

PRACTICAL 2

Question 1. How many “Homo sapiens” sequences deposited in DDBJ? Compare with Genbank and EMBL

<http://www.ddbj.nig.ac.jp/>

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

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

Solution.

DDBJ	Genbank	EMBL
5,711,070	28,472,631	40,620,497

Organism

Search

- ☐ Severe acute respiratory syndrome coronavirus 2 13753745
- ☐ Homo sapiens 5711070
- ☐ Mus musculus 3238163
- ☐ soil metagenome 1052953
- ☐ human gut metagenome 972410
- ☐ metagenome 727997
- ☐ Escherichia coli 500017

Results by taxon

Top Organisms [\[Tree\]](#)

Homo sapiens (28472631)
Severe acute respiratory syndrome-related coronavirus (8723757)
Escherichia coli (4569390)
Klebsiella pneumoniae (2334738)
Acinetobacter baumannii (1596327)
Pseudomonas aeruginosa (1377091)
Staphylococcus aureus (1025654)
Human immunodeficiency virus 1 (846692)
uncultured bacterium (777777)
Mycobacterium tuberculosis (636949)
Salmonella enterica (593067)
Enterococcus faecium (580056)

Text Search

Uses EBI Search to perform a free text search across

Search term: "Homo sapiens"

Search results for "Homo sapiens" No result

Refer to

- Assembly
 - Assembly (572,928)
- Sequence
 - Sequence (40,620,497)
 - Sequence (CON) (565,646)
 - Sequence (Standard) (40,054,851)

(A) DDBJ

(B) Genbank

(C) EMBL

FIGURE 1. Number of “Homo sapiens” sequences deposited in each databases

Question 2. What is the GC-content of the AY330867?

Solution. Firstly I searched for the sequence, whose accession number is AY330867, on the **ENA (European Nucleotide Archive)**. On the right portion of the search, there is an option to download the **FASTA file**, which I downloaded. Later, I uploaded it to the **Seq2Feature** tool to compute the nucleotide content. The GC-content of **AY330867** is **46.527778%**

ENA
European Nucleotide Archive

Home Submit Search Rulespace About Support

Enter text search terms Search

Examples: Histone, B1000005

AY330867 View

Examples: Taxon:2606, B1000005, PRJEB402

Sequence: AY330867.1

Synthetic construct human lysozyme mRNA, complete cds.

Organism: synthetic construct

Mol Type: mRNA

Topology: linear

Base Count: 432

Dataclass: STD

Tax Division: SYN

Accession: AY330867

View: EMBL FASTA

Download: EMBL FASTA

Navigation: Show

Sequence Versions: View

FIGURE 2. Text search on the ENA browser

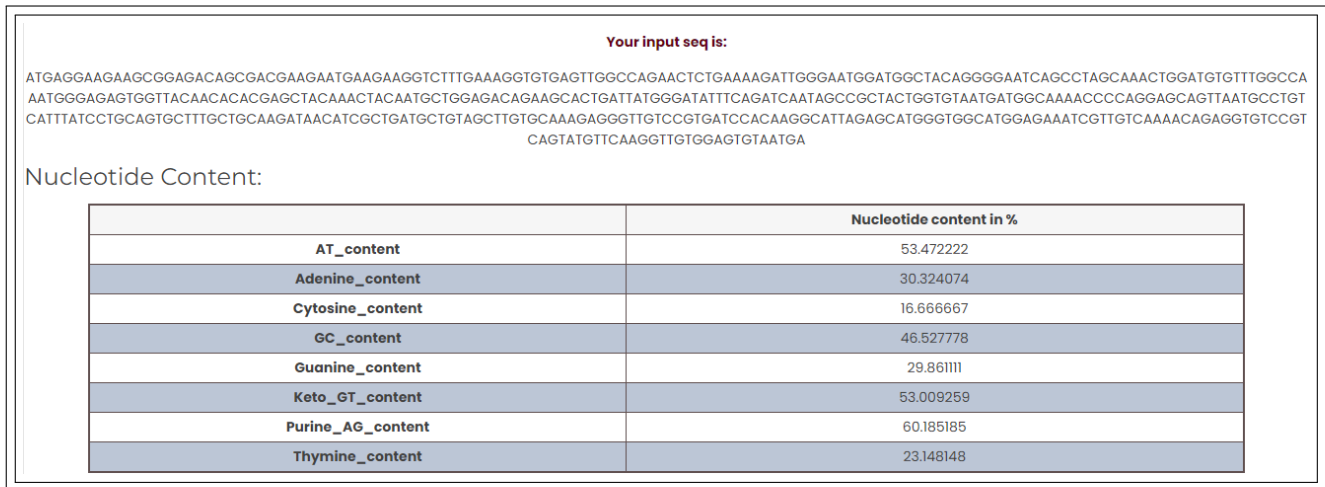


FIGURE 3. Nucleotide content on the Seq2Feature tool

Question 3. Compare the contents in DDBJ, Genbank and EMBL

Solution. The contents in DDBJ database are as follows:

