



# An Introduction to Molecular Phylogenetic Inference



**Kirill Vinnikov**

PhD Student,  
Department of Biology

ITS Cyberinfrastructure Workshop Series  
September 9, 2016

- CASE STUDY dataset
- Introduction to NCBI's sequence resources
- Using different alignment approaches
- Sequence file formats
- Other sequence databases

- Text editor (e.g. TextWrangler, Notepad++)
- MEGA 7.0 ([www.megasoftware.net](http://www.megasoftware.net))
- Sequence Matrix ([gaurav.github.io/taxondna/](http://gaurav.github.io/taxondna/))

Available online 1 April 2004

# CASE STUDY: extant bear species



Polar bear  
*Ursus maritimus*



Brown bear  
*Ursus arctos*



Sloth bear  
*Melursus ursinus*



Panda  
*Ailuropoda melanoleuca*



American black bear  
*Ursus americanus*



Spectacled bear  
*Tremarctos ornatus*



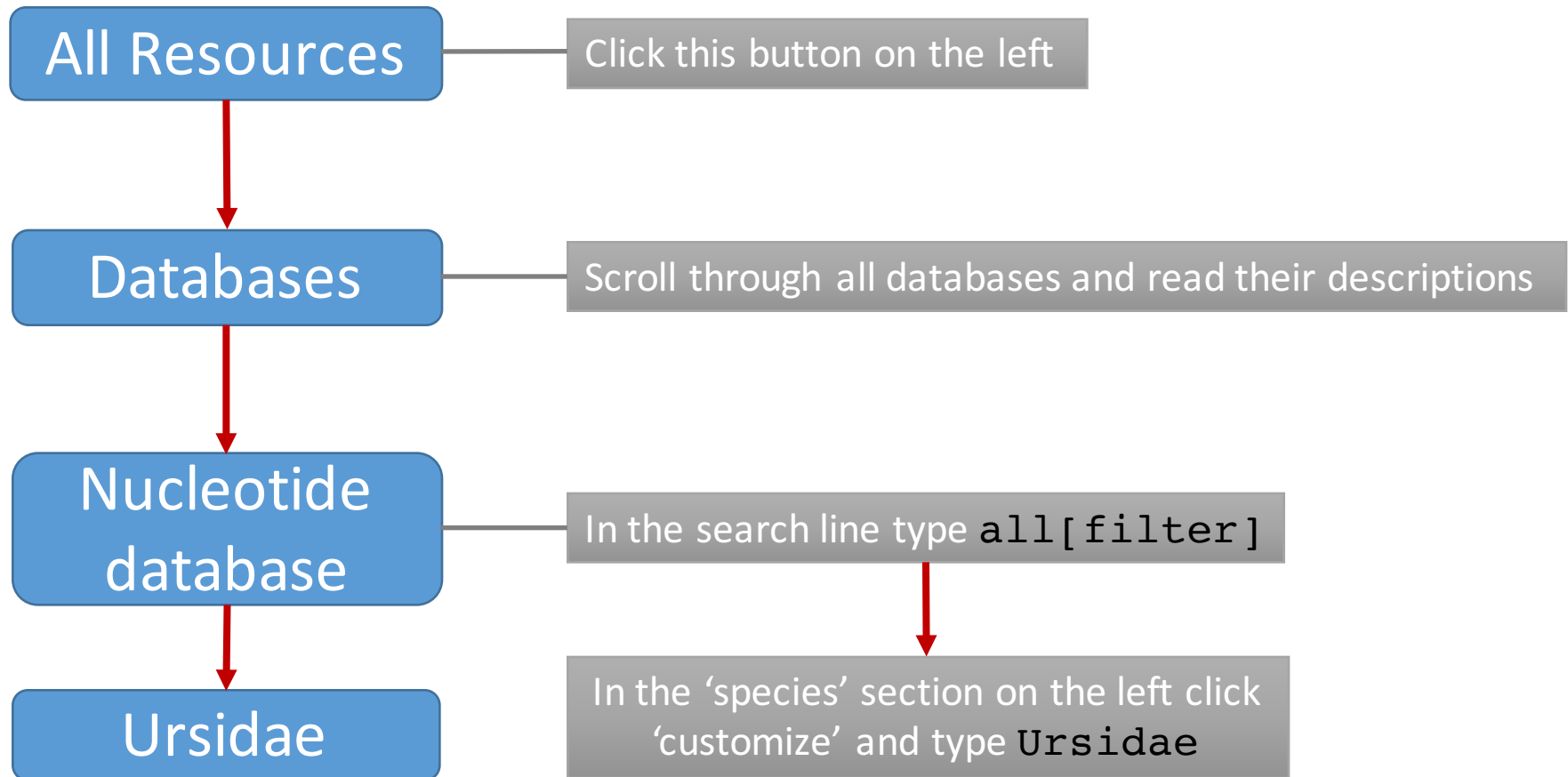
Sun bear  
*Helarctos malayanus*



Asian black bear  
*Ursus thibetanus*

**Keywords:** Interphotoreceptor retinoid binding protein; Transthyretin; Evolution; Phylogenetic analysis







**Nucleotide**

Nucleotide  Search

Create alert Advanced Help

---

**Species**

- Animals (110,055,153)
- Plants (34,873,683)
- Fungi (6,106,011)
- Protists (4,697,825)
- Bacteria (26,489,369)
- Archaea (526,967)
- Viruses (2,287,703)
- Ursidae (309,546)
- Customize ...

**Molecule types**

- genomic
- DNA/RNA (111,736,208)
- mRNA (39,508,054)
- rRNA (176,233)
- Customize ...

**Source databases**

- INSDC (GenBank) (172,853,511)
- RefSeq (35,924,195)
- Customize ...

**Genetic compartments**

- Chloroplast (888,551)
- Mitochondrion (3,689,309)
- Plasmid (152,833)
- Plastid (976,374)

**Sequence length**

Custom range...

**Release date**

Custom range...

**Revision date**

Custom range...

