# NCBI Resources in the Data Science Era!

**Making the Transition from Sharing Data to Sharing Knowledge**

Ben Busby, Ph.D.

Genomics Outreach Coordinator, Bioinformatics Training Lead

NCBI

Founder, Department of Bioinformatics and Data Science

FAES

ben.busby@nih.gov

NCBI

# NCBI



**Search NCBI databases**                                                  Help

all[sb]                                                    ⊗   Search

**Results found in 38 databases for "all[sb]"**

**Literature**

| Books | 543,504 | books and reports |
|---|---|---|
| MeSH | 266,733 | ontology used for PubMed indexing |
| NLM Catalog | 1,555,942 | books, journals and more in the NLM Collections |
| PubMed | 26,734,411 | scientific & medical abstracts/citations |
| PubMed Central | 4,179,853 | full-text journal articles |

**Health**

| ClinVar | 173,675 | human variations of clinical significance |
|---|---|---|
| dbGaP | 225,011 | genotype/phenotype interaction studies |
| GTR | 48,738 | genetic testing registry |
| MedGen | 293,754 | medical genetics literature and links |
| OMIM | 24,895 | online mendelian inheritance in man |
| PubMed Health | 63,536 | clinical effectiveness, disease and drug reports |

**Genomes**

| Assembly | 102,316 | genome assembly information |
|---|---|---|
| BioProject | 207,505 | biological projects providing data to NCBI |
| BioSample | 5,568,573 | descriptions of biological source materials |
| Clone | 38,170,166 | genomic and cDNA clones |
| dbVar | 6,206,480 | genome structural variation studies |
| Genome | 21,144 | genome sequencing projects by organism |
| GSS | 39,765,380 | genome survey sequences |
| Nucleotide | 222,391,803 | DNA and RNA sequences |
| Probe | 32,405,068 | sequence-based probes and primers |
| SNP | 825,828,843 | short genetic variations |
| SRA | 3,481,910 | high-throughput DNA and RNA sequence read archive |
| Taxonomy | 1,644,293 | taxonomic classification and nomenclature catalog |

**Genes**

| EST | 76,324,331 | expressed sequence tag sequences |
|---|---|---|
| Gene | 26,043,141 | collected information about gene loci |
| GEO DataSets | 2,110,951 | functional genomics studies |
| GEO Profiles | 128,414,055 | gene expression and molecular abundance profiles |
| HomoloGene | 141,268 | homologous gene sets for selected organisms |
| PopSet | 262,192 | sequence sets from phylogenetic and population studies |
| UniGene | 6,473,284 | clusters of expressed transcripts |

**Proteins**

| Conserved Domains | 52,411 | conserved protein domains |
|---|---|---|
| Protein | 342,326,582 | protein sequences |
| Protein Clusters | 820,546 | sequence similarity-based protein clusters |
| Structure | 124,173 | experimentally-determined biomolecular structures |

**Chemicals**

| BioSystems | 932,719 | molecular pathways with links to genes, proteins and chemicals |
|---|---|---|
| PubChem BioAssay | 1,218,723 | bioactivity screening studies |
| PubChem Compound | 92,574,428 | chemical information with structures, information and links |
| PubChem Substance | 225,315,243 | deposited substance and chemical information |

NCBI

# Three (and a half) ways to literature!

# EUtils (Search API) Command Line EDirect

s://github.com/NCBI-Hackathons/EDirect_EUtils_API_Cookbook

## EDirect Scripts

### Gene Aliases

Description (optional):
Written by: NCBI Folks (12/14/2016)
Confirmed by:
Databases: gene

```
esearch -db gene -query "Liver cancer AND Homo sapiens" | \
efetch -format docsum | \
xtract -pattern DocumentSummary -element Name OtherAliases OtherDesignations
```

### Genomic sequence fastas from RefSeq assembly for specified taxonomic designation

Description (optional):
Written by: NCBI Folks (12/14/2016)
Confirmed by: Peter Cooper (NCBI) and Wayne Matten (NCBI) (12/29/2016, v6.00)
Databases: assembly

```
wget `esearch -db assembly -query "Leptospira alstonii[ORGN] AND latest[SB]" | \
efetch -format docsum | \
xtract -pattern DocumentSummary -element FtpPath_RefSeq | \
awk -F"/" '{print $0"/"$NF"_genomic.fna.gz"}'`
```

(For larger sets of data the above may fail as wget may not accept a very large number of arguments.
The command below should work for all.)

Google for
EDirect Cookbook

NIH ⟩ U.S. National Library of Medicine

NCBI

# EDirect Local Caching!

# PubMed and PMC (Open) FTP

ncbiinsights.ncbi.nlm.nih.gov/2017/06/22/pubmed-available-for-download-without-license/?campaign=announce-06222017

## PubMed is now available for download without a license and can be updated every day!

★★★★★ ❶ 4 Votes

*This blog post is directed toward PubMed users.*

Did you know you can download the entire PubMed database, and keep this dataset current with our daily update files? These data are available for free from our FTP site and no longer require a license agreement, whether you're interested in text mining, or want to create your own database for searching and analytics.

Each year in December, NLM releases a comprehensiv... XML format for download. Every day, incremental upda... include new, revised and deleted citations. Please see... information and contact info@ncbi.nlm.nih.gov with qu...

---

NCBI    Resources ⊡    How To ⊡                                                                                                          Sign in to

**PMC**
US National Library of Medicine
National Institutes of Health

PMC ⬍  [                                      ]  Search

Advanced    Journal list

About PMC ▾    For Publishers ▾    Related Resources ▾

### Open Access Subset

The PMC Open Access Subset ⬚ is a part of the total collection of articles in PMC. The articles in the OA Subset are made available under a Creative Commons or similar license that generally allows more liberal redistribution and reuse than a traditional copyrighted work.

To preview the articles or get a current count of articles in the OA Subset, do a search for open access[filter] in PMC. As of 2015, there were over 1 million articles available in this collection.

Please note the following:

- The license terms are not identical for all of the articles in this subset. Please refer to the license statement in each article for specific terms of use.

- The majority of the articles in PMC are subject to traditional copyright restrictions and are not part of this subset.