- Locus
- Name
- Source
- Accession number
- Version number
- Keywords
- Authors
- Reference
- Journal
- PUBMED index
- Features that include exon, variation, gaps, mRNA
- Nucleotide sequence
- Frequency of bases (A,T,C,G)
- Protein sequence (translated)

```

LOCUS       HUM1621FD                470 bp    DNA    linear    HUM 13-AUG-1993
DEFINITION  Homo sapiens DNA sequence.
ACCESSION   U22438
VERSION     L22438.1
KEYWORDS    repeat polymorphism; single strand sequencing.
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
            Catarrhini; Hominidae; Homo.
  REFERENCE 1 (bases 1 to 470)
  AUTHORS   Weber, J.L.
  JOURNAL   Unpublished
  COMMENT   Original source text: Homo sapiens male blood DNA.
  FEATURES   Location/Qualifiers
             source          1..470
                        /organism="Homo sapiens"
                        /mol_type="genomic DNA"
                        /db_xref="taxon:9606"
                        /sex="male"
                        /tissue_type="blood"
BASE COUNT  142 a                82 c                85 g                161 t
ORIGIN
1 tggcttcttc ccccgaaaa ttaaaagac acagactgag aaactgaatg tgggagaaag
61 cagtgagatta tgctgtcttc agggctacc tcaggttgg agctattctt tcaaatcac
121 cactctcagg cactctcagg agagagagag tggatcttaa cttttttctt tgcagcatt
181 tggtaggggt gttttatga ccaaatatgt tccacacacc tagttttttg tgactgacta
241 aatatatagg agtttaaat tttttctcca aaatttatag aaattttttg ttgaagaatg
301 acctcatatc tatgtcttgc tttttctctt cttattaaat acattgtctg atgaagaac
361 aaggttcag aatgaagatt atttgtcttc ttgctatgga atttatatac ggctactctt
421 ctttaggtc gtatcaagg acatatctttt acctattaaa aggaagatc
//

```

FIGURE 4. DDBJ

The contents in Genbank database are as follows:

- Locus
- Definition
- Source
- Accession number
- Version number
- Keywords
- Title
- Authors
- Reference
- Journal
- PMID
- Features that include exon, variation, gaps, mRNA
- CDs (Coding sequence)
- Frequency of bases (A,T,C,G)
- Protein sequence (translated)

```

LOCUS       KJ3946236                441 bp    DNA    linear    PRI 16-SEP-2014
DEFINITION  Homo sapiens Kidd blood group protein (SLC14A1) gene, exons 4, 5
            and partial cds.
ACCESSION   KJ3946236
VERSION     KJ3946236.1
KEYWORDS    .
SOURCE      Homo sapiens (human)
  ORGANISM  Homo sapiens
            Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi;
            Mammalia; Eutheria; Euarchontoglires; Primates; Haplorrhini;
            Catarrhini; Hominidae; Homo.
  REFERENCE 1 (bases 1 to 441)
  AUTHORS   Zhang, A. and Chi, Q.
  TITLE     Direct Submission
  JOURNAL   Submitted (07-JUN-2014) Laboratory Department, Fujian Provincial
            Blood Center, No. 28 Xihuan South Road, Fuzhou, Fujian 350004,
            China
  COMMENT   ##Assembly-Data=START##
            Sequencing Technology :: Sanger dideoxy sequencing
            ##Assembly-Data=END##
  FEATURES   Location/Qualifiers
             gap            152..251
                        /estimated_length=unknown
             exon            252..341
                        /gene="SLC14A1"
                        /number=5
             variation       320
                        /gene="SLC14A1"
                        /replace="a"
ORIGIN
1 atggagggaca gccccactat ggttagagtg gacagcccca ctatggttag aggtgaaac
61 caggtttgga cagtcacagg ggaagagtg tttcccaagg cttttgacta tgtcacggt
121 gacatgaaa aactgtccaa ccagctaaa g
[gap 100 bp] Expand_Hs

```

FIGURE 5. Genbank

The contents in EMBL database are as follows:

