



[www.esaunggul.ac.id](http://www.esaunggul.ac.id)

# PENGANTAR BIOINFORMATIKA

## IBT 431

*By Seprianto S.Pi, M.Si*



Pertemuan 9

# Analisis BLAST

# Sasaran Perkuliahan

- Mahasiswa mampu menganalisis hasil sekuens dengan Teknik BLAST (BLASTn, BLASTx, proteinBLAST)
- Mahasiswa Mampu menerjemahkan hasil penelusuran BLAST dari nilai Query Covarage, E-Value dan Maximum identity

# What is BLAST?

Free, online service from National Center for Biotechnology Information (NCBI)

**BLAST** Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

My NCBI [Sign In] [Register]

► NCBI/BLAST Home

BLAST finds regions of similarity between biological sequences. [more...](#)

**New** Aligning Multiple Protein Sequences? Try the [COBALT Multiple Alignment Tool](#). [Go](#)

**BLAST Assembled Genomes**

Choose a species genome to search, or [list all genomic BLAST databases](#).

- [Human](#)
- [Mouse](#)
- [Rat](#)
- [Arabidopsis thaliana](#)
- [Oryza sativa](#)
- [Bos taurus](#)
- [Danio rerio](#)
- [Drosophila melanogaster](#)
- [Gallus gallus](#)
- [Pan troglodytes](#)
- [Microbes](#)
- [Apis mellifera](#)

**Basic BLAST**

Choose a BLAST program to run.

<a href="#">nucleotide blast</a>	Search a <b>nucleotide</b> database using a <b>nucleotide</b> query <i>Algorithms:</i> blastn, megablast, discontinuous megablast
<a href="#">protein blast</a>	Search <b>protein</b> database using a <b>protein</b> query <i>Algorithms:</i> blastp, psi-blast, phi-blast
<a href="#">blastx</a>	Search <b>protein</b> database using a <b>translated nucleotide</b> query
<a href="#">tblastn</a>	Search <b>translated nucleotide</b> database using a <b>protein</b> query
<a href="#">tblastx</a>	Search <b>translated nucleotide</b> database using a <b>translated nucleotide</b> query

**News**

[BLAST 2.2.23 release](#)

A new version of the stand-alone applications is available.  
Mon, 22 Mar 2010 15:00:00 EST

[More BLAST news...](#)

**Tip of the Day**

[How to do Batch BLAST jobs.](#)

BLAST makes it easy to examine a large group of potential gene candidates.

[More tips...](#)

<http://blast.ncbi.nlm.nih.gov/Blast.cgi>

# What is BLAST?

**BLAST** : Nucleotide/Protein  
Sequence Databases

Basic Local Alignment Search Tool

as

**Google** : Internet

# Some Uses for BLAST

- Identify an unknown sequence
- Build a homology tree for a protein and Nucleotide
- Get clues about protein structure by finding similar proteins with known structures
- Map a sequence in a genome
- Etc., etc.

# **BLAST** helps you to find homologous genes and proteins

## **Homologous Proteins (or genes)**

- **Have a common ancestor (they're related)**
  - **Have similar structures**
  - **Have similar functions**

# BLAST is Search Tool

- By aligning query sequence against all sequences in a database, alignment can be used to search database for similar sequences
- But alignment algorithms are slow



# Why do we need local alignments?

- To compare a short sequence to a large one.
- To compare a single sequence to an entire database
- To compare a partial sequence to the whole.
- Identify newly determined sequences
- Compare new genes to known ones
- Guess functions for entire genomes full of ORFs of unknown function

# Pairwise Alignment

## Global

- Best score from among alignments of full-length sequences
- Needleman-Wunch algorithm

## Local

- Best score from among alignments of partial sequences
- Smith-Waterman algorithm

# BLAST programs

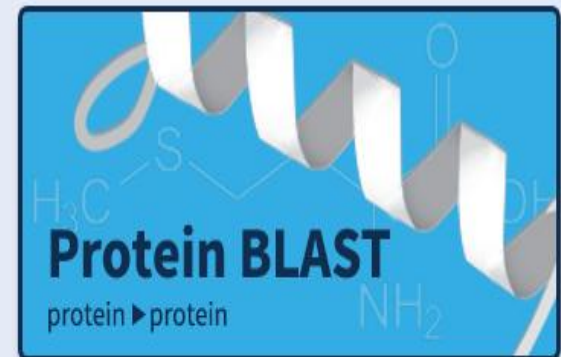
Program	Description
blastp	Compares an amino acid query sequence against a protein sequence database.
blastn	Compares a nucleotide query sequence against a nucleotide sequence database.
blastx	Compares a nucleotide query sequence translated in all reading frames against a protein sequence database. You could use this option to find potential translation products of an unknown nucleotide sequence.
tblastn	Compares a protein query sequence against a nucleotide sequence database dynamically translated in all reading frames.