- Users are directly and solely responsible for compliance with copyright restrictions and are expected to adhere to the terms and conditions defined by the copyright holder (see the PMC Copyright Notice).

# Instead of PubMed FTP…
# An automated tool (alpha)



Pubrunner.org

# For more information go to:
## ncbi.nlm.nih.gov/learn

# NGS (c. 2013) in 90 Seconds for non-bioinformaticians



© Martine Zilversmit 2013

# Cross-Data-Type Descriptors

# Cross-Data-Type Descriptors

# Reporting

# BioSample

# Labels that can be used!

Pathogen or Virus packages.

**Model organism or animal sample**

Use for multicellular samples or cell lines derived from common laboratory model organisms, e.g., mouse, rat, Drosophila, worm, fish, frog, or large mammals includi...

**Metagenome or environmental sample**

Use for metagenomic and environmental samples when i... packages.

⦿ **Invertebrate**

Use for any invertebrate sample.

**Human sample**

WARNING: Only use for human samples or cell lines that have no privacy concerns. For all studies involving human subjects, it is the submitter's responsibility to ensure that the information supplied protects participant privacy in accordance with all applicable laws, regulations and institutional policies. Make sure to remove any direct personal identifiers from your submission. If there are patient privacy concerns regarding making data fully public, please submit samples and data to NCBI's dbGaP database. dbGaP has controlled access mechanisms and is an appropriate resource for hosting sensitive patient data.

For samples isolated from humans use the Pathogen, Microbe or appropriate MIxS package.

**Plant sample**

Use for any plant sample or cell line.

**Virus sample**

Use for all virus samples not directly associated with disease. Viral pathogens should be submitted using the Pathogen: Clinical or host-associated pathogen package.

Google for "BioSample Template"

NCBI

# BioSample

| *sample_name | sample_title | bioproject_accession | *organism | isolate | breed | host | isolation_source | *collection_date | *geo_loc_name | *tissue | age | altitude | bio |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |
| | | | | | | | | | | | | | |

# But wait, how do I find data if the metadata is insufficient?

# Now all of SRA is taxonomically indexed!



*COOL*

*THING*

*#1 !*

NIH U.S. National Library of Medicine

NCBI

# Now all of SRA is taxonomically indexed!

# Now all of SRA is taxonomically indexed!

# But how do I extract the data?

# That said, raw data is messy and hard to dump (plus its huge)

# So why do that?

# Wait, what, that's too simple/primitive?

# Still too simple

## Standalone and API BLAST

**↓ Download BLAST**
Get BLAST databases and executables

**▪ Use BLAST API**
Call BLAST from your application

**☁ Use BLAST in the cloud**
Start an instance at a cloud provider

## Specialized searches

| **SmartBLAST** 🔍 | **Primer-BLAST** 🔍 | **Global Align** 🔍 | **CD-search** 🔍 |
|---|---|---|---|
| Find proteins highly similar to your query | Design primers specific to your PCR template | Compare two sequences across their entire span (Needleman-Wunsch) | Find conserved domains in your sequence |

| **GEO** 🔍 | **IgBLAST** 🔍 | **VecScreen** 🔍 | **CDART** 🔍 |
|---|---|---|---|
| Find matches to gene expression profiles | Search immunoglobulins and T cell receptor sequences | Search sequences for vector contamination | Find sequences with similar conserved domain architecture |

# MAGIC!

# Just download a binary…

# Find something to BLAST into:

PubMed.gov
US National Library of Medicine
National Institutes of Health

[ PubMed ⬍ ]    [                                        ]
Advanced

Format: Abstract ⌄                                                                                    Send to ⌄

## Identification of anti-filarial leads against aspartate semialdehyde dehydrogenase of Wolbachia endosymbiont of Brugia malayi: combined molecular docking and molecular dynamics approaches.

Amala M[1], Rajamanikandan S[2], Prabhu D[1], Surekha K[1], Jeyakanthan J[1].

⊕ Author information

**Abstract**

Lymphatic filariasis is a debilitating vector borne parasitic disease that infects human lymphatic system by nematode Brugia malayi. Currently available anti-filarial drugs are effective only on the larval stages of parasite. So far, no effective drugs are available for humans to treat filarial infections. In this regard, aspartate semialdehyde dehydrogenase (ASDase) in lysine biosynthetic pathway from Wolbachia endosymbiont Brugia malayi represents an attractive therapeutic target for the development of novel anti-filarial agents. In this present study, molecular modeling combined with molecular dynamics simulations and structure-based virtual screening were performed to identify potent lead molecules against ASDase. Based on Glide score, toxicity profile, binding affinity and mode of interactions with the ASDase, five potent lead molecules were selected. The molecular docking and dynamics results revealed that the amino acid residues Arg103, Asn133, Cys134, Gln161, Ser164, Lys218, Arg239, His246, and Asn321 plays a crucial role in effective binding of Top leads into the active site of ASDase. The stability of the ASDase-lead complexes was confirmed by running the 30 ns molecular dynamics simulations. The pharmacokinetic properties of the identified lead molecules are in the acceptable range. Furthermore, density functional theory and binding free energy calculations were performed to rank the lead molecules. Thus, the identified lead molecules can be used for the development of anti-filarial agents to combat the pathogenecity of Brugia malayi.

