq

- Created by NCBI data curators and computational algorithms
- Used by researchers & collaborators as a reference standard
- For selected eukaryotes, represent all molecules in the central dogma
  - Genomic (DNA), Transcripts (mRNA), Proteins
- Distinct accession series – with an underscore!
  - Genomic: NC\_, AC\_, NG\_, NT\_, NW\_
  - Transcripts: NM\_, NR\_, XM\_, XR\_
  - Proteins: NP\_, XP\_

# RefSeq

## DSG2 desmoglein 2 [ *Homo sapiens* (human) ]

Gene ID: 1829, updated on 5-Mar-2017

Genomic context

### Summary

**Official Symbol** DSG2 provided by [HGNC](#)  
**Official Full Name** desmoglein 2 provided by [HGNC](#)  
**Primary source** [HGNC:HGNC:3049](#)  
**See related** [Ensembl:ENSG00000046604](#) [MIM:125671](#) [Vega:OTTHUMG000000](#)

### RefSeqs of Annotated Genomes: Homo sapiens Annotation Release 108 [details...](#)

The following sections contain reference sequences that belong to a specific genome build. [Explain](#)

#### Reference GRCh38.p7 Primary Assembly

##### Genomic

##### 1. NC\_000018.10 Reference GRCh38.p7 Primary Assembly

Range	31498004..31549008
Download	<a href="#">GenBank</a> , <a href="#">FASTA</a> , <a href="#">Sequence Viewer (Graphics)</a>

a; Euteleostomi  
Hominidae; Homo

family and cadherin cell adhesion molecule superfamily of  
membrane glycoprotein components of desmosomes,  
and other cell types. The encoded preproprotein is  
glycoprotein. This gene is present in a gene cluster with  
some 18. Mutations in this gene have been associated  
initial, 10. [provided by RefSeq, Jan 2016]

Ontology

General protein information

NCBI Reference Sequences (RefSeq)

Related sequences

Additional links

Locus-specific Databases

### NCBI Reference Sequences (RefSeq)

#### RefSeqs maintained independently of Annotated Genomes

These reference sequences exist independently of genome builds. [Explain](#)

##### Genomic

##### 1. NG\_007072.3 RefSeqGene

Range	4823..55610
Download	<a href="#">GenBank</a> , <a href="#">FASTA</a> , <a href="#">Sequence Viewer (Graphics)</a> , <a href="#">LRG_397</a>

##### mRNA and Protein(s)

##### 1. [NM\\_001943.4](#) → [NP\\_001934.2](#) desmoglein-2 preproprotein

[See identical proteins and their annotated locations for NP\\_001934.2](#)



I am studying the globin gene cluster on chromosome 11 and I need to get the protein sequences for all gene family members.

I am interested in the promoter region of genes that don't code for proteins. How can I get this information?

### 3. What variations are present in the gene and are they associated with disease?

#### C9orf72 chromosome 9 open reading frame 72 [ *Homo sapiens* (human) ]

Gene ID: 203228, updated on 3-Sep-2017

##### Summary

**Official Symbol** C9orf72 provided by [HGNC](#)  
**Official Full Name** chromosome 9 open reading frame 72 provided by [HGNC](#)  
**Primary source** [HGNC:HGNC:28337](#)  
**See related** [Ensembl:ENSG00000147894](#) [MIM:614260](#); [Vega:OTTHUMG00000019716](#)  
**Gene type** protein coding  
**RefSeq status** REVIEWED  
**Organism** [Homo sapiens](#)  
**Lineage** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo  
**Also known as** ALSFTD; FTDALS; DENNL72; FTDALS1  
**Summary** The protein encoded by this gene plays an important role in the regulation of endosomal trafficking, and has been shown to interact with Rab proteins that are involved in autophagy and endocytic transport. Expansion of a GGGGCC repeat from 2-22 copies to 700-1600 copies in the intronic sequence between alternate 5' exons in transcripts from this gene is associated with 9p-linked ALS (amyotrophic lateral sclerosis) and FTD (frontotemporal dementia) (PMID: 21944778, 21944779). Studies suggest that hexanucleotide expansions could result in the selective stabilization of repeat-containing pre-mRNA, and the accumulation of insoluble dipeptide repeat protein aggregates that could be pathogenic in FTD-ALS patients (PMID: 23393093). Alternative splicing results in multiple transcript variants encoding different isoforms. [provided by RefSeq, Jul 2016]  
**Orthologs** [mouse](#) [all](#)