- Identification line
- Source
- Accession number
- Sequence version
- Date of first entry
- Description
- Keywords
- Organism species
- Organism classification
- References
- Author
- CDs (Coding sequence)
- Repeat regions
- Frequency of bases (A,T,C,G)
- Protein sequence (translated)

```

ID X2304236; SV 1; linear; genomic DNA; STD; HMM; 442 BP.
XX
AC X2304236;
XX
DT 18-SEP-2014 (Ref. 122, Created)
DT 18-SEP-2014 (Ref. 122, Last updated, Version 1)
XX
DE Hm sapiens K14d blood group protein (SLC14A4) gene, exon 4, 5 and
DE partial cds.
XX
KW
XX
KW
XX
KW sapiens (human)
XX
CC Eukaryota; Metazoa; Chordata; Craniata; Vertebrata; Euteleostomi; Mammalia;
CC Eutheria; Lauriarchontia; Primates; Euarchontii; Carnivora;
CC Homo.
XX
RN [1]
RP 1-442
RA Zhang A., Chi Q.;
RT ;
RL Submitted (07-JUN-2014) to the INSDC.
RL Humanoid Department, Fujian Provincial Blood Center, No. 28 Xinhua South
RL Road, Fuzhou, Fujian 350001, China
XX
DR M5; J060737654/F068011203/0680457.
XX
XX
XX ##Assembly-Data-START##
CC Sequencing Technology : Sanger dideoxy sequencing
XX ##Assembly-Data-END##
XX
FH key Location/Qualifiers
XX
FT source 1. 441

```

FIGURE 6. EMBL

Question 4. Get the papers about **discrimination of beta-barrel membrane proteins**. <https://pubmed.ncbi.nlm.nih.gov/>

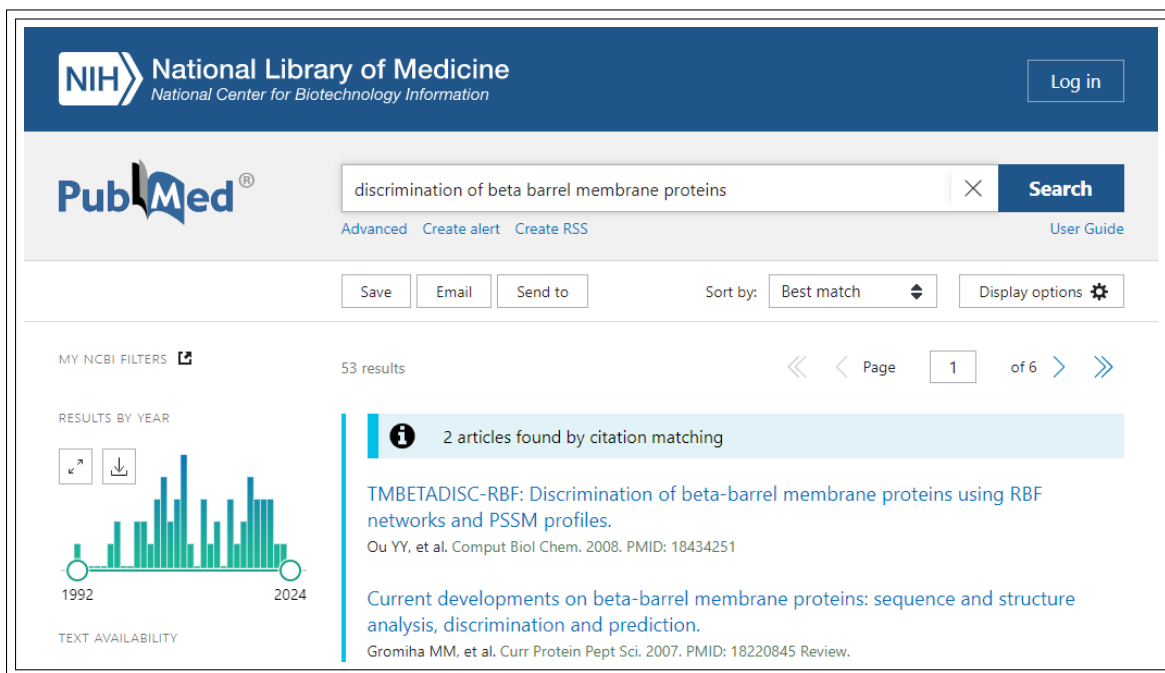


FIGURE 7. PUBMED papers about "discrimination of beta barrel membrane proteins

Solution. The list of papers is:

- TMBETADISC-RBF: Discrimination of beta-barrel membrane proteins using RBF networks and PSSM profiles
 - Current developments on beta-barrel membrane proteins: sequence and structure analysis, discrimination and prediction
 - BetAware-Deep: An Accurate Web Server for Discrimination and Topology Prediction of Prokaryotic Transmembrane beta-barrel Proteins
 - Current developments on beta-barrel membrane proteins: sequence and structure analysis, discrimination and prediction.
- and many more**