# This is where I go really fast.

Please check out slideshare to grab the details.

https://www.slideshare.net/benbusby

## Search NCBI databases

aspartate semialdehyde dehydrogenase wolbachia   ✖   **Search**

Results found in 11 databases for **aspartate semialdehyde dehydrogenase wolbachia**

### Literature

| | | |
|---|---|---|
| **Books** | 0 | books and reports |
| **MeSH** | 0 | ontology used for PubMed indexing |
| **NLM Catalog** | 0 | books, journals and more in the NLM Collections |
| **PubMed** | 1 | scientific and medical abstracts/citations |
| **PubMed Central** | 16 | full-text journal articles |

### Health

| | | |
|---|---|---|
| **ClinVar** | 2 | human variations of clinical significance |
| **dbGaP** | 1 | genotype/phenotype interaction studies |
| **GTR** | 0 | genetic testing registry |
| **MedGen** | 0 | medical genetics literature and links |
| **OMIM** | 0 | online mendelian inheritance in man |
| **PubMed Health** | 0 | clinical effectiveness, disease and drug reports |

### Genomes

| | | |
|---|---|---|
| **Assembly** | 0 | genome assembly information |

### Genes

| | | |
|---|---|---|
| **EST** | 0 | expressed sequence tag sequences |
| **Gene** | 6 | collected information about gene loci |
| **GEO DataSets** | 0 | functional genomics studies |
| **GEO Profiles** | 0 | gene expression and molecular abundance profiles |
| **HomoloGene** | 0 | homologous gene sets for selected organisms |
| **PopSet** | 0 | sequence sets from phylogenetic and population studies |
| **UniGene** | 23 | clusters of expressed transcripts |

### Proteins

| | | |
|---|---|---|
| **Conserved Domains** | 0 | conserved protein domains |
| **Identical Protein Groups** | 30 | protein sequences grouped by identity |
| **Protein** | 65 | protein sequences |
| **Protein Clusters** | 1 | sequence similarity-based protein clusters |
| **Sparcle** | 0 | functional categorization of proteins by domain architecture |
| **Structure** | 0 | experimentally-determined biomolecular structures |

Gene    | Gene ÷ | | | Search |
                 Advanced

Full Report ▾                                                                    Send to: ▾

# WD_RS04305  aspartate-semialdehyde dehydrogenase [ *Wolbachia endosymbiont of Drosophila melanogaster* ]

Gene ID: 29555238, updated on 13-Feb-2018

## ▴ Summary

| | |
|---|---|
| **Gene symbol** | WD_RS04305 |
| **Gene description** | aspartate-semialdehyde dehydrogenase |
| **Locus tag** | WD_RS04305 |
| **Gene type** | protein coding |
| **Organism** | Wolbachia endosymbiont of Drosophila melanogaster (strain: wMel, nat-host: Drosophila melanogaster, other: Wolbachia pipientis wMel) |
| **Lineage** | Bacteria; Proteobacteria; Alphaproteobacteria; Rickettsiales; Anaplasmataceae; Wolbachieae; Wolbachia |
| **Old locus tag** | WD0954 |

## ▴ Genomic context

**Sequence:**  NC_002978.6 (913553..914587, complement)

NC_002978.6

[ 911303 ▶ ]                                    [ 915778 ▶ ]
              WD_R094295                WD_R094305 ◄         WD_R9H6115
                     WD_R094300                     WD_R094313

## ▴ Genomic regions, transcripts, and products

**Genomic Sequence:**   NC_002978.6

Go to reference sequence details

Go to nucleotide:   Graphics   FASTA   GenBank

NC_002978.6 ▾ | Find:               ⌄ ⇦ ⇨ ⊖ ——— ⊕ ⚓ ⊞        ✕ Tools ▾  ☰  ● Tracks  ⟳ ? ▾

| 914,700 | 914,600 | 914,500 | 914,400 | 914,300 | 914,200 | 914,100 | 914 K | 913,900 | 913,800 | 913,700 | 913,600 | 913,500 |

Genes
                                                                WD_RS04305
  UP_009629451

  UP_015509289i
STS Markers

| 914,700 | 914,600 | 914,500 | 914,400 | 914,300 | 914,200 | 914,100 | 914 K | 913,900 | 913,800 | 913,700 | 913,600 | 913,500 |

NC_002978.6: 915K..913K (1.3Kbp) C                                     ✎  ● Tracks shown: 3/6

Sequence: NC_002978.6 (913553..914587, complement)

## Genomic regions, transcripts, and products

Genomic Sequence: NC_002978.6

Go to reference sequence details
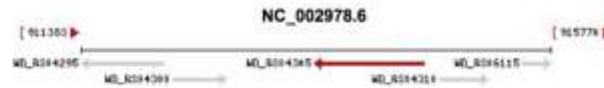
Go to nucleotide: Graphics  FASTA  GenBank

### WD_RS04305

Gene: WD_RS04305
Location: complement(913,553..914,587)
Length: 1,035
Position: 1,014,486
[Qualifiers]
old_locus_tag: WD0954

CDS: WP_010962945.1
Title: aspartate-semialdehyde dehydrogenase
Location: complement(913,553..914,587)
[Length]
Span: 1,035
Product: 344
Position: 1,014,486
[Qualifiers]
inference: COORDINATES: similar to AA sequence:RefSeq:WP_010962945.1

Download: WP_010962945.1

Links & Tools
View GeneID: 29555238 (WD_RS04305)

BLAST Genomic: NC_002978.6 (913,553..914,587)
   FASTA View: NC_002978.6 (913,553..914,587)
GenBank View: NC_002978.6 (913,553..914,587)
BLAST Protein: WP_010962945.1
   FASTA View: NC_002978.6 (913,553..914,587), WP_010962945.1
GenBank View: NC_002978.6 (913,553..914,587), WP_010962945.1
Graphical View: WP_010962945.1

## Bibliography

### Related articles in PubMed

1. Phylogenomics of the reproductive parasite Wolbachia pipientis wMel: a streamlined genome overrun by mobile genetic elements
   Wu M, et al. PLoS Biol, 2004 Mar. PMID 15024419, Free PMC Article

### GeneRIFs: Gene References Into Functions   What's a GeneRIF?

Submit:  New GeneRIF   Correction

## Pathways from BioSystems

2-Oxocarboxylic acid metabolism, organism-specific biosystem (from KEGG)
2-Oxocarboxylic acid metabolism, conserved biosystem (from KEGG)
Biosynthesis of amino acids, organism-specific biosystem (from KEGG)
Biosynthesis of amino acids, conserved biosystem (from KEGG)
Biosynthesis of antibiotics, organism-specific biosystem (from KEGG)
Biosynthesis of secondary metabolites, organism-specific biosystem (from KEGG)

You can now filter BLAST+ Databases like nr, nt and refseq_genomic by taxonomy. Check out https://ftp.ncbi.nlm.nih.gov/blast/db/v5/blastdbv5.pdf for details!