##### Genomic context

##### Table of contents

Summary
Genomic context
Genomic regions, transcripts, and products
Expression
Bibliography
Phenotypes
Variation
Pathways from BioSystems
Interactions
General gene information Markers, Clone Names, Homology, Gene Ontology
General protein information
NCBI Reference Sequences (RefSeq)
Related sequences
Additional links Locus-specific Databases

# dbSNP, dbVar, ClinVar and their overlap

- dbSNP: Database of Short Genetic Variations (Single-nucleotide polymorphism); 53 organisms
- dbVar: Database of Genomic Structural Variations (insertions and deletions, tri-/ hexa- nucleotide repeat expansion)
- ClinVar: Database of relationships between human variations and phenotypes (observed health status)

# Variation and Disease

## Phenotypes

[Find tests for this gene in the NIH Genetic Testing Registry \(GTR\)](#)

[Review eQTL and phenotype association data in this region using PheGenI](#)

Associated conditions

### Description

[Amyotrophic lateral sclerosis](#)

MedGen: [C0002736](#), GeneReviews: [Amyotrophic Lateral Sclerosis Overview](#)

[Amyotrophic lateral sclerosis and/or frontotemporal dementia 1](#)

MedGen: [C1862937](#), OMIM: [105550](#), GeneReviews: [Amyotrophic Lateral Sclerosis Overview](#), [C9orf72-Related Amyotrophic Lateral Sclerosis and Frontotemporal Dementia](#)

## Variation

[See variants in ClinVar](#)

[See studies and variants in dbVar](#)

[See Variation Viewer \(GRCh37.p13\)](#)

[See Variation Viewer \(GRCh38\)](#)

What diseases can be caused by variations in the tyrosine hydroxylase gene? Can I get a list of all disease causing single nucleotide variants that affect the coding regions with their positions? Are there any common protein variants in this gene?

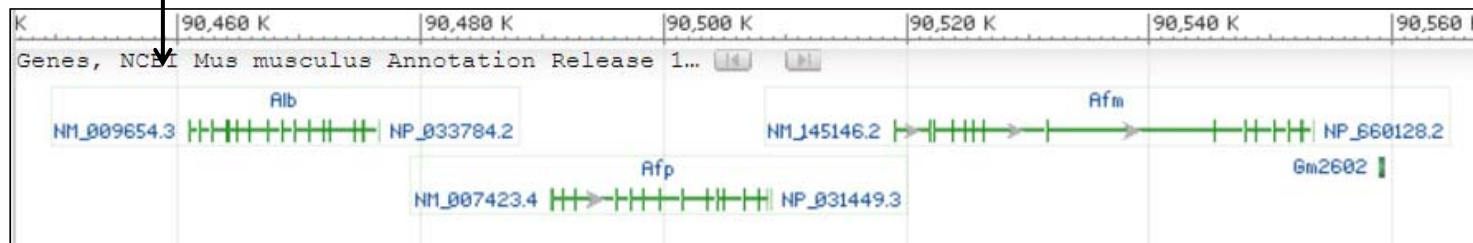
4. What are the equivalent genes (homologs) in other species?

# Orthologs and Paralogs

Many gene families have multiple members; Serum albumin gene family: albumin (Alb), alpha-fetoprotein (AFP) and afamin (Afm).



- Human ALB, AFP and AFM are **paralogous** genes.
- **Paralogs are derived by gene duplication within a species (or ancestral species).**



- ALB (human) and Alb (mouse) are **orthologous** genes.
- **Orthologs are derived by speciation events (homologs between species).**

# 4. What are the equivalent genes (homologs) in other species?

## C9orf72 chromosome 9 open reading frame 72 [ *Homo sapiens* (human) ]

Gene ID: 203228, updated on 3-Sep-2017

### Summary

**Official Symbol** C9orf72 provided by [HGNC](#)  
**Official Full Name** chromosome 9 open reading frame 72 provided by [HGNC](#)  
**Primary source** [HGNC:HGNC:28337](#)  
**See related** [Ensembl:ENSG00000147894](#) [MIM:614260](#); [Vega:OTTHUMG00000019716](#)  
**Gene type** protein coding  
**RefSeq status** REVIEWED  
**Organism** [Homo sapiens](#)  
**Lineage** Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini; Catarrhini; Hominidae; Homo  
**Also known as** ALSFTD; FTDALS; DENNL72; FTDALS1  
**Summary** The protein encoded by this gene plays an important role in the regulation of endosomal trafficking, and has been shown to interact with Rab proteins that are involved in autophagy and endocytic transport. Expansion of a GGGGCC repeat from 2-22 copies to 700-1600 copies in the intronic sequence between alternate 5' exons in transcripts from this gene is associated with 9p-linked ALS (amyotrophic lateral sclerosis) and FTD (frontotemporal dementia) (PMID: 21944778, 21944779). Studies suggest that hexanucleotide expansions could result in the selective stabilization of repeat-containing pre-mRNA, and the accumulation of insoluble dipeptide repeat protein aggregates that could be pathogenic in FTD-ALS patients (PMID: 23393093). Alternative splicing results in multiple transcript variants encoding different isoforms. [provided by RefSeq, Jul 2016]  
**Orthologs** [mouse](#) [all](#)

