

# **Secondary Nucleotide Databases**

**Shatakshi Kulkarni**

# Introduction

- ❑ Secondary nucleotide databases are databases that curate and organize nucleotide sequences from primary databases, such as GenBank, EMBL, and DDBJ.
- ❑ They provide added value through annotation, analysis, and the integration of additional information.
- ❑ These databases often offer user-friendly interfaces and tools for sequence alignment, similarity searches, and functional predictions.
- ❑ These databases often focus on specific types of sequences, specific organisms, or specific types of analysis.
- ❑ E.g., Annotations and functional information about the identified proteins is given.

# Secondary Nucleotide Databases

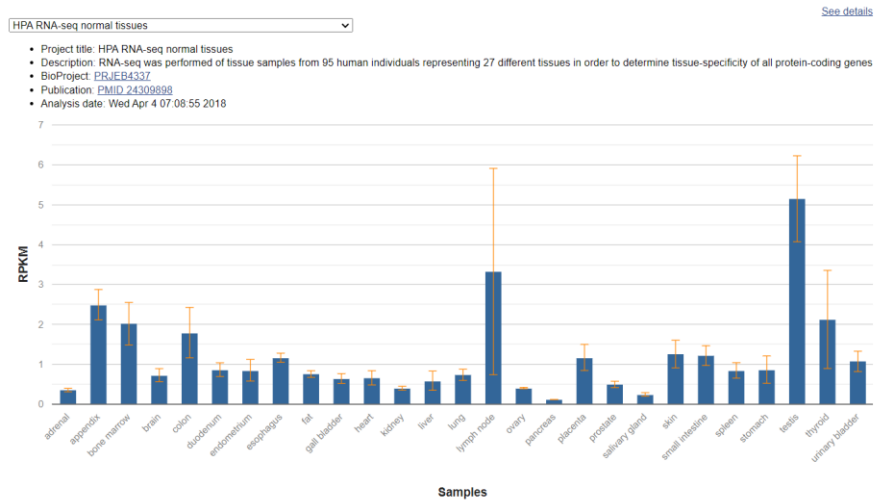
- ❑ **Gene**: concentrate on individual genes and their related data, including sequences, structure, function, and expression.
- ❑ **Genome**: Genome secondary databases focus on entire genomes, providing comprehensive resources for the analysis and comparison of whole-genome sequences.
- ❑ **EST**: Expressed sequence tags. Maintains expressed sequence tags (ESTs) and short, single-pass reads from mRNA (cDNA)
- ❑ **STS**: Sequenced Tagged Site. Is a relatively short, easily PCR-amplified sequence (200 to 500 bp) which can be specifically amplified by PCR and detected in the presence of all other genomic sequences and whose location in the genome is mapped.
- ❑ **GSS**: Genome Survey Sequence. The GSS division of GenBank is similar to the EST division, with the exception that most of the sequences are genomic in origin, rather than cDNA (mRNA).

# Gene

- ❑ <https://www.ncbi.nlm.nih.gov/gene/?term=>
- ❑ The Gene database is a resource of the National Center for Biotechnology Information (NCBI) that centralizes gene-related information into individual records.
- ❑ The NCBI Gene database has information on gene sequences, gene alleles and mutations, genomes, amino acid sequences for proteins, and much more genetic data on humans, as well as many other animal species.
- ❑ RefSeq and other primary databases provide data to Gene database. They integrate, annotate, and often reanalyze data from primary databases like GenBank, EMBL, and DDBJ to provide additional layers of information, context, and utility for researchers.
- ❑ RefSeq provides a comprehensive, integrated, non-redundant set of sequences, including genomic DNA, transcripts, and proteins. It also provides sequences curated giving annotated functional elements, gene expression studies.