Edit and Resubmit   Save Search Strategies   ▸ Formatting options   ▸ Download

## DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

### Job title: ref|WP_010962945.1| (344 letters)

**RID** EKDMMTCY014 (Expires on 05-03 21:01 pm)

**Query ID** WP_010962945.1
**Description** MULTISPECIES: aspartate-semialdehyde dehydrogenase [Wolbachia]
**Molecule type** amino acid
**Query Length** 344

**Database Name** nr
**Description** All non-redundant GenBank CDS translations+PDB+SwissProt+PIR+PRF excluding environmental samples from WGS projects
**Program** BLASTP 2.8.0+ ▸Citation

Other reports: ▸ Search Summary [Taxonomy reports] [Distance tree of results] [Multiple alignment] [MSA viewer]
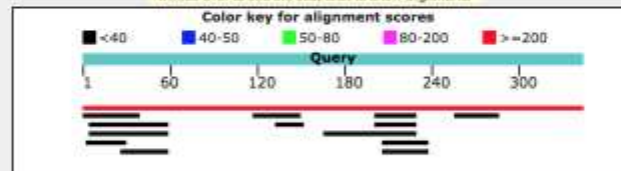
**New** Analyze your query with SmartBLAST

## ⊟Graphic Summary

⊟Show Conserved Domains

Putative conserved domains have been detected, click on the image below for detailed results.

Query seq.
Specific hits
Superfamilies                    Asd superfamily

Distribution of the top 14 Blast Hits on 14 subject sequences ⓘ
Mouse over to see the title, click to show alignments

**Color key for alignment scores**
■ <40   ■ 40-50   ■ 50-80   ■ 80-200   ■ >=200

Query
1      60      120      180      240      300

## ⊟Descriptions

Run PSI-Blast iteration 2 with max [500]   [Go]

⊟Sequences producing significant alignments with E-value BETTER than threshold

Select: All None Selected 0

| Description | Max score | Total score | Query cover | E value | Ident | Accession | Select for PSI blast | Used to build PSSM |
|---|---|---|---|---|---|---|---|---|
| ☐ aspartate-semialdehyde dehydrogenase [Wolbachia endosymbiont of Brugia malayi] | 567 | 567 | 100% | 0.0 | 86% | WP_011256244.1 | ☑ | |

Run PSI-Blast iteration 2 with max [500]   [Go]

▣ Questions/comments

aspartate-semialdehyde dehydrogenase [Wolbachia endosymbiont of Brugia malayi]
Sequence ID: WP_011256244.1  Length: 347  Number of Matches: 1
▶ See 1 more title(s)

Range 1: 1 to 347 GenPept   Graphics                          ▼ Next Match  ▲ Previous Match

| Score | | Expect | Method | Identities | Positives | Gaps |
|---|---|---|---|---|---|---|
| 567 bits(1462) | | 0.0 | Composition-based stats. | 300/347(86%) | 325/347(93%) | 3/347(0%) |

```
Query   1    MRYKIAVIGATGRVGREVLSTLAEFQDE---AIDCVIALASKKSEGKKVSFGDKELTVLC   57
             M  KIAV+GATGRVG EVLS LAEFQDE   +ID VI  ASKKS+GKKVSFG++ELTVLC
Sbjct   1    MGQKIAVVGATGRVGHEVLSILAEFQDEGKISIDSVITFASKKSKGKKVSFGNEELTVLC   60

Query   58   LEDYDFVGTNVAIFCAGSHVSEEYVPIATQAGCIVIDNSSHFRMKEGVPLIIPEINKEKI   117
             LE+YDF    ++AIFCAG HVSE+YVPIAT+AGCIVIDNSS+FRMKEGVPLIIPEINKEKI
Sbjct   61   LENYDFADIDIAIFCAGYHVSEKYVPIATEAGCIVIDNSSYFRMKEGVPLIIPEINKEKI   120

Query   118  MEYKNHNIISNPNCTTIQMLLVLHLLHQKAKIKRIVASTYQSTSGAGKAAMDELYDQTKK   177
             MEYKNHNIISNPNCT IQMLLVLHLL+QKAKIKRIVASTYQSTSGAGKAAMDELY+QTKK
Sbjct   121  MEYKNHNIISNPNCTIIQMLLVLHLLYQKAKIKRIVASTYQSTSGAGKAAMDELYNQTKK   180

Query   178  IFMNEAKKPKIFSKQIAFNCIPHVGEFMENGSTEEEWKMQEETKKILEEDIKVTATCVRV   237
             IF NEAKKP+IF KQIAFNCIPH+GEFME+GST+EEWKMQEETKKILE DIKVTATCVRV
Sbjct   181  IFTNEAKKPEIFPKQIAFNCIPHIGEFMEDGSTKEEWKMQEETKKILEADIKVTATCVRV   240

Query   238  PVFIGHAMAVNVEFDQHITEEQAREVLSEAEDSGVLVYNRREDSEYITQIDVVQENAVYV   297
             PVFIGHA+AVNVEF QHITEEQARE+LSE ED+G+LVY+RR+D +YITQIDVVQENAVYV
Sbjct   241  PVFIGHAIAVNVEFYQHITEEQAREMLSEVEDTGILVYDRRKDGKYITQIDVVQENAVYV   300

Query   298  SRIRRDNTVEHGLNMWIVADNLRKGAALNIVQILEILIREHLSIKCI      344
             SRIR+DNTVEHGLNMWIVADNLRKGAALNIVQILEIL REHLSIKCI
Sbjct   301  SRIRKDNTVEHGLNMWIVADNLRKGAALNIVQILEILTREHLSIKCI      347
```

*The aspartate-semialdehyde dehydrogenases are not dramatically different in Drosophila melanogaster wolbachia and Brugia malayi wolbachia*

Identical Protein Groups

| Identical Protein ⬍ | |
| --- | --- |
| | Advanced |

Identical Protein Groups ▾                                                    Send to: ▾