### Homology

[Homologs of the C9orf72 gene](#): The C9orf72 gene is conserved in chimpanzee, Rhesus monkey, dog, cow, mouse, rat, chicken, and zebrafish.  
[Orthologs from Annotation Pipeline](#): 216 organisms have orthologs with human gene C9orf72  
[The Hierarchical Catalog of Orthologs](#)

### Table of contents

Summary
Genomic context
Genomic regions, transcripts, and products
Expression
Bibliography
Phenotypes
Variation
Pathways from BioSystems
Interactions
General gene information
Markers, Clone Names, Homology, Gene Ontology
General protein information
NCBI Reference Sequences (RefSeq)
Related sequences
Additional links



## 4. What are the equivalent genes (homologs) in other species?

### HomoloGene:10137. Gene conserved in Euteleostomi

#### Genes

Genes identified as putative homologs of one another during the construction of HomoloGene.

C9orf72, *H.sapiens*  
chromosome 9 open reading frame 72  
C9H9orf72, *P.troglodytes*  
chromosome 9 open reading frame, human C9orf72  
C15H9orf72, *M.mulatta*  
chromosome 9 open reading frame 72 ortholog  
C11H9orf72, *C.lupus*  
chromosome 11 open reading frame, human C9orf72  
C8H9orf72, *B.taurus*  
chromosome 8 open reading frame, human C9orf72  
3110043O21Rik, *M.musculus*  
RIKEN cDNA 3110043O21 gene  
RGD1359108, *R.norvegicus*  
similar to RIKEN cDNA 3110043O21  
C9ORF72, *G.gallus*  
chromosome Z open reading frame, human C9orf72  
zgc:100846, *D.rerio*  
zgc:100846

#### Protein Alignments

Protein multiple alignment, pairwise similarity scores and evolutionary distances.

Show Multiple Alignment

#### Proteins

Proteins used in sequence comparisons and their conserved domain architectures.

NP\_001242983.1 —————  
481 aa  
XP\_003951454.1 —————  
481 aa  
XP\_002800118.1 —————  
481 aa  
XP\_005626755.1 —————  
481 aa  
NP\_001096558.1 —————  
481 aa  
NP\_001074812.1 —————  
481 aa  
NP\_001007703.1 —————  
481 aa  
XP\_424945.2 —————  
481 aa  
NP\_991166.1 —————  
462 aa

#### Conserved Domains

Conserved Domains from CDD found in protein sequences by rpsblast searching.

C9orf72-like (pfam15019)

■ C9orf72-like protein family.

1. A lab report has shown that a patient has a deletion of chromosome 5p which corresponds to nts 204,700-5,500,000. How do I find out what genes are involved?
2. I am studying the globin gene cluster on 11 and I need to get the protein sequences for all gene family members.
3. I am interested in the promoter region of genes that don't code for proteins. How can I get this information?
4. What diseases can be caused by variations in the tyrosine hydroxylase gene? Can I get a list of all disease causing single nucleotide variants that affect the coding regions with their positions? Are there any common protein variants in this gene?



# Data formats

- GenBank DNA sequence entry
- EMBL sequence entry
- FASTA sequence format
- ASN.1: Abstract Syntax Notation
- XML
- PDB
- Phylip: Phylogenetic Inference package
- Clustal
- MSF: multiple sequence format

# One format can be converted to another

- READSEQ: <http://iubio.bio.indiana.edu/cgi-bin/readseq.cgi>
  - <http://iubio.bio.indiana.edu/soft/molbio/readseq/java/Readseq2-help.html>
- EMBOSS Seqret: [www.ebi.ac.uk/Tools/sfc/emboss\\_seqret/](http://www.ebi.ac.uk/Tools/sfc/emboss_seqret/)
- SEQIO: <http://search.cpan.org/dist/BioPerl/Bio/SeqIO.pm>