# Gene Database Results

- ☐ Summary: Pathophysiological role of gene, disease mechanism and molecular pathways
- ☐ Expression: gene expression of gene in different tissues



Sample	BioSample	RPKM <a href="#">?</a>	Count <a href="#">?</a>	Links
<a href="#">+</a> adrenal	3 samples	0.355 ± 0.042	29579	<a href="#">SRA</a> , <a href="#">BioSample</a>
<a href="#">+</a> appendix	3 samples	2.495 ± 0.381	181185	<a href="#">SRA</a> , <a href="#">BioSample</a>
<a href="#">+</a> bone marrow	4 samples	2.018 ± 0.532	432170	<a href="#">SRA</a> , <a href="#">BioSample</a>
<a href="#">+</a> brain	3 samples	0.73 ± 0.163	72250	<a href="#">SRA</a> , <a href="#">BioSample</a>
<a href="#">+</a> colon	5 samples	1.794 ± 0.632	396724	<a href="#">SRA</a> , <a href="#">BioSample</a>
<a href="#">+</a> duodenum	2 samples	0.867 ± 0.173	44311	<a href="#">SRA</a> , <a href="#">BioSample</a>
<a href="#">+</a> endometrium	3 samples	0.852 ± 0.271	89581	<a href="#">SRA</a> , <a href="#">BioSample</a>
<a href="#">+</a> esophagus	3 samples	1.166 ± 0.113	160671	<a href="#">SRA</a> , <a href="#">BioSample</a>

# Gene Database Results

- ❑ Orthologs: comparison of searched gene in different species with possible common ancestry. For each species, all transcript variants and the transcripts coding for a protein are provided. No. of Amino acids, Architecture (1 transcript containing different exons part (amino acids) coding for different proteins).

SEARCH THE TAXONOMY TREE

Enter taxonomic name

vertebrates  
birds  
turtles  
alligators and others  
lizards & snakes  
mammals  
amphibians  
lampreys  
cartilaginous fishes

0 selected

					Previous	Next
Species	Gene	Architecture	aa			
<input type="checkbox"/> <i>Homo sapiens</i> human	BRCA1 BRCA1 DNA repair associated		1,863	▼		
<input type="checkbox"/> <i>Mus musculus</i> house mouse	Brca1 breast cancer 1, early onset		1,812	▼		
<input type="checkbox"/> <i>Rattus norvegicus</i> Norway rat	Brca1 BRCA1, DNA repair associated		1,817	▼		
<input type="checkbox"/> <i>Canis lupus familiaris</i> dog	BRCA1 BRCA1 DNA repair associated		1,934	▼		
<input type="checkbox"/> <i>Gallus gallus</i> chicken	BRCA1 BRCA1 DNA repair associated		1,750	▼		
<input type="checkbox"/> <i>Bos taurus</i> cattle	BRCA1 BRCA1 DNA repair associated		1,849	▼		

☐

*Homo sapiens*  
human

BRCA1  
BRCA1 DNA repair associated

1,863



RefSeq transcripts (368)	RefSeq proteins (368)	Architecture	aa
NM_007294.4	NP_009225.1		1,863
NM_007297.4	NP_009228.2		1,816
NM_007298.4	NP_009229.2		759
NM_007299.4	NP_009230.2		699
NM_007300.4	NP_009231.2		1,884
NM_001407571.1	NP_001394500.1		1,792
NM_001407581.1	NP_001394510.1		1,885
NM_001407582.1	NP_001394511.1		1,885
NM_001407583.1	NP_001394512.1		1,884

# Gene Database Results

- ❑ Gene Table: Provides access to ortholog sets and transcripts table
- ❑ Transcript Table: For the particular species that is searched, it provides information about how many different mRNA molecules can be transcribed from that particular gene, each transcript variant eventually leading to distinct protein isoforms. The variants arise due to alternative splicing, alternative promoter usage, post-transcriptional modification, etc.

1 Gene									
Download ▾		Edit genes ▾		Select columns		Rows per page 20 ▾		1-1 of 1 < >	
<input type="checkbox"/>	Gene ID	Symbol	Gene name	Scientific name	Common Name	Gene type	Transcripts	Ortholog set	Input
<input type="checkbox"/>	672	BRCA1	BRCA1 DNA repair associated	Homo sapiens (human)	human	protein-coding	368	672	672