Question 5. Find the papers published by any author (E.g.) Kihara D

Solution. The author chosen is **Daisuke Kihara**

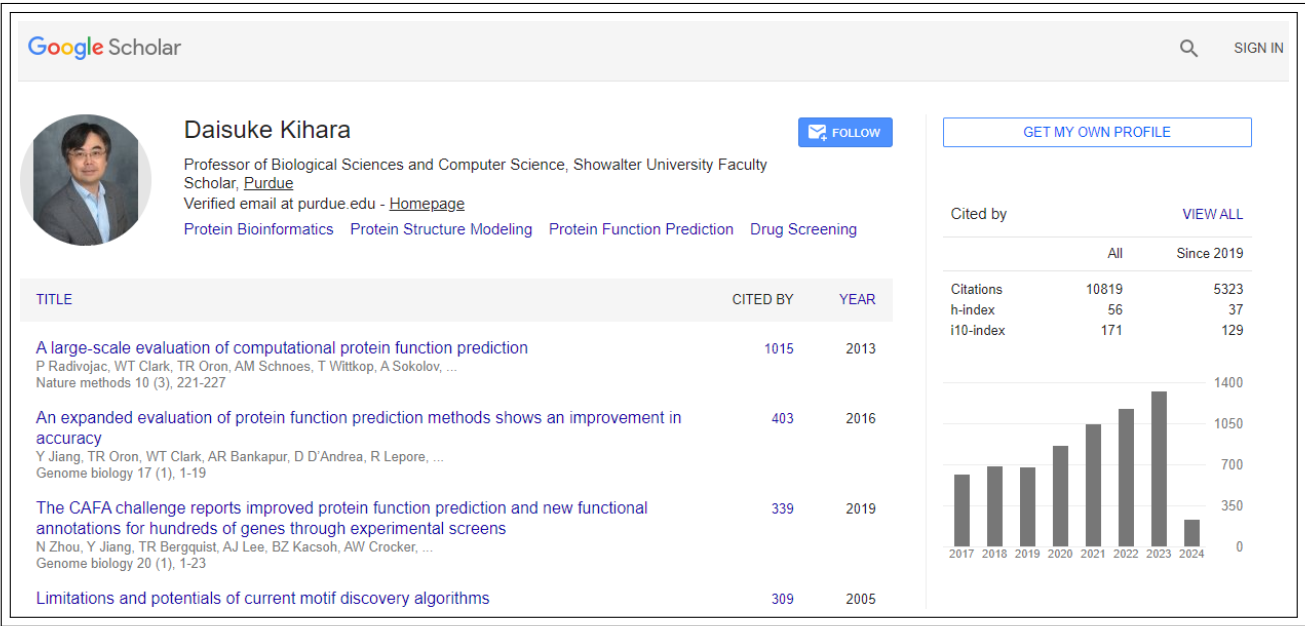


FIGURE 8. Google scholar search for Daisuke Kihara

The list of his papers is:

- A large-scale evaluation of computational protein function prediction
- An expanded evaluation of protein function prediction methods shows an improvement in accuracy
- Limitations and potentials of current motif discovery algorithms
- Protein-protein docking using region-based 3D Zernike descriptors
- and many more

Question 6. How many related articles are listed in PUBMED for the paper, **Cell 2008 Dec 26;135(7):1158-9?**



FIGURE 9. Similar articles to Ceel 2008 Dec26;135(7):1158-9?

Solution. The PMID number is 19109882. There are 27 similar articles. Some of them have been listed below:

- Evolution. Tinkering inside the organelle.
 - Transport proteins (carriers) of mitochondria.
 - Systematic analysis of the twin cx(9)c protein family.
 - Mitochondrial permeability transition pore opening as a promising therapeutic target in cardiac diseases.
- and many more**

Question 7. List the papers published in the journal “**Nature**” for the year 2024. Check the list in **SCOPUS** and **PUBMED**