## aspartate-semialdehyde dehydrogenase

GenPept    FASTA    Graphics    BLAST

*Name:* aspartate-semialdehyde dehydrogenase
*RefSeq Selected Product:* WP_011256244.1, 347 amino acids
*Taxonomic Group:* a-proteobacteria
*Assembly Accessions:* 2
*Protein Accessions:* 2
*CDS Regions:* 2
*Total Rows:* 2

| Source | CDS Region in Nucleotide | Protein | Name | Organism | Strain | Assembly |
| --- | --- | --- | --- | --- | --- | --- |
| RefSeq | NC_006833.1 54209-55252 (-) | WP_011256244.1 | aspartate-semialdehyde dehydrogenase | Wolbachia endosymbiont strain TRS of Brugia malayi | | GCF_000008385.1 |
| INSDC | AE017321.1 54209-55252 (-) | AAW70634.1 | Aspartate-semialdehyde dehydrogenase | Wolbachia endosymbiont strain TRS of Brugia malayi | | GCA_000008385.1 |

Nucleotide

| Nucleotide ⬍ | |
Advanced

FASTA ▾                                                        Send to: ▾

# Wolbachia endosymbiont strain TRS of Brugia malayi, complete genome

NCBI Reference Sequence: NC_006833.1

GenBank    Graphics

```
>NC_006833.1:c55252-54209 Wolbachia endosymbiont strain TRS of Brugia malayi,
complete genome
ATGGGACAAAAAATTGCTGTTGTTGGAGCAACCGGTAGAGTAGGACACGAAGTACTAAGCATACTTGCTG
AGTTCCAAGACGAGGGAAAAATTTCGATAGATTCTGTTATTACATTTGCATCAAAAAAATCAAAGGGAAA
AAAGGTGAGTTTTGGTAACGAAGAATTAACTGTTTTATGCCTTGAAAATTATGACTTTGCTGATATTGAT
ATAGCCATCTTCTGTGCTGGGTACCATGTTTCGGAAAAGTACGTACCGATTGCAACTGAAGCTGGATGTA
TCGTAATAGATAACAGCTCTTATTTTAGGATGAAAGAAGGTGTACCACTAATCATTCCAGAAATTAACAA
AGAAAAAATCATGGAATACAAAAACCACAACATAATATCCAATCCAAACTGTACTATAATACAGATGCTG
TTAGTGCTACATTTATTATACCAAAAAGCAAAAATAAAGAGAATCGTTGCTTCAACTTATCAATCAACCT
CTGGTGCAGGCAAAGCAGCAATGGATGAACTCTATAATCAGACAAAAAAAATCTTCACAAATGAAGCCAA
AAAGCCTGAAATATTCCCTAAGCAAATAGCATTCAATTGCATTCCTCATATAGGAGAGTTTATGGAAGAT
GGTTCTACAAAAGAGGAATGGAAAATGCAAGAGGAAACAAAAAAAATCCTAGAGGCAGATATAAAAGTTA
CTGCAACTTGTGTAAGGGTGCCCGTTTTTATTGGTCATGCTATAGCAGTAAATGTAGAGTTTTACCAGCA
CATAACTGAAGAACAAGCTCGTGAAATGCTAAGTGAAGTTGAAGATACCGGAATTTTAGTGTATGATAGG
CGAAAAGACGGCAAATATATAACCCAAATTGATGTCGTACAGGAGAATGCAGTATACGTATCACGTATTA
GAAAAGACAATACTGTTGAACATGGATTAAATATGTGGATAGTGGCTGATAATCTACGCAAAGGTGCGGC
ACTGAATATAGTACAAATTTTGGAGATCTTGACGAGGGAGCATTTATCAATCAAGTGCATATAG
```

# So here's some data…

# …and the metadata isn't horrible!

# Make a BLAST db, and go!

```
busbybr@ncbimacbook2244:~/magicblast2$ tar -xvzf ncbi-magicblast-1.3.0-x64-macosx\ \(1\).tar.gz
x ncbi-magicblast-1.3.0/
x ncbi-magicblast-1.3.0/bin/
x ncbi-magicblast-1.3.0/bin/makeblastdb
x ncbi-magicblast-1.3.0/bin/magicblast
x ncbi-magicblast-1.3.0/ncbi_package_info
x ncbi-magicblast-1.3.0/README
x ncbi-magicblast-1.3.0/ChangeLog
busbybr@ncbimacbook2244:~/magicblast2$ ./ncbi-magicblast-1.3.0/bin/makeblastdb -in ../Desktop/Wolbachia_Bm_ASD.fasta -dbtype nu
cl -parse_seqids

Building a new DB, current time: 05/02/2018 09:25:51
New DB name:   /Users/busbybr/Desktop/Wolbachia_Bm_ASD.fasta
New DB title:  ../Desktop/Wolbachia_Bm_ASD.fasta
Sequence type: Nucleotide
Keep MBits: T
Maximum file size: 1000000000B
Adding sequences from FASTA; added 1 sequences in 0.0219159 seconds.
busbybr@ncbimacbook2244:~/magicblast2$ ./ncbi-magicblast-1.3.0/bin/magicblast -db ../Desktop/Wolbachia_Bm_ASD.fasta
Wolbachia_Bm_ASD.fasta      Wolbachia_Bm_ASD.fasta.nin  Wolbachia_Bm_ASD.fasta.nsd  Wolbachia_Bm_ASD.fasta.nsq
Wolbachia_Bm_ASD.fasta.nhr  Wolbachia_Bm_ASD.fasta.nog  Wolbachia_Bm_ASD.fasta.nsi
busbybr@ncbimacbook2244:~/magicblast2$ ./ncbi-magicblast-1.3.0/bin/magicblast -db ../Desktop/Wolbachia_Bm_ASD.fasta -no_unalign
ed -splice F -num_threads 2 -sra SRR3111492 -out SRR3111492 into Wol Bm ASD.sam &
[1] 6726
busbybr@ncbimacbook2244:~/magicblast2$ 
```