## Transcripts and Proteins

BRCA1 – BRCA1 DNA repair associated

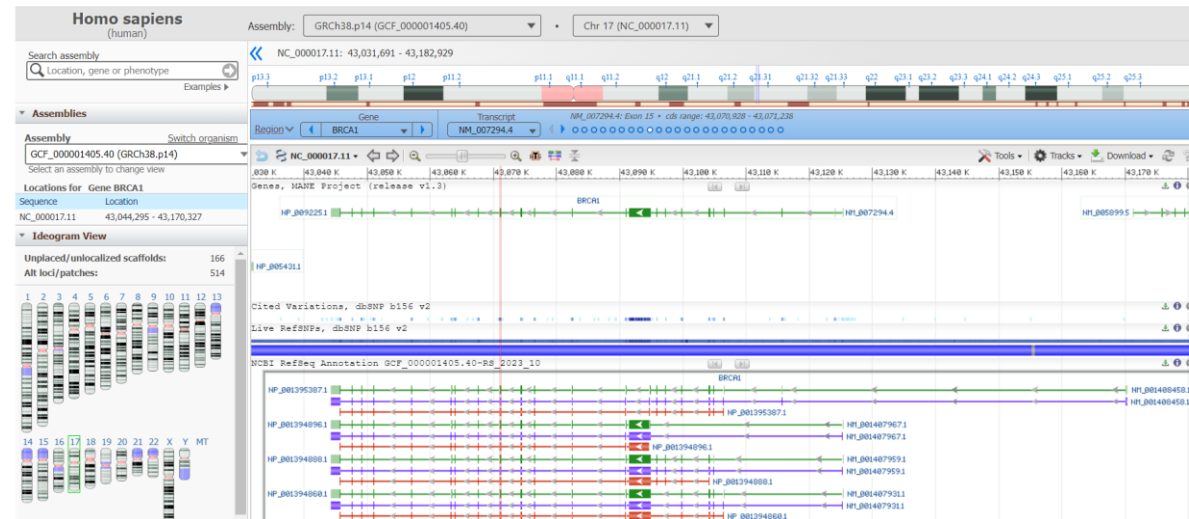
Homo sapiens (human)

368 Transcripts 368 selected							
Download ▾		Select columns					
Gene ID	Gene symbol	Transcript	Length (nt)	Protein	Length (aa)	Protein name	Isoform
672	BRCA1	NM_001408458.1	3785	NP_001395387.1	712	breast cancer type 1 susceptibil...	116
672	BRCA1	NM_001407967.1	6390	NP_001394896.1	1566	breast cancer type 1 susceptibil...	85
672	BRCA1	NM_001407959.1	6851	NP_001394888.1	1736	breast cancer type 1 susceptibil...	78
672	BRCA1	NM_001407931.1	6803	NP_001394860.1	1774	breast cancer type 1 susceptibil...	59
672	BRCA1	NM_001407747.1	6926	NP_001394676.1	1815	breast cancer type 1 susceptibil...	37
672	BRCA1	NM_001407962.1	6902	NP_001394891.1	1735	breast cancer type 1 susceptibil...	80
672	BRCA1	NM_001408512.1	3473	NP_001395441.1	592	breast cancer type 1 susceptibil...	145

# Gene Database Results

- ❑ Genomic Context: Chromosome location, reference assembly, Genomic regions, transcripts, and products

Genome Data Viewer



- ❑ Bibliographic Data: Related research articles
- ❑ Phenotypes: Associated disease conditions provides information on diseases in which the gene is related