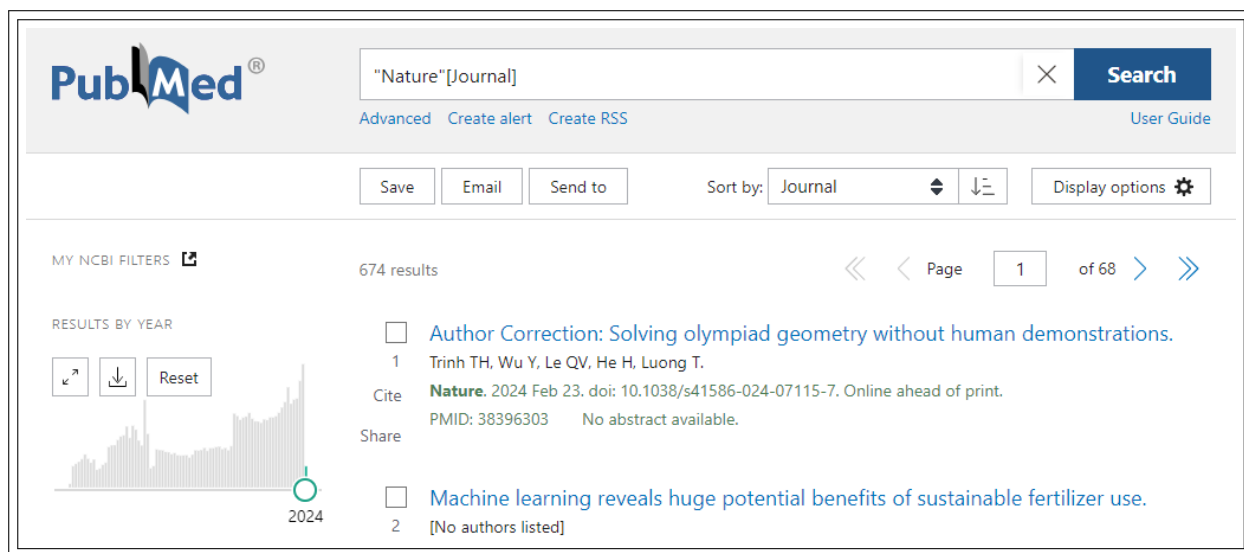


FIGURE 10. Nature publications on PUBMED

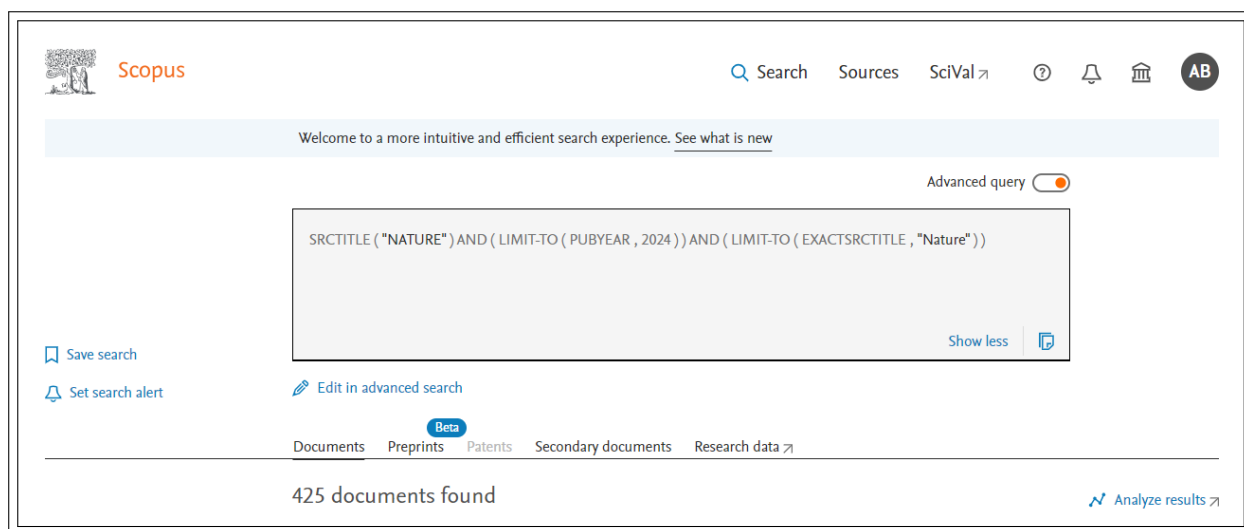


FIGURE 11. Nature publications on SCOPUS

Advanced query search

PUBMED has **674** publications
SCOPUS has **425** publications

Question 8. Find the h-index and number of citations for “Burkhard Rost”

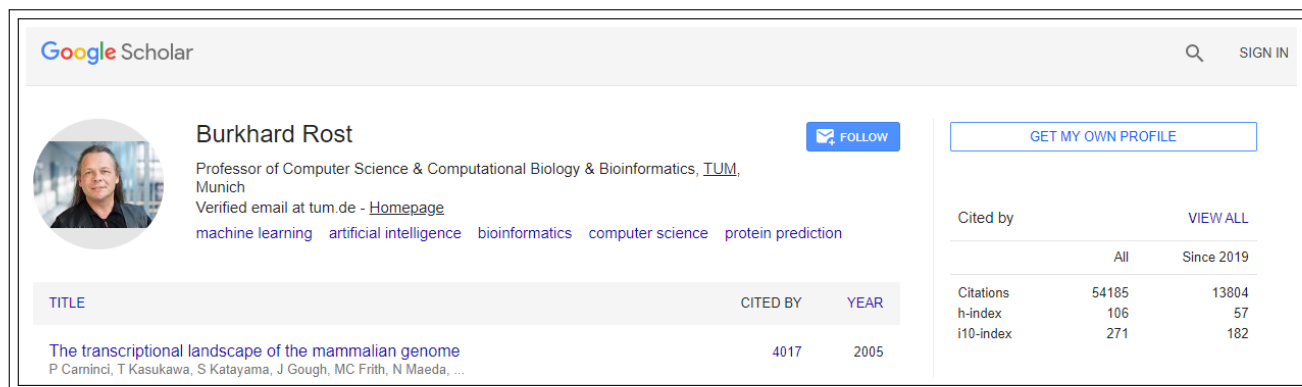


FIGURE 12. Google Scholar search for Burkhard Rost

Solution. The h-index of Burkhard Rost is **106**

The total number of citations is **54185**

Question 9. Find the class of the enzyme EC 1.7.2.3 and its function

<http://www.brenda-enzymes.org/>

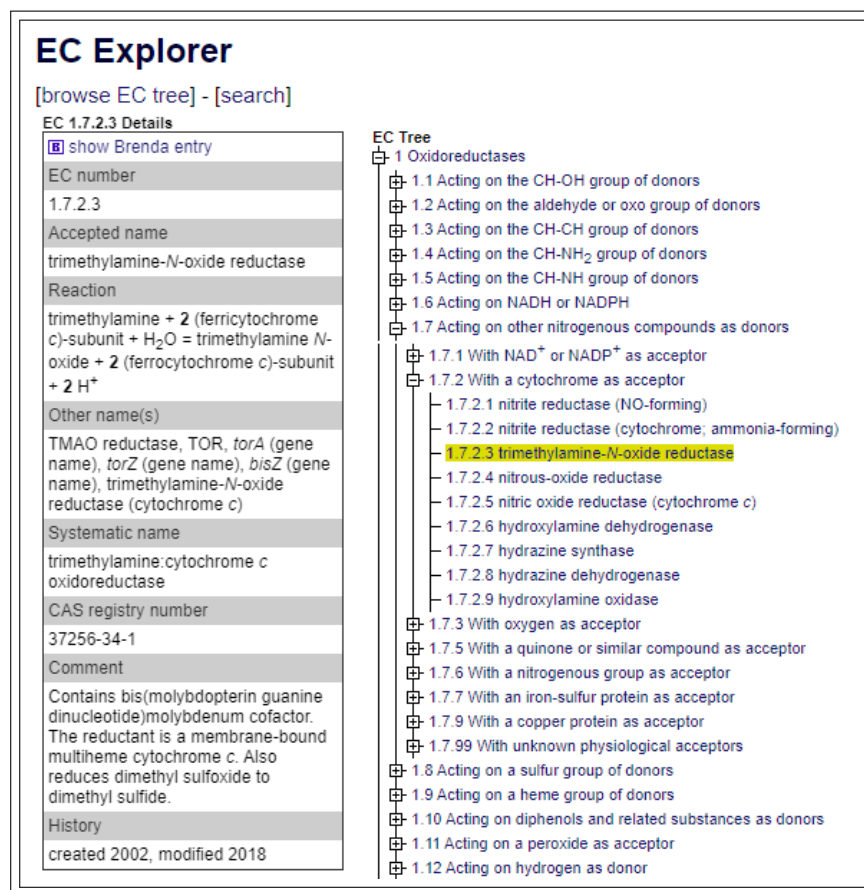
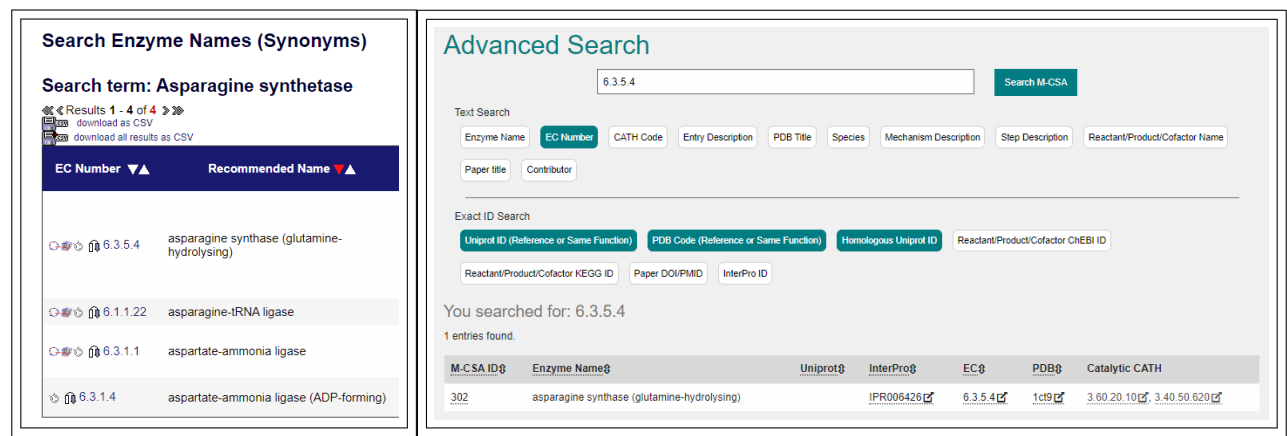


FIGURE 13. BRENDA search for enzyme EC 1.7.2.3

Solution. EC 1.7.2.3 is the EC number for the enzyme **trimethylamine-N-oxide reductase**. It is an **Oxidoreductase**. It contains a bis (molybdopterin guanine dinucleotide) molybdenum cofactor. The reductant is a membrane-bound multiheme cytochrome c. It reduces dimethyl sulfoxide to dimethyl sulfide.

Question 10. Find the catalytic site residues in **Asparagine synthetase**. *Hint:* Find the EC number and search in Catalytic site atlas <https://www.ebi.ac.uk/thornton-srv/m-csa/>



(A) BRENDA enzyme search (B) Catalytic site atlas search

FIGURE 14. Seq2Feature tool for nucleotide content

Solution. First I searched for the enzyme on the BRENDA to get its EC number. The EC number of **Asparagine synthetase** is **6.3.5.4**. On searching the EC number on the catalytic site atlas, in the mechanism sector, I found the catalytic site residues. They are

- Cys2 (N-term)
- Leu51 (main-C)
- Thr322, Arg325
- Cys2
- Gly76 (main-N), Asn75

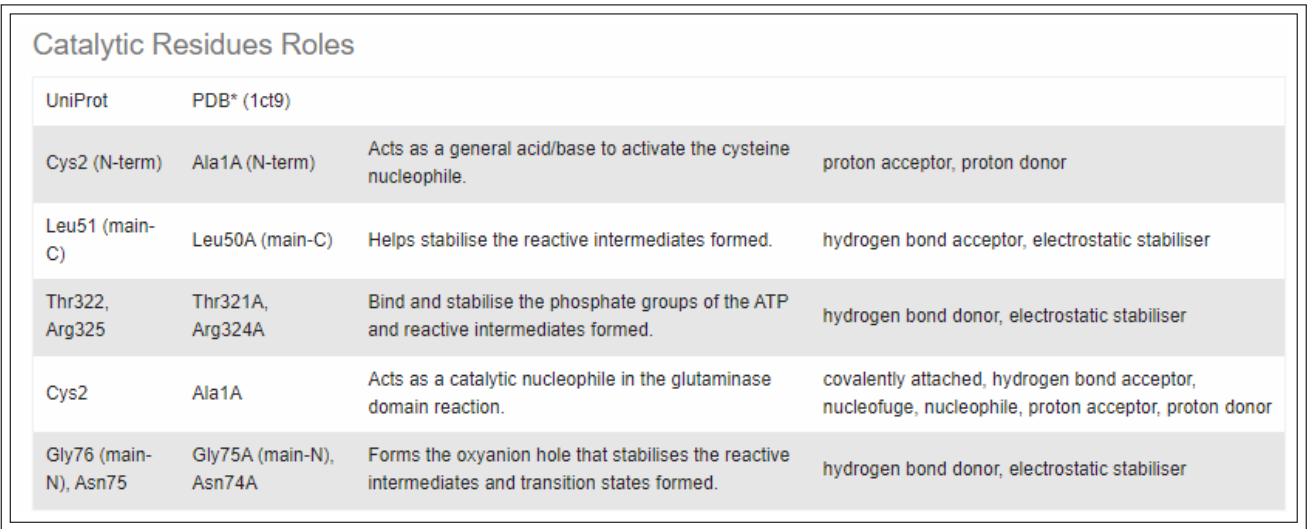


FIGURE 15. Catalytic site atlas residue for EC 6.3.5.4

Question 11. Find the scientific name, the taxonomy ID, and the number of chromosomes for the following organisms:
<http://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?mode=Root>
Organisms: **Human, Cat, Dog, domestic guinea pig, and Thale cress**

Solution. Below, I have tabulated each of the properties for each of the given organisms.

Human

Scientific Name: *Homo sapiens*
Taxonomy ID: 9606
Number of chromosomes: 46

Cat

Scientific Name: *Felis catus*
Taxonomy ID: 9685
Number of chromosomes: 38

Dog

Scientific Name: *Canis lupus familiaris*
Taxonomy ID: 9615
Number of chromosomes: 78

Domestic guinea pig

Scientific Name: *Cavia porcellus*
Taxonomy ID: 10141
Number of chromosomes: 64

Thale cress

Scientific Name: *Arabidopsis thaliana*
Taxonomy ID: 3702
Number of chromosomes: 10

The properties have been computed by the **Taxonomy browser**. The number of chromosomes could be extracted from any other database too.

[Homo sapiens](#)

Taxonomy ID: 9606 (for references in articles please use NCBI:txid9606)

current name

Homo sapiens Linnaeus, 1758

Genbank common name: **human**

FIGURE 16. Taxonomy Browser: Human

Felis catus

Taxonomy ID: 9685 (for references in articles please use NCBI:txid9685)

current name

Felis catus Linnaeus, 1758

homotypic synonym: ***Felis silvestris catus***

includes: ***Korat cats*** L.

Genbank common name: **domestic cat**

FIGURE 17. Taxonomy Browser: Cat

Canis lupus familiaris

Taxonomy ID: 9615 (for references in articles please use NCBI:txid9615)

current name

Canis lupus familiaris Linnaeus, 1758

homotypic synonym: Canis familiaris Linnaeus, 1758

includes: ☐ beagle dog

Genbank common name: dog

FIGURE 18. Taxonomy Browser: Dog

Cavia porcellus

Taxonomy ID: 10141 (for references in articles please use NCBI:txid10141)

current name

Cavia porcellus

basionym: Mus porcellus Linnaeus, 1758

homotypic synonym: Cavia aperea porcellus

Genbank common name: domestic guinea pig

FIGURE 19. Taxonomy Browser: domestic guinea pig

Arabidopsis thaliana

Taxonomy ID: 3702 (for references in articles please use NCBI:txid3702)

current name

Arabidopsis thaliana (L.) Heynh., 1842

basionym: Arabis thaliana L., 1753

Genbank common name: thale cress

FIGURE 20. Taxonomy Browser: Thale cress

Question 12. What are NCBI E-utilities? Give the syntax for fetching a record in FASTA format using E-utilities.

<http://www.ncbi.nlm.nih.gov/books/NBK25500/>

Solution. NCBI E-utilities are a set of tools that allow you to access and manipulate data from the **NCBI databases**, such as **PubMed**, **PMC**, **Gene**, **Nuccore** and **Protein**. You can use E-utilities to perform tasks such as searching, linking, downloading, and converting data formats. E-utilities use a fixed URL syntax for their operations, and can be called from a web browser, a command line, or a custom program. Syntax for fetching a record in FASTA format:

Fetch a record in FASTA format

```

1 efetch.fcgi?db=<database>&id=<uid_list>&rettype=fasta&retmode=<
  retrieval_mode>
2 # Entrez database (&db)
3 # List of UIDs (&id)
4 # Retrieval type (&rettype), [here fasta]
5 # Retrieval mode (&retmode)
6
7 # Example
8 efetch.fcgi?db=nuccore&id=34577062,24475906&rettype=fasta&retmode=text
9 # Database is nuccore
10 # UIDs are 34577062 and 24475906
11 # rettype is fasta
12 # retmode is text

```

Question 13. List two databases under each of the following category.

- Protein properties
- Small molecules (Structure related)
- Cancer gene databases

Hint: Use Nucleic Acids Research (NAR) – ‘database category list’

<https://www.oxfordjournals.org/our-journals/nar/database/c/>

Solution. Below are the databases available in each of the categories specified in the question:

Protein properties	Small molecules	Cancer gene databases
AAindex BindingDB Cybase dbPTM eProS HHMD iProLINK MALISAM MegaMotifbase MemMoRF Minimotif miner MobiDB MP:PD PFD - Protein Folding Database PIDD PINT PPD PPT-DB Proteome-pI 2.0 ProTherm PRTAD REFOLD SwissSidechain TOPPR	BitterDB ChEBI - Chemical Entities of Biological ChemBank ChemDB Crystallography Open Database CSD - Cambridge Structural Database EDULISS (currently unavailable) FragmentStore Hemolytik Het-PDB Navi HIC-Up MMsINC mVOC PDB-Ligand PubChem R.E.D.D.B. SCRIPDB SuperDrug SuperScent SuperToxic	ArrayMap BCCTBbp BreCAN-DB Cancer RNA-Seq Nexus Cancer3D CancerGenes CancerPPD Candidate Cancer Gene Database CanGEM CanSAR CaSNP cBioPortal CCDB ccmGDB CellLineNavigator ChimerDB CIViCdb ClinVar CMPD Colorectal Cancer Atlas COLT-Cancer COSMIC CTDatabase dbDEMC

(A) Protein properties

(B) Small molecules

(C) Cancer gene

FIGURE 21. Database Category Lists

Protein Properties

AAindex
BindingDB
Cybase

Small molecules (structure related)

BitterDB
ChemBank
Crystallography Open Database

Cancer gene databases

BCCTBbp
BreCAN-DB
Cancer RNA-Seq Nexus