software carpentry

# Make a BLAST db, and go!

**tar -xvzf ncbi-magicblast...**

**makeblastdb -dbtype nucl -in <fasta> -parse_seqids**

**magicblast -db <fasta> -sra SRR... -splice F -no_unaligned**

45

NCBI

# It runs pretty quick…

**magicblast -db <fasta> -sra SRR... -splice F -no_unaligned**

*-num_threads X*

```
43052  Google Chrom 0.0  02:25.01 15   0   149  34M    0B   152M   647   647   sleeping *0[1]        0.00000 0.00000  11102
40069  syspolicyd   0.0  00:00.42 2    1   36   16K    0B   2572K  40069 1     sleeping  0[5]        0.00000 0.00000  0
40068  Microsoft Po 0.1  04:23.05 17   5   282  12M    0B   119M   40068 1     sleeping *44[8]       0.00000 0.00000  11102
29024  Slack Helper 0.0  01:18.13 17   0   117  2768K  0B   64M    29018 29018 sleeping *0[1]        0.00000 0.00000  11102
29022  crashpad_han 0.0  00:00.13 4    0   28   4096B  0B   960K   29021 1     sleeping *0[1]        0.00000 0.00000  11102
29019  Slack Helper 0.0  02:27.00 5    0   225  1856K  0B   55M    29018 29018 sleeping *0[1]        0.00000 0.00000  11102
29018  Slack        0.0  11:28.36 43   1   481  49M    0B   60M    29018 1     sleeping *0[2754]     0.00000 0.00000  11102
26366  screencaptur 0.0  00:00.20 4    2   56   2292K  20K  0B     661   661   sleeping *0[1]        0.00000 0.00000  11102
26218  trivial-rewr 0.0  00:00.01 1    0   18   744K   0B   0B     26214 26214 sleeping *0[1]        0.00000 0.00000  27
26217  cleanup      0.0  00:00.02 1    0   18   808K   0B   0B     26214 26214 sleeping *0[1]        0.00000 0.00000  27
26216  qmgr         0.0  00:00.01 1    0   18   752K   0B   0B     26214 26214 sleeping *0[1]        0.00000 0.00000  27
26215  pickup       0.0  00:00.01 1    0   18   784K   0B   0B     26214 26214 sleeping *0[1]        0.00000 0.00000  27
26214  master       0.0  00:00.02 1    0   18   804K   0B   0B     26214 1     sleeping *0[1]        0.00000 0.00000  0
26211  magicblast   8.4  00:08.42 2    0   15   309M+  0B   0B     26211 5031  sleeping *0[1]        0.00000 0.00000  11102
26203  top          4.9  00:06.48 1/1  0   28   4792K  0B   964K   26203 26156 running  *0[1]        0.00000 0.00000  0
```

# Even on a plane!

| 43052 | Google Chrom | 0.0 | 02: | | | | | | | | | | | 0.00000 | 0.00000 | 11102 |
| 40069 | syspolicyd | 0.0 | 00: | | | | | | | | | | | 0.00000 | 0.00000 | 0 |
| 40068 | Microsoft Po | 0.1 | 04:23.05 | 17 | 5 | 282 | 12M | 0B | 119M | 40068 | 1 | | sleeping | *44[8] | 0.00000 | 0.00000 | 11102 |
| 29024 | Slack Helper | 0.0 | 01:18.13 | 17 | 0 | 117 | 2768K | 0B | 64M | 29018 | 29018 | sleeping | *0[1] | | 0.00000 | 0.00000 | 11102 |
| 29022 | crashpad_han | 0.0 | 00:00.13 | 4 | 0 | 28 | 4096B | 0B | 960K | 29021 | 1 | | sleeping | *0[1] | 0.00000 | 0.00000 | 11102 |
| 29019 | Slack Helper | 0.0 | 02:27.00 | 5 | 0 | 225 | 1856K | 0B | 55M | 29018 | 29018 | sleeping | *0[1] | | 0.00000 | 0.00000 | 11102 |
| 29018 | Slack | 0.0 | 11:28.36 | 43 | 1 | 481 | 49M | 0B | 60M | 29018 | 1 | | sleeping | *0[2754] | 0.00000 | 0.00000 | 11102 |
| 26366 | screencaptur | 0.0 | 00:00.20 | 4 | 2 | 56 | 2292K | 20K | 0B | 661 | 661 | sleeping | *0[1] | | 0.00000 | 0.00000 | 11102 |
| 26218 | trivial-rewr | 0.0 | 00:00.01 | 1 | 0 | 18 | 744K | 0B | 0B | 26214 | 26214 | sleeping | *0[1] | | 0.00000 | 0.00000 | 27 |
| 26217 | cleanup | 0.0 | 00:00.02 | 1 | 0 | 18 | 808K | 0B | 0B | 26214 | 26214 | sleeping | *0[1] | | 0.00000 | 0.00000 | 27 |
| 26216 | qmgr | 0.0 | 00:00.01 | 1 | 0 | 18 | 752K | 0B | 0B | 26214 | 26214 | sleeping | *0[1] | | 0.00000 | 0.00000 | 27 |
| 26215 | pickup | 0.0 | 00:00.01 | 1 | 0 | 18 | 784K | 0B | 0B | 26214 | 26214 | sleeping | *0[1] | | 0.00000 | 0.00000 | 27 |
| 26214 | master | 0.0 | 00:00.02 | 1 | 0 | 18 | 804K | 0B | 0B | 26214 | 1 | | sleeping | *0[1] | 0.00000 | 0.00000 | 0 |
| 26211 | magicblast | 8.4 | 00:08.42 | 2 | 0 | 15 | 309M+ | 0B | 0B | 26211 | 5031 | sleeping | *0[1] | | 0.00000 | 0.00000 | 11102 |
| 26203 | top | 4.9 | 00:06.48 | 1/1 | 0 | 28 | 4792K | 0B | 964K | 26203 | 26156 | running | *0[1] | | 0.00000 | 0.00000 | 0 |

# Bam! (well, .sam)

# Similar searches on the web!

# Works with some other software!



Secure | https://trace.ncbi.nlm.nih.gov/Traces/sra/sra.cgi?view=software

Main | Browse | Search | Download | Submit | **Software** | Trace Archive | Trace Assembly | Trace BLAST

**Download** | Toolkit Documentation | XML Schema

## NCBI SRA Toolkit

Please consult SRA Toolkit Documentation for help.
Below are the latest releases of various tools and release checksum file.