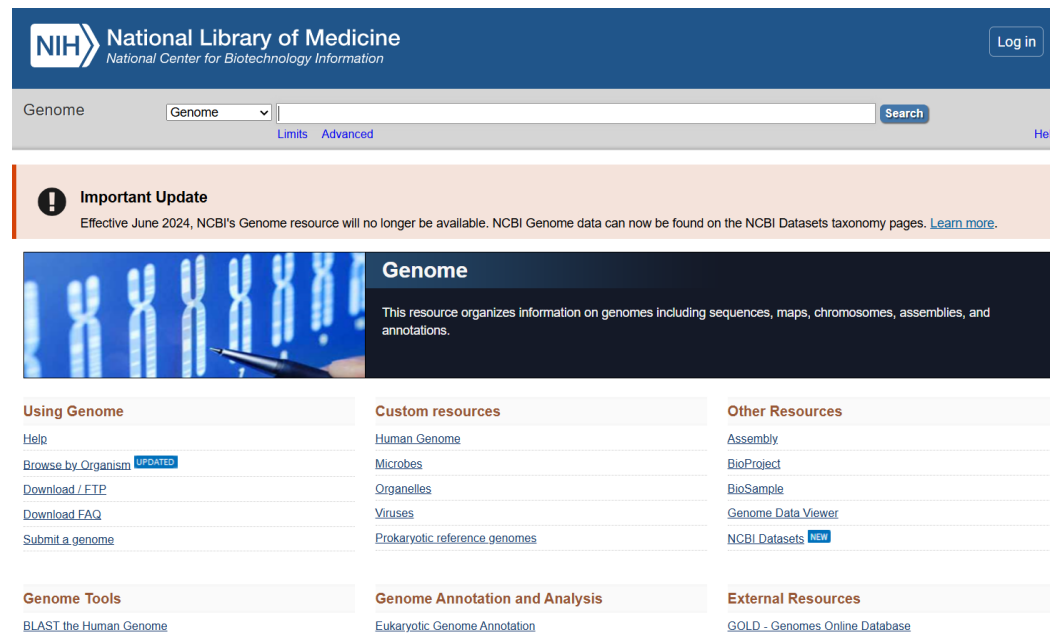


# Gene Database Results

- ❑ Variations: ClinVar, dbVar
- ❑ Pathways: Cellular and molecular pathways in which they are involved. E.g., cell cycle, cell proliferation, reproduction
- ❑ Interactions: Interaction of searched gene with other genes
- ❑ General Gene Information:
  - **Gene Ontology**: info about pathways (molecular, cellular) pathway in which it is involved. E.g., DNA repair, checkpoint in cell cycle
- ❑ General Protein Information: info related to proteins, enzymes, activities of that particular gene
- ❑ NCBI RefSeq: information of each protein coded by transcript variants of that particular gene

# Genome

- ❑ [Home - Genome - NCBI \(nih.gov\)](https://www.ncbi.nlm.nih.gov/genome/)
- ❑ Genome is a Secondary nucleotide databases of NCBI
- ❑ Whole genome sequence of an organism can be obtained



The screenshot shows the NCBI Genome homepage. At the top is the NIH logo and the text "National Library of Medicine National Center for Biotechnology Information". A search bar is present with a dropdown menu set to "Genome" and a "Search" button. Below the search bar is an "Important Update" banner stating that the NCBI Genome resource will no longer be available as of June 2024, with a link to "Learn more". The main content area features a large image of chromosomes and a section titled "Genome" with a description: "This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and annotations." Below this, there are three columns of links. The first column, "Using Genome", includes links for Help, Browse by Organism (marked as UPDATED), Download / FTP, Download FAQ, and Submit a genome. The second column, "Custom resources", includes links for Human Genome, Microbes, Organelles, Viruses, and Prokaryotic reference genomes. The third column, "Other Resources", includes links for Assembly, BioProject, BioSample, Genome Data Viewer, and NCBI Datasets (marked as NEW). At the bottom, there are three more sections: "Genome Tools" with a link to BLAST the Human Genome, "Genome Annotation and Analysis" with a link to Eukaryotic Genome Annotation, and "External Resources" with a link to GOLD - Genomes Online Database.

NIH National Library of Medicine  
National Center for Biotechnology Information

Genome Genome Search Limits Advanced Help

**Important Update**  
Effective June 2024, NCBI's Genome resource will no longer be available. NCBI Genome data can now be found on the NCBI Datasets taxonomy pages. [Learn more.](#)

**Genome**  
This resource organizes information on genomes including sequences, maps, chromosomes, assemblies, and annotations.

**Using Genome**  
[Help](#)  
[Browse by Organism](#) **UPDATED**  
[Download / FTP](#)  
[Download FAQ](#)  
[Submit a genome](#)

**Custom resources**  
[Human Genome](#)  
[Microbes](#)  
[Organelles](#)  
[Viruses](#)  
[Prokaryotic reference genomes](#)

**Other Resources**  
[Assembly](#)  
[BioProject](#)  
[BioSample](#)  
[Genome Data Viewer](#)  
[NCBI Datasets](#) **NEW**

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

**Genome Annotation and Analysis**  
[Eukaryotic Genome Annotation](#)

**External Resources**  
[GOLD - Genomes Online Database](#)

# Genome List

## Genome

Download a genome data package including genome, transcript and protein sequence, annotation and a data report