**How do I input a query into BLAST?**

# Choose which “flavor” of BLAST to use

- **BLAST comes in many “flavors”**
  - Protein BLAST (BLASTp)
    - Compares a protein query with sequences in GenBank protein database
  - Nucleotide BLAST (BLASTn)
  - BLASTx

## Web BLAST



# Enter your “query” sequence

- A sequence can be input as a (an)
  - FASTA format sequence
  - Accession number
  - Choose file
  - Protein blast can only accept amino acid sequences

The screenshot shows the BLASTN web interface. At the top, there are tabs for 'blastn', 'blastp', 'blastx', 'tblastn', and 'tblastx'. The main heading is 'Enter Query Sequence'. Below this, there is a text input field for 'Enter accession number(s), gi(s), or FASTA sequence(s)'. To the right of this field are 'Clear' and 'Query subrange' options. The 'Query subrange' section includes 'From' and 'To' input fields. Below the main input field, there is an 'Or, upload file' section with a 'Choose File' button and the text 'No file chosen'. At the bottom, there is a 'Job Title' section with an input field and the text 'Enter a descriptive title for your BLAST search'.

blastn blastp blastx tblastn tblastx

BLASTN programs search nucleotide databases using a nucleotide query. [more...](#) [Reset page](#) [Bookmark](#)

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s)

Clear Query subrange

From

To

Or, upload file  No file chosen

Job Title

Enter a descriptive title for your BLAST search


# Choose search set

- **Choose which database to search**
  - Default is non-redundant protein sequences (nr)
    - Searches all databases that contain protein sequences

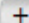
## Choose Search Set


**Database**

Human genomic + transcript  Mouse genomic + transcript  Others (nr etc.):

Nucleotide collection (nr/nt) 

**Organism**  
Optional

Exclude 

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown 


**Exclude**  
Optional


Models (XM/XP)  Uncultured/environmental sample sequences

**Limit to**  
Optional

Sequences from type material

**Entrez Query**  
Optional

 [Create custom database](#)

Enter an Entrez query to limit search 

# Choose organism

- Default is all organisms represented in databases
- Use this to limit your search to one organism (eg. Yeast)

Choose Search Set

Database  Human genomic + transcript  Mouse genomic + transcript  Others (nr etc.):  
Nucleotide collection (nr/nt)

Organism Optional  
Exclude Optional  
Limit to Optional  
Entrez Query Optional

Program Selection  
Optimize for

**BLAST**

[+ Algorithm parameter](#)

esche

- Escherichia (taxid:561)
- Escherichia coli (taxid:562)
- Escherichia sp. 3\_2\_53FAA (taxid:562)
- Escherichia sp. MAR (taxid:562)
- Escherichia coli O157:H7 (taxid:83334)
- Escherichia coli O26:H11 (taxid:244319)
- Escherichia freundii (taxid:546)
- Escherichia coli O104:H4 (taxid:1038927)
- Escherichia coli O111:NM (taxid:373045)
- Escherichia coli K-12 (taxid:83333)
- Escherichia coli O121:H19 (taxid:991915)
- Escherichia albertii (taxid:208962)
- Escherichia coli O157 (taxid:1045010)
- Escherichia coli O145:NM (taxid:991919)
- Escherichia coli O111:H8 (taxid:991910)
- Escherichia fergusonii (taxid:564)
- Escherichia adecarboxylata (taxid:83655)
- Escherichia coli MG1655 (taxid:511145)
- Escherichia coli str. K-12 substr. MG1655 (taxid:511145)
- Escherichia coli O111:H11 (taxid:1163390)

Exclude +

Exclude

[Create custom database](#)

Optimize for highly similar sequences



# BLAST off!!

- Choose Highly similar sequence (MegaBLAST)
- Click on the BLAST button at the bottom of the page!

## Program Selection

Optimize for

- Highly similar sequences (megablast)  
 More dissimilar sequences (discontiguous megablast)  
 Somewhat similar sequences (blastn)

Choose a BLAST algorithm 

**BLAST**

Search database **Nucleotide collection (nr/nt)** using **Megablast (Optimize for highly similar sequences)**

Show results in a new window

[+ Algorithm parameters](#)

## Basic Local Alignment Search Tool


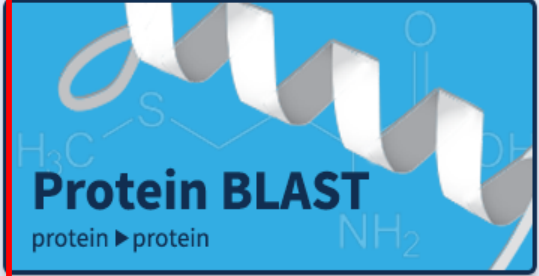
**BLAST** finds regions of similarity between biological sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance. [Learn more](#)

The Basic Local Alignment Search Tool (BLAST) finds regions of local similarity between sequences. The program compares nucleotide or protein sequences to sequence databases and calculates the statistical significance of matches. BLAST can be used to infer functional and evolutionary relationships between sequences as well as help identify members of gene families.

**NEWS**

**IgBLAST 1.9.0 released**  
 IgBLAST now supports AIRR rearrangement reports.  
 Fri, 18 May 2018 08:00:00 EST [More BLAST news...](#)

## Web BLAST

 <p><b>Nucleotide BLAST</b> nucleotide ▶ nucleotide</p>	<p><b>blastx</b> translated nucleotide ▶ protein</p> <p><b>tblastn</b> protein ▶ translated nucleotide</p>	 <p><b>Protein BLAST</b> protein ▶ protein</p>
--	--	--

## BLAST Genomes

**blastn** blastp blastx tblastn tblastx

BLASTN programs search nucleotide databases using a nucleotide query. [more...](#)

[Reset page](#) [Bookmark](#)

### Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s)

[Clear](#)

Query subrange

From

To

Or, upload file

Choose File

No file chosen

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

### Choose Search Set

Database

Human genomic + transcript  Mouse genomic + transcript  Others (nr etc.):

Nucleotide collection (nr/nt)

Limit by

Organism  BioProjectID  WGS Project

Organism  
Optional

Exclude +

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown.

Exclude  
Optional

Models (XM/XP)  Uncultured/environmental sample sequences

Limit to  
Optional

Sequences from type material

Entrez Query  
Optional

[YouTube](#) [Create custom database](#)

Enter an Entrez query to limit search

### Program Selection

Optimize for

Highly similar sequences (megablast)

More dissimilar sequences (discontiguous megablast)

Somewhat similar sequences (blastn)

Choose a BLAST algorithm

Choose Search Set

Database  Human genomic + transcript  Mouse genomic + transcript  Others (nr etc.):

Nucleotide collection (nr/nt)

Organism Optional   Exclude +

Enter organism common name, binomial, or tax id. Only 20 top taxa will be shown

Exclude Optional  Models (XM/XP)  Uncultured/environmental sample sequences

Limit to Optional  Sequences from type material

Entrez Query Optional  [YouTube](#) [Create custom database](#)

Enter an Entrez query to limit search

Program Selection

Optimize for  Highly similar sequences (megablast)

More dissimilar sequences (discontiguous megablast)

Somewhat similar sequences (blastn)

Choose a BLAST algorithm

**BLAST**

Search database Nucleotide collection (nr/nt) using Megablast (Optimize for highly similar sequences)

Show results in a new window

+ [Algorithm parameters](#)

BLAST is a registered trademark of the National Library of Medicine

[Support center](#) [Mailing list](#) [YouTube](#)

## Specialized searches

### SmartBLAST



Find proteins highly similar to your query

### Primer-BLAST



Design primers specific to your PCR template

### Global Align



Compare two sequences across their entire span (Needleman-Wunsch)

### CD-search



Find conserved domains in your sequence

### GEO



Find matches to gene expression profiles

### IgBLAST



Search immunoglobulins and T cell receptor sequences

### VecScreen



Search sequences for vector contamination

### CDART



Find sequences with similar conserved domain architecture

### Targeted Loci



Search markers for phylogenetic analysis

### Multiple Alignment



Align sequences using domain and protein constraints

### BioAssay



Search protein or nucleotide targets in PubChem BioAssay

### MOLE-BLAST



Establish taxonomy for uncultured or environmental sequences

**How do I interpret the results of a BLAST search?**

# BLAST creates **local** alignments

- **What is a local alignment?**
  - **BLAST** looks for similarities between **regions** of two sequences

```
Global  FGFTALILLAVKV  
        F--TAL-LLA--V
```

```
Local   FGFTALILL-AVKAV  
        --FTAL-LLAAV---
```

# The Graphic Display

## 1. How good is the match?

- **Red = excellent!**
- **Pink = pretty good**
- **Green = OK, but look at other factors**
- **Blue = bad**
- **Black = really bad!**

## 2. How long are the matched segments?

**Longer = better**



**Query ID** | Query\_55165  
**Description** | SPB  
**Molecule type** | nucleic acid  
**Query Length** | 1310

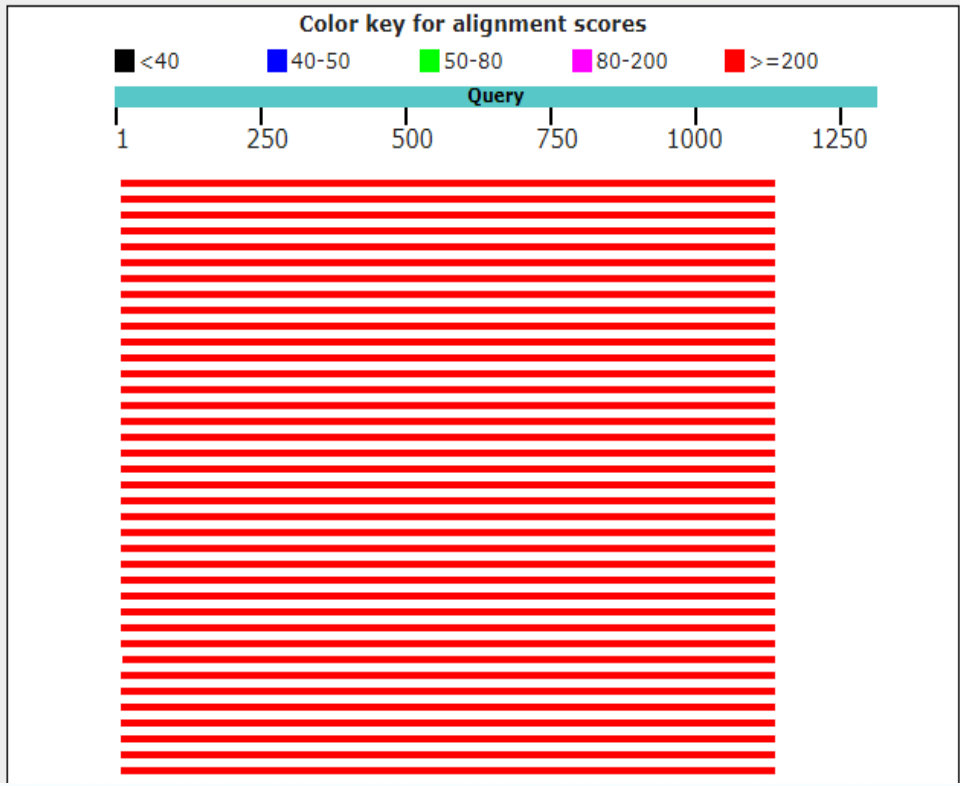
**Database Name** | nr  
**Description** | Nucleotide collection (nt)  
**Program** | BLASTN 2.8.0+ [▶ Citation](#)

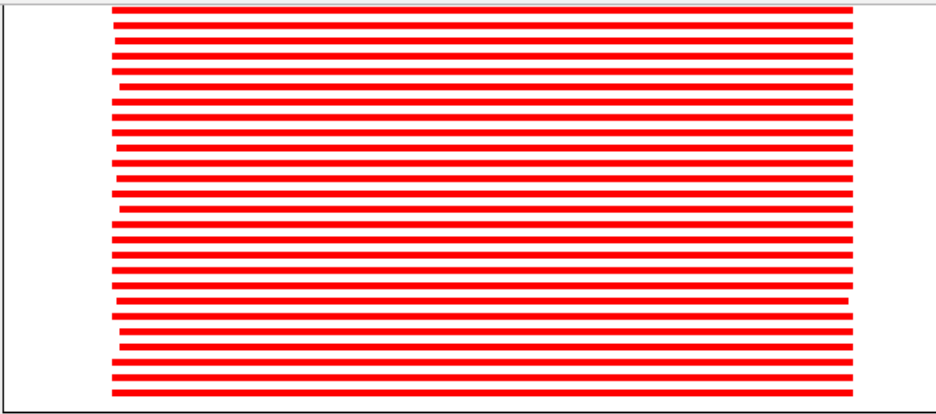
Other reports: [▶ Search Summary](#) [\[Taxonomy reports\]](#) [\[Distance tree of results\]](#) [\[MSA viewer\]](#)

☐ **Graphic Summary**

Distribution of the top 119 Blast Hits on 100 subject sequences 📄

Mouse over to see the title, click to show alignments





Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

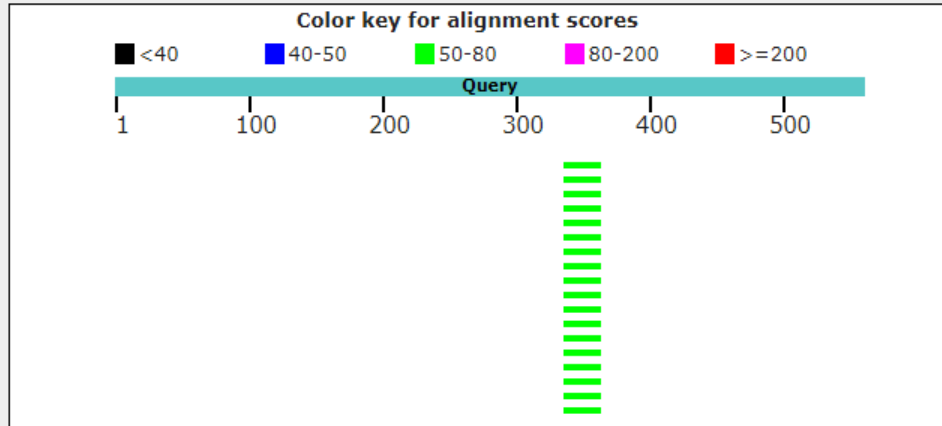
Alignments Download GenBank Graphics Distance tree of results

	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	<a href="#">Streptomyces aurantiacogriseus</a> gene for 16S rRNA, partial sequence, strain: NBRC 13668	2584	2584	98%	0.0	99%	<a href="#">AB184451.2</a>
<input type="checkbox"/>	<a href="#">Streptomyces</a> sp. B7 16S ribosomal RNA gene, partial sequence	2579	2579	98%	0.0	99%	<a href="#">AF332543.1</a>
<input type="checkbox"/>	<a href="#">Streptomyces luridus</a> strain YIM 131011 16S ribosomal RNA gene, partial sequence	2555	2555	96%	0.0	99%	<a href="#">KX502956.1</a>
<input type="checkbox"/>	<a href="#">Streptomyces</a> sp. W26 16S ribosomal RNA gene, partial sequence	2553	2553	99%	0.0	99%	<a href="#">EU596428.1</a>
<input type="checkbox"/>	<a href="#">Streptomyces</a> sp. CPE275 16S ribosomal RNA gene, partial sequence	2536	2536	98%	0.0	99%	<a href="#">JN969009.1</a>
<input type="checkbox"/>	<a href="#">Streptomyces setonensis</a> strain HBUM82838 16S ribosomal RNA gene, partial sequence	2536	2536	98%	0.0	99%	<a href="#">EU841570.1</a>
<input type="checkbox"/>	<a href="#">Streptomyces cyaneogriseus</a> strain 179 16S ribosomal RNA gene, partial sequence	2534	2534	98%	0.0	99%	<a href="#">EF063457.1</a>
<input type="checkbox"/>	<a href="#">Streptomyces kurssanovii</a> strain ICTA165 16S ribosomal RNA gene, partial sequence	2531	2531	98%	0.0	99%	<a href="#">EU841570.1</a>
<input type="checkbox"/>	<a href="#">Streptomyces</a> sp. YH4 16S ribosomal RNA gene, partial sequence	2531	2531	98%	0.0	99%	<a href="#">EU841570.1</a>

Questions/comments

Distribution of the top 18 Blast Hits on 18 subject sequences

Mouse over to see the title, click to show alignments



Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

Alignments Download GenBank Graphics Distance tree of results

Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/> <a href="#">Mycobacterium marinum ATCC 927 DNA, complete genome</a>	52.8	52.8	5%	0.011	100%	<a href="#">AP018496.1</a>
<input type="checkbox"/> <a href="#">Mycobacterium pseudoshottsii JCM 15466 DNA, complete genome</a>	52.8	52.8	5%	0.011	100%	<a href="#">AP018410.1</a>
<input type="checkbox"/> <a href="#">Actinosynnema pretiosum strain X47 chromosome, complete genome</a>	52.8	52.8	5%	0.011	100%	<a href="#">CP023445.1</a>
<input type="checkbox"/> <a href="#">Mycobacterium aurum isolate liquid genome assembly, chromosome:1</a>	52.8	52.8	5%	0.011	100%	<a href="#">LT899394.1</a>
<input type="checkbox"/> <a href="#">Cupriavidus sp. USMAA1020 chromosome 2, complete sequence</a>	52.8	52.8	5%	0.011	100%	
<input type="checkbox"/