### SRA Toolkit

Compiled binaries of March 14, 2017, version 2.8.2-1 release:

- CentOS Linux 64 bit architecture
- Ubuntu Linux 64 bit architecture
- MacOS 64 bit architecture
- MS Windows 64 bit architecture

### Magic-BLAST

Magic-BLAST is a tool for mapping large next-generation RNA or DNA sequencing runs against a whole genome or transcriptome.

- Magic-BLAST executables for LINUX, MacOSX, and Windows as well as the source files are available on the FTP site
- Read more about Magic BLAST on the FTP site

### Third Party Software

Builds of Third Party Software Tools with SRA support ( NGS 1.3.0 release ):

- Genome Analysis Toolkit (GATK) version 3.6-6-ngs.1.3.0 - including direct support of SRA
- HISAT2 version 2.0.6-ngs.1.3.0 - graph-based alignment of next generation sequencing reads to a population of genomes with direct support of SRA, built for Linux 64 bit architecture

### Latest Source Code

- NGS Software Development Kit – October 7, 2016, version 1.3.0 release
- NCBI VDB Software Development Kit – March 7, 2017, version 2.8.2 release
- NCBI SRA Toolkit – March 14, 2017, version 2.8.2-1 release



NIH〉U.S. National Library of Medicine

50

NCBI

# This is how…

# Prokaryotic Genome Annotation

Genome

Genome  ⬍ | [                                                                    ] | Search

Limits   Advanced

Prokaryotic Annotation Home | Documentation ▼ | Complete Genome Submission ▼ | WGS Genome Submission ▼

## NCBI Prokaryotic Genome Annotation Pipeline

NCBI Prokaryotic Genome Annotation Pipeline is designed to annotate bacterial and archaeal genomes (chromosomes and plasmids).

Genome annotation is a multi-level process that includes prediction of protein-coding genes, as well as other functional genome units such as structural RNAs, tRNAs, small RNAs, pseudogenes, control regions, direct and inverted repeats, insertion sequences, transposons and other mobile elements.

NCBI has developed an automatic prokaryotic genome annotation pipeline that combines *ab initio* gene prediction algorithms with homology based methods. The first version of NCBI Prokaryotic Genome Automatic Annotation Pipeline (PGAAP; see Pubmed Article) developed in 2005 has been replaced with an upgraded version that is capable of processing a larger data volume. NCBI's annotation pipeline depends on several internal databases and is not currently available for download or use outside of the NCBI environment.

Related documentation:

- Annotation process
- Annotation standards
- Pipeline Release notes

### GenBank

The NCBI prokaryotic annotation pipeline is available as a service for GenBank submitters. The pipeline is capable of annotating both complete genomes and draft WGS genomes consisting of multiple contigs. You can request PGAP annotation when you submit your genome to GenBank.

Both WGS and non-WGS genomes, including gapless complete bacterial chromosomes, can be submitted via the Submission Portal. You will be asked to choose whether the genome being submitted is considered WGS or not. The differences for GenBank purposes are: non-WGS    Each chromosome is in a single sequence and there are no extra sequences    Each sequence in the genome must be assigned to a chromosome or plasmid or organelle    Plasmids and organelles can still be in multiple pieces. WGS    One or more chromosomes are in multiple pieces and/or some sequences are not assembled into chromosomes In both cases:    There can still be gaps within the sequences; you will supply that information in the submission    Plasmids and organelles can still be in multiple pieces.    Internal sequences must be arranged in the correct order and orientation.    Sequences concatenated in unknown order are not allowed. Submission is through the Genome Submission Portal. See the genome submission instructions page for details.

### Refseq

All RefSeq bacterial and archaeal genomes, with the exception of RefSeq Prokaryotic Reference Genomes, are annotated using NCBI's prokaryotic genome annotation pipeline. Additional information on this policy is available here:

# Prokaryotic Genome Annotation

# EUtils (Search API) Command Line EDirect

s://github.com/NCBI-Hackathons/EDirect_EUtils_API_Cookbook

## EDirect Scripts

### Gene Aliases

Description (optional):
Written by: NCBI Folks (12/14/2016)
Confirmed by:
Databases: gene

```
esearch -db gene -query "Liver cancer AND Homo sapiens" | \
efetch -format docsum | \
xtract -pattern DocumentSummary -element Name OtherAliases OtherDesignations
```

### Genomic sequence fastas from RefSeq assembly for specified taxonomic designation

Description (optional):
Written by: NCBI Folks (12/14/2016)
Confirmed by: Peter Cooper (NCBI) and Wayne Matten (NCBI) (12/29/2016, v6.00)
Databases: assembly

```
wget `esearch -db assembly -query "Leptospira alstonii[ORGN] AND latest[SB]" | \
efetch -format docsum | \
xtract -pattern DocumentSummary -element FtpPath_RefSeq | \
awk -F"/" '{print $0"/"$NF"_genomic.fna.gz"}'`
```

(For larger sets of data the above may fail as wget may not accept a very large number of arguments. The command below should work for all.)

Google for
EDirect Cookbook

# EUtils (Search API) Command Line EDirect

## Get all SRA runs for a given BioProject

Description (optional):

Written by: Bob Sanders (3/22/2017)

Confirmed by:

Databases: SRA, BioProject

Google for
EDirect Cookbook

```
esearch -db bioproject -query "PRJNA356464" | elink -target sra | efetch -format docsum | \
xtract -pattern DocumentSummary -ACC @acc -block DocumentSummary -element "&ACC"
```

## Get latitiude and longitude for SRA Datasets (e.g. outbreaks and metagenomes)

Description (optional):

Written by: BB, Mike D, Rob Edwards (4/12/2017)

Confirmed by:

Databases: SRA, BioSample

```
for i in $(cat sra_ids.txt); do ll=$(esearch -db sra -query $i | \
elink -target biosample | efetch -format docsum | \
xtract -pattern DocumentSummary -block Attribute -if Attribute@attribute_name -equals lat_lon -element Attribute); \
echo -e "$i\t$ll"; done
```

# More Information…

# SNP Calling on the Fly!

Polygenic SNP Search Tool
https://github.com/NCBI-Hackathons/PSST

## https://github.com/NCBI-Hackathons/GenomicRobots

**Access a local repository of sequencing data**

**Input requirements for precomputed analysis:**
- VCFs
- AAFs for each rsID in VCFs
- rsIDs and MAFs from public data (e.g. GNOMAD)

- FASTQs

- Paths and names of input files
- Precomputed/on-the-fly switch

**Obfuscate genotypes to protect sample identity**

**Strategic flipping**
Based on AF in local and public data, decides whether to flip a SNP from alt to ref for reporting purposes

**Random flipping**
SNPs with MAF<1% in your data are randomly flipped from alt to ref for reporting purposes

- Flipping defaults to on but can be turned off (e.g. for queries by owners of the data)

**Report results for each queried SNP**

- **Yes** reported if SNP is in database and has not been flipped to reference
- **No** reported if SNP is not in database or has been flipped to reference

- Future work: add an option for verbose reporting (e.g. including sample names and zygosity) for queries by owners of the data

# Novel Virus Discovery!



```
cd Virus_Detection_SRA/cwl/tools
wget ftp://ftp.ncbi.nlm.nih.gov/genomes/all/GCF/000/848/505/GCF_000848505.1_ViralProj14703/GCF_000848505.1_ViralProj1
gunzip GCF_000848505.1_ViralProj14703_genomic.fna.gz
makeblastdb -dbtype nucl -in GCF_000848505.1_ViralProj14703_genomic.fna -out ebolazaire -parse_seqids
export BLASTDB=$BLASTDB:`pwd`
```

These steps downloaded the Ebola virus genome and uncompressed it. Using the Ebola virus genome, a BLAST
database was created with `makeblastdb`. Then your local directory was added to the BLASTDB environmental variable.

```
sidearm.cwl sidearm.SRR1553459.ebola.yml
```

https://github.com/NCBI-Hackathons/ViruSpy

https://github.com/NCBI-Hackathons/EndoVir

**Initial run**

RefSeq Virus → magicBLAST → Assembler → Motif screen → Rejected contigs

SRA

**BUD**

Overlapper ← magicBLAST ← Flanker ← Contigs

Extender → Contig inspector ← Motif screen

Contig merger

External tools
Database
--- Piping (no disk I/O)

# NCBI Hackathons!

# Antibiotic Resistance Detection!

# Antibiotic Resistance Detection!

```
tblastn_vdb -query mdr_sequence.aa -db "SRR1427233 SRR515906" -outfmt 6 -max_target_seqs 2500 -out
sra_tblastn.tab
```

https://www.ncbi.nlm.nih.gov/core/assets/sra/files/Factsheet_SRA.pdf

## Metagenomics Discovery Challenge (MDC) Press Release

CUNY Academic Commons

HOME          METAGENOMICS DISCOVERY CHALLENGE (MDC) PRESS RELEASE

## Metagenomics Discovery Challenge (MDC)

Leave a reply

Search

Starting March 2nd, Kingsborough Community College (**KCC**) of the City University of New York (**CUNY**) will be hosting Spring 2018 CUNY-wide *Metagenomics Discovery Challenge* (**MDC CUNY**). **MDC CUNY** – it's a certificate-based course. The independent study course is offered in a hybrid format and it will be using the Team-Based Learning (**TBL**) pedagogy and Open Educational Resources (**OER**). It doesn't have any prerequisites, but a background in biological sciences or computer science would be helpful. Students will learn to do bio-surveillance on public metagenomic datasets*, identifying fungal (and other micro-eukaryotic), bacterial, archaeal and viral

**RECENT POSTS**

Metagenomics Discovery Challenge (MDC)

**RECENT COMMENTS**

# An Educational Resource for RNAseq

Available to

anyone!

# Part of an Online Workshop

## First 5 lectures

## now available

## on YouTube

## Load a fasta file for use in Biopython

In this step, we want to load the yakuba.fa sequence into a variable that can be used in our blast search. To to this we create a variable called `fasta_file` and use Python's `open()` function to read the file. As shown above, the yakuba file is in a folder called files at `./files/yakuba.fa`

```
In [ ]:  # Complete this code by entering the name of your file. The filename and
         # filepath should be in quotes

         fasta_file = open().read()
```

```
In [2]:  fasta_file = open('./files/yakuba.fa').read()
```

We can preview what was read into the fasta file by printing it:

```
In [ ]:  print(fasta_file)
```

## Preform a BLAST search using Biopython

As mentioned in the introduction, BLAST is a tool for similarity searching. This is done by taking your **query** sequence (the sequence you want to find matches for), as well as **search parameters** (some optional adjustments to the way you wish to limit or expand your search) and searching a **database** (a repository of known DNA sequences).

First, we will load the appropriate Biopython module for doing a BLAST search over the Internet. The NCBIWWW module has a variety of features we will explore in a moment.

```
In [4]:  from Bio.Blast import NCBIWWW
```

We will do our first BLAST using this piece of Biopython code.

> tip: Since this is a real BLAST search, you will get an 'In [*]' in the cell below for up to several minutes as the search is executed. Don't proceed in the notebook until the '*' turns into a number.

# Collaboration!

https://ncbi-hackathons.github.io/GeneExpressionAging/ideogram

# deSRA/viSRA (prototype)

# VIRGO (prototype)

# Translating from Bioinformatics and Clinical Informatics

| EMR with variants (JSON format) | → | "Useful" (i.e. deleterious, pathogenic) variants | → | Associated genotypes and phenotypes | → | EMR with variants and associated genotypes and phenotypes learned |
|---|---|---|---|---|---|---|

| Tools: | | SnpEff and Annovar | | GWAS catalog | | FHIR-compliant JSON |
|---|---|---|---|---|---|---|



Use case 2

Use case 0

Use case 1

Bioinformatics Software

Docker

CWL

NCBI

# Other People's Hackathons

# Communication



@DCGenomics

# Creating a Community



https://ncbi-hackathons.github.io

# Creating a Community

Come work at NCBI for 4-6 weeks!
Email bioinformatics-training@ncbi.nlm.nih.gov
for more information!

# Creating a Community



https://ncbi-hackathons.github.io