Selected taxa  
Escherichia coli (E. coli) Enter one or more taxonomic names

Filters

Download Select columns 278,189 Genomes Rows per page 20 1-20 of 278,189

<input type="checkbox"/> Assembly	GenBank	RefSeq	Scientific name	Modifier	Annotation	Action
<input type="checkbox"/> ASM584v2	GCA_000005845.2	GCF_000005845.2	Escherichia coli str. K-12 substr...	K-12 substr. MG165...	NCBI RefSeq Submitter	⋮
<input type="checkbox"/> ASM886v2	GCA_000008865.2	GCF_000008865.2	Escherichia coli O157:H7 str. S...	Sakai substr. RIMD ...	NCBI RefSeq Submitter	⋮
<input type="checkbox"/> ASM285371v1	GCA_002853715.1	GCF_002853715.1	Escherichia coli (E. coli)	14EC020 (strain)	NCBI RefSeq Submitter	⋮
<input type="checkbox"/> ASM1326v1	GCA_000013265.1	GCF_000013265.1	Escherichia coli UTI89	UTI89 (strain)	NCBI RefSeq Submitter	⋮

Level of assembly  
(complete, scaffold or contig)

Level	Release ...	WGS accession	Scaffolds count
Complete	Sep, 2013		1
Complete	Jun, 2018		3
Complete	Jan, 2018		3
Complete	Apr, 2006		2

# Genome Assembly

## Genome assembly ASM584v2 reference

<a href="#">Download</a>	<a href="#">datasets</a>	<a href="#">URL</a>	<a href="#">FTP</a>	Actions
NCBI RefSeq assembly	GCF_000005845.2			⋮
Submitted GenBank assembly	GCA_000005845.2			⋮
Taxon	<a href="#">Escherichia coli str. K-12 substr. MG1655</a>			
Strain	K-12 substr. MG1655			
Submitter	Univ. Wisconsin			
Date	Sep 26, 2013			
<a href="#">View the legacy Assembly page</a>				

### Assembly statistics

	RefSeq	GenBank
Genome size	4.6 Mb	4.6 Mb
Total ungapped length	4.6 Mb	4.6 Mb
Number of chromosomes	1	1
Number of scaffolds	1	1
Scaffold N50	4.6 Mb	4.6 Mb
Scaffold L50	1	1
Number of contigs	1	1
Contig N50	4.6 Mb	4.6 Mb
Contig L50	1	1
GC percent	51	51
Assembly level	Complete Genome	Complete Genome
View sequences	<a href="#">view RefSeq sequences</a>	<a href="#">view GenBank sequences</a>

### Annotation details

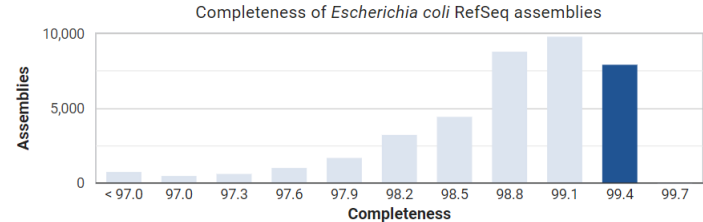
	RefSeq	GenBank
Provider	NCBI RefSeq	Univ. Wisconsin
Name	Annotation submitted by NCBI RefSeq	Annotation submitted by Univ. Wisconsin
Date	Mar 9, 2022	Nov 8, 2022
Genes	4,639	4,639
Protein-coding	4,288	4,288
	<a href="#">View RefSeq annotation</a>	<a href="#">View GenBank annotation</a>

### Quality analysis

#### CheckM analysis (v1.2.2)

Completeness: 99.48% (89th Percentile, dark blue bar)

Contamination: 0.15%



# Genome Annotation

## Genome Annotation

Genes annotated on [Escherichia coli str. K-12 substr. MG1655 ASM584v2 \(GCF\\_000005845.2\)](#)

Annotation Name: Annotation submitted by NCBI RefSeq (March 9, 2022)

Filters

Download

Select columns

4,639 Genes

Rows per page201-20 of 4,639

<input type="checkbox"/>	n	Name	Symbol	Locus Tag	Gene ID	Gene type	Proteins	Length (aa)	Action
<input type="checkbox"/>		thr operon lea...	thrL	b0001	944742	protein-coding	1 NP_414542.1	21	...
<input type="checkbox"/>		fused aspartat...	thrA	b0002	945803	protein-coding	1 NP_414543.1	820	...
<input type="checkbox"/>		homoserine ki...	thrB	b0003	947498	protein-coding	1 NP_414544.1	310	...
<input type="checkbox"/>		threonine synt...	thrC	b0004	945198	protein-coding	1 NP_414545.1	428	...
<input type="checkbox"/>		DIIE2502 dom	vaaY	b0005	944747	protein-coding	1 NP_414546.1	98	...