[Clear all](#)

[Show additional filters](#)

**Summary** ▾ 20 per page ▾ Sort by Default order ▾

**Items: 1 to 20 of 209053924**

<< First < Prev Page 1 of 10452697 Next > Last >>

**Found 324883430 nucleotide sequences.** Nucleotide (209053924) EST (76257001) GSS (39572505)

- ☐ [Macellibacteroides sp. HH-ZS, whole genome shotgun sequencing project](#)  
**4,081,830 rc other DNA**  
This entry is the master record for a whole genome shotgun sequencing project and contains no sequence data.  
 Accession: NZ\_LZEK00000000.1 GI: 1056764880  
[GenBank](#)
- ☐ [Macellibacteroides sp. HH-ZS contig\\_00067, whole genome shotgun sequence](#)  
**348 bp linear DNA**  
 Accession: NZ\_LZEK01000067.1 GI: 1056764864  
[GenBank](#) [FASTA](#) [Graphics](#)
- ☐ [Macellibacteroides sp. HH-ZS contig\\_00066, whole genome shotgun sequence](#)  
**397 bp linear DNA**  
 Accession: NZ\_LZEK01000066.1 GI: 1056764861  
[GenBank](#) [FASTA](#) [Graphics](#)
- ☐ [Macellibacteroides sp. HH-ZS contig\\_00065, whole genome shotgun sequence](#)  
**397 bp linear DNA**  
 Accession: NZ\_LZEK01000065.1 GI: 1056764853  
[GenBank](#) [FASTA](#) [Graphics](#)
- ☐ [Macellibacteroides sp. HH-ZS contig\\_00064, whole genome shotgun sequence](#)  
**540 bp linear DNA**  
 Accession: NZ\_LZEK01000064.1 GI: 1056764850  
[GenBank](#) [FASTA](#) [Graphics](#)
- ☐ [Macellibacteroides sp. HH-ZS contig\\_00063, whole genome shotgun sequence](#)  
**629 bp linear DNA**  
 Accession: NZ\_LZEK01000063.1 GI: 1056764845  
[GenBank](#) [FASTA](#) [Graphics](#)

**Filter your results:**

All (209053924)

[Manage Filters](#)

---

**Results by taxon**

---

**Find related data**

Database:

[Find items](#)

---

**Search details**

all[filter]

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

---

**Recent activity**

[all\[filter\] \(209053924\)](#) Nucleotide

[all\[filter\] AND \("Ursidae"\[ORGN\]\) \(309546\)](#) Nucleotide

[all \(0\)](#) Nucleotide

[all\[filter\] AND \(animals\[filter\] AND biomol\\_genomic\[PROP\]\) \(47722949\)](#) Nucleotide

[all\[filter\] AND \(animals\[filter\]\) \(110055153\)](#) Nucleotide



## Search Field Descriptors for NCBI Nucleotide Database

Search Field	Short	Description
<b>Accession</b>	ACCN	Sequence accession number in NCBI
<b>Author</b>	AUTH	Authors of the publication
<b>Feature Key</b>	FKEY	See Biological Features in the Feature Table
<b>Filter</b>	FILT	Filters subsets of the database
<b>Gene Name</b>	GENE	Gene name
<b>Issue</b>	ISS	Journal issue
<b>Journal</b>	JOUR	Journal name
<b>Keyword</b>	KYWD	Keywords
<b>Organism</b>	ORGN	Organism name (scientific or common)
<b>Properties</b>	PROP	Molecular type, source database, and other properties of the sequence record
<b>Protein Name</b>	PROT	Protein name
<b>Publication Date</b>	PDAT	Sequence publication date in NCBI
<b>Sequence Length</b>	SLEN	Total sequence length or the range
<b>Title</b>	TITL	Words and phrases in the title of the sequence record
<b>Volume</b>	VOL	Journal Volume Number

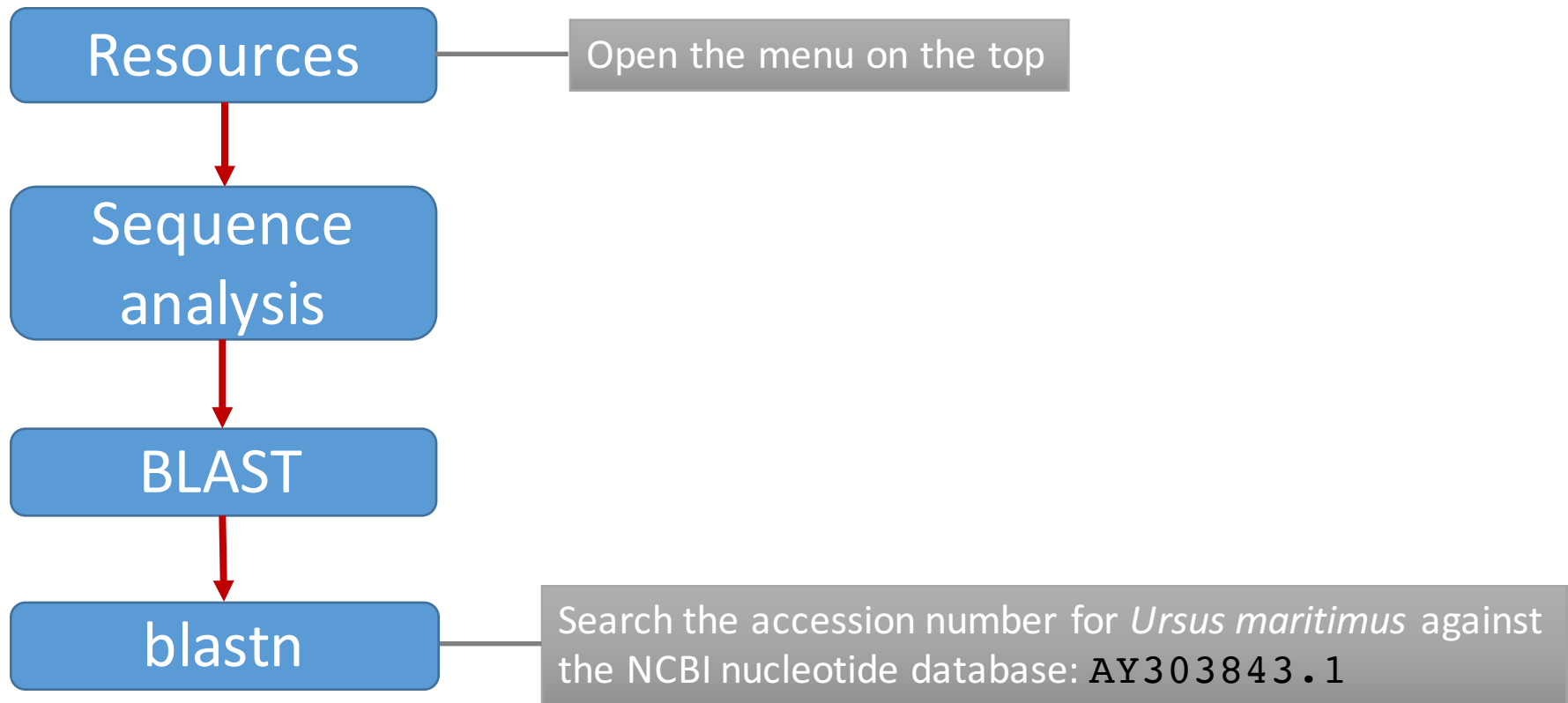
# Nucleotide database

```
"Yu L"[AUTH] AND "Ursidae"[ORGN] AND "irbp"[TITL]
```

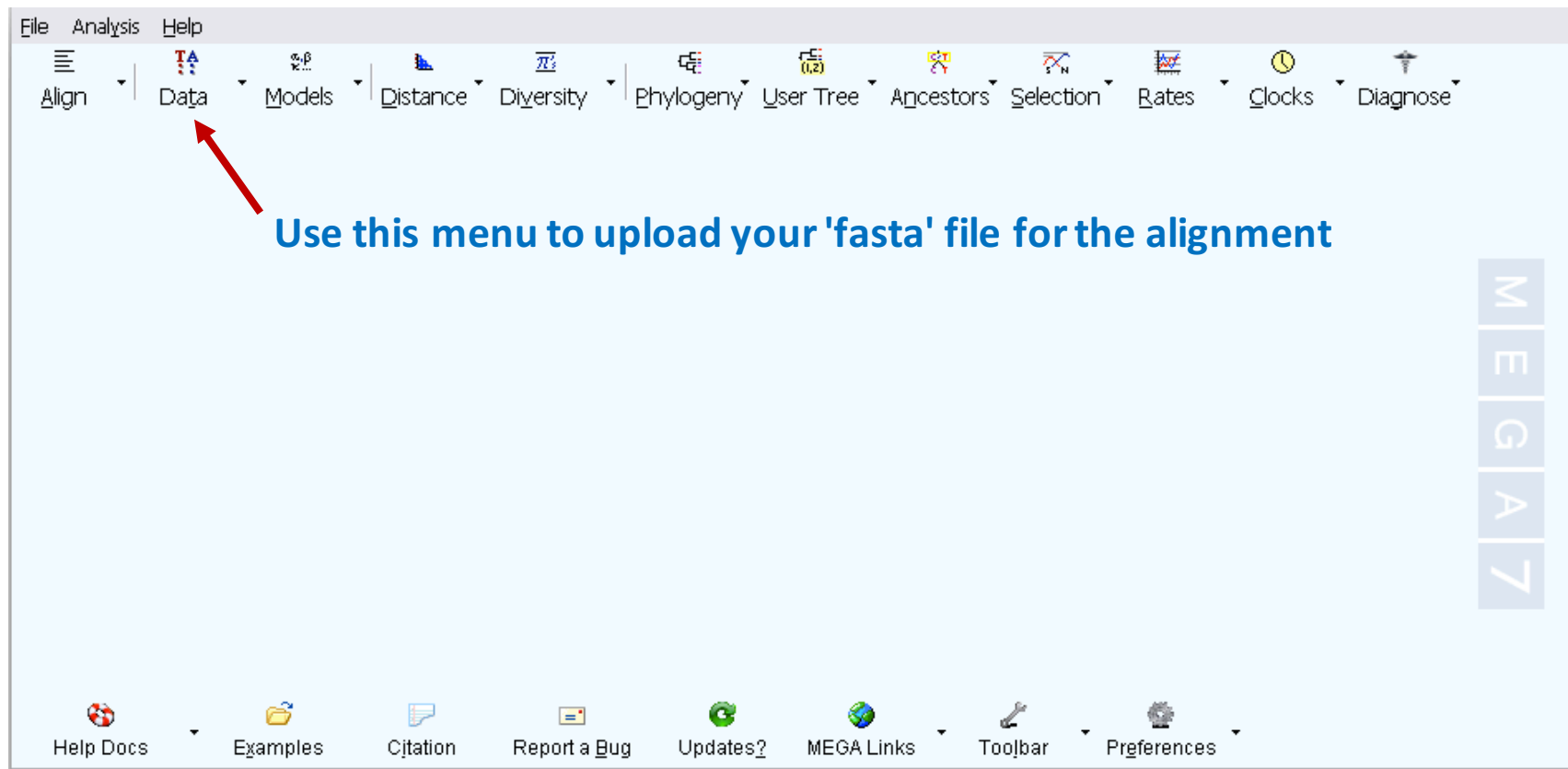
```
"mitochondrion"[filter] AND "Ursidae"[ORGN] AND "Talbot  
SL"[Author] AND 1140[SLEN]
```

## Save search results into 'cytb-ursidae.fasta' file

[See more...](#)



Choose one sequence and save its 15000:16500 region into 'cytb-outgroup.fasta' file

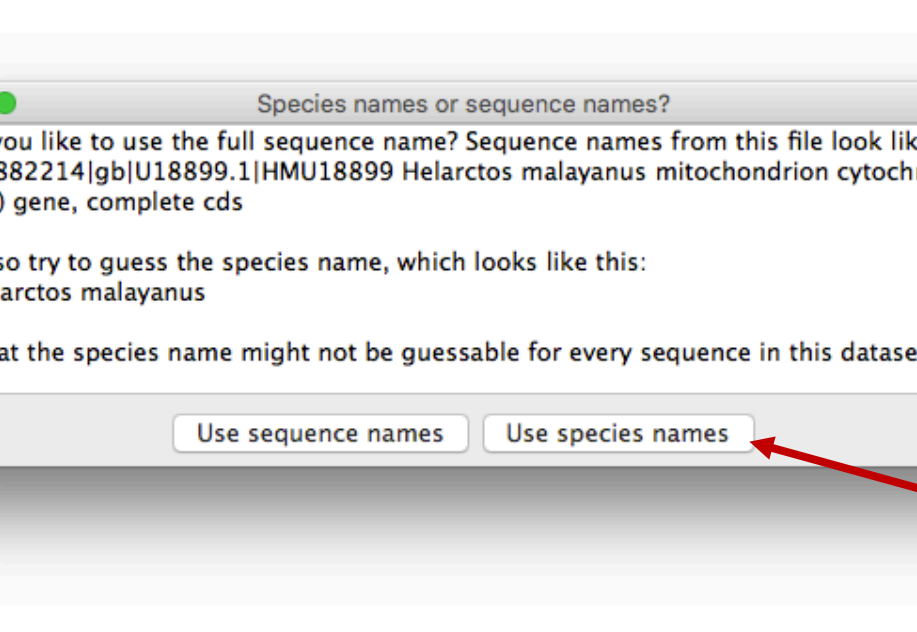


- 1. Trim columns with gaps on the both ends of your alignment!**
- 2. Check and correct ORF by translating your DNA sequences**



## Export your alignment into fasta file

## Sequence Matrix



Species names or sequence names?

Would you like to use the full sequence name? Sequence names from this file look like this:  
gi|882214|gb|U18899.1|HMU18899 Helarctos malayanus mitochondrion cytochrome  
b (cyt b) gene, complete cds

I can also try to guess the species name, which looks like this:  
Helarctos malayanus

Note that the species name might not be guessable for every sequence in this dataset.

Use sequence names    Use species names

No sequences loaded. Sorted by name.

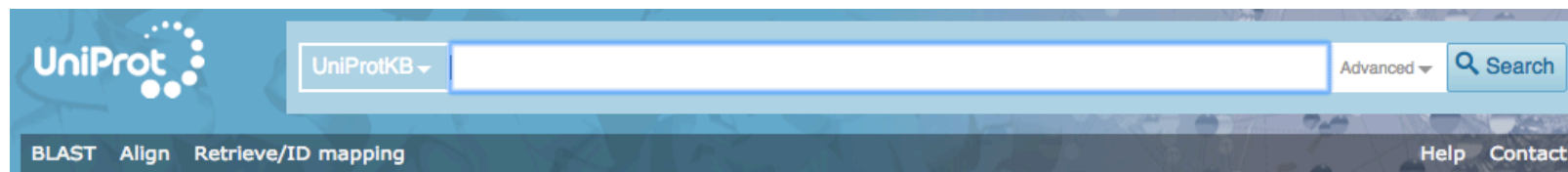
**Click to choose  
one sequence  
per species**

**Repeat the same format conversion for *cytochrome b***

bears.phy

# Other databases: UniProt example

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



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.

**UniProtKB**  
UniProt Knowledgebase

Swiss-Prot  
(551,705)  
 Manually annotated and reviewed.

TrEMBL  
(65,378,749)  
 Automatically annotated and not reviewed.

**UniRef**  
Sequence clusters

**UniParc**  
Sequence archive

**Proteomes**

**Supporting data**

Literature citations 	Taxonomy 	Subcellular locations 
Cross-ref. databases 	Diseases <b>XXX</b>	Keywords 

**News**

**Forthcoming changes**  
Planned changes for UniProt

---

**UniProt release 2016\_07**  
(Bacterial) immigration under control

---

**UniProt release 2016\_06**  
Strength through unity | Removal of the cross-references to NextBio | Change of URIs for neXtProt

---

**News archive**

## Other databases: UniProt example

Type in the search line:

```
name:"interphotoreceptor retinoid binding"
```

Type in the search line:

```
name:"interphotoreceptor retinoid binding" AND taxonomy:"Ursidae"
```

## UniProtKB results

Filter by<sup>i</sup>

Reviewed (9)

Swiss-Prot

Unreviewed (4,256)

TrEMBL

Popular organisms

Rat (7)

Mouse (3)

Zebrafish (2)

Bovine (1)

Human (1)

Other organisms

Search terms

Filter "interphotoreceptor retinoid binding" as:

protein name

View by

BLAST

Align

Download

Add to basket

Columns

1 to 25 of 4,265

Show 25

<input type="checkbox"/>	Entry	Entry name		Protein names	Gene names	Organism	Length	
<input type="checkbox"/>	P10745	RET3_HUMAN		Retinol-binding protein 3	RBP3	Homo sapiens (Human)	1,247	
<input type="checkbox"/>	Q7SZI7	RET3_XENLA		Retinol-binding protein 3	rbp3	Xenopus laevis (African clawed frog)	1,219	
<input type="checkbox"/>	P12661	RET3_BOVIN		Retinol-binding protein 3	RBP3	Bos taurus (Bovine)	1,286	
<input type="checkbox"/>	P49194	RET3_MOUSE		Retinol-binding protein 3	Rbp3	Mus musculus (Mouse)	1,234	
<input type="checkbox"/>	P12664	RET3_RABIT		Retinol-binding protein 3	RBP3	Oryctolagus cuniculus (Rabbit)	23	
<input type="checkbox"/>	P12662	RET3_PIG		Retinol-binding protein 3	RBP3	Sus scrofa (Pig)	25	
<input type="checkbox"/>	P12663	RET3_SHEEP		Retinol-binding protein 3	RBP3	Ovis aries (Sheep)	24	
<input type="checkbox"/>	P12666	RET3_CAVPO		Retinol-binding protein 3	RBP3	Cavia porcellus (Guinea pig)	19	
<input type="checkbox"/>	P12665	RET3_CRISP		Retinol-binding protein 3	RBP3	Cricetidae sp. (Hamster)	15	
<input type="checkbox"/>	Q9R0I0	Q9R0I0_RAT		Interphotoreceptor retinoid binding...	Rbp3 irbp	Rattus norvegicus (Rat)	261	
<input type="checkbox"/>	O57689	O57689_DANRE		Interphotoreceptor retinoid-binding...	irbp IRBP	Danio rerio (Zebrafish) (Brachydanio rerio)	628	