Gene

Gene

Advanced

Full Report

Send to:

thrL thr operon leader peptide [Escherichia coli str. K-12 substr. MG1655]

Gene ID: 944742, updated on 2-May-2024

Download Database

Summary

Gene symbol

thrL

Gene description

thr operon leader peptide

Primary source

ASAP ABE-0000006

Locus tag

b0001

See related

ECCOYC:EG11277

Gene type

protein coding

RefSeq status

PROVISIONAL

Organism

Escherichia coli str. K-12 substr. MG1655 (strain: K-12\_substrain\_MG1655)

Lineage

Bacteria; Pseudomonadota; Gammaproteobacteria; Enterobacterales; Enterobacteriaceae; Escherichia

Also known as

ECK0001

Summary

The ThrL leader peptide controls by attenuation the expression of the thrLABC operon, which encodes four out of the five enzymes of threonine biosynthesis pathway, in response to the threonine and isoleucine levels. [More information is available at EcoCyc: EG11277].

NEW

Try the new Gene table

Try the new Transcript table

Protein

Protein

Advanced

GenPept

thr operon leader peptide [Escherichia coli str. K-12 substr. MG1655]

NCBI Reference Sequence: NP\_414542.1

Identical ProteinsFASTAGraphics

Go to

LOCUSNP\_41454221 aa linearCON 09-MAR-2022

DEFINITIONthr operon leader peptide [Escherichia coli str. K-12 substr. MG1655].

ACCESSIONNP\_414542

VERSIONNP\_414542.1

DBLINKBioProject: PRJNA57779

BioSample: SAMN02684021

REFSEQ: accession NC\_000913.3

DBSOURCERefSeq

KEYWORDS

SOURCEEscherichia coli str. K-12 substr. MG1655

ORGANISM

Escherichia coli str. K-12 substr. MG1655

Bacteria; Pseudomonadota; Gammaproteobacteria; Enterobacterales; Enterobacteriaceae; Escherichia.

REFERENCE1 (residues 1 to 21)

AUTHORS

Riley,M., Abe,T., Arnaud,M.B., Berlyn,M.K., Blattner,F.R., Chaudhuri,R.R., Glasner,J.D., Horiiuchi,T., Kessler,I.M., Kosuge,T., Mori,H., Perna,M.T., Plunkett,G. III, Rudd,K.E., Serres,M.H., Thomas,G.H., Thomson,M.R., Wishart,D. and Wanner,B.L.

TITLE

Escherichia coli K-12: a cooperatively developed annotation snapshot--2005

JOURNAL

Nucleic Acids Res. 34 (1), 1-9 (2006)

# Genome Search Results

## ❑ Genomes List:

(Obtained from different strains, samples of same species and derived from different experimental methodologies)

Assembly, GenBank, RefSeq ID, Scientific name, Modifier (strain), Annotations (submitter), level(Indicates the completeness of the genome assembly [e.g., complete, scaffold, contig]), scaffold count (No. of scaffolds, gives an idea about gaps in the sequence and hence completeness)

## → Assembly:

### ➤ Reference genomes: (green tick)

High quality, well annotated genome assemblies that serve as a standard for the species. They are derived from multiple sources to create a consensus sequence. This reference genome provides a baseline for comparison which facilitates genetic variation such as mutations, insertions, deletions, SNPs/single nucleotide polymorphisms. Reference genomes are used in a wide range of applications, including gene discovery, comparative genomics, evolutionary studies, and medical research.

### ➤ Other genome assemblies: (remaining)

Help in identifying genetic variation, evolutionary relationships and functional genomics when compared with reference genome.

# Genome Search Results

## ❑ Exploring Assembly (Reference or others):

→ **Assembly statistics:** Genome size, GC%, Assembly level, View Sequences (RefSeq or GenBank)

→ **Sample details**

→ **Quality Analysis:** % of contaminants like adapters, primers given which gives an idea about completeness of sequencing

→ **Annotations details:** View annotations of genes present in the genome (RefSeq, GenBank)

- Genome location and name of gene. Each gene can be explored in Gene secondary nucleotide database
- Protein: Protein coded by that particular gene. Each protein can be explored in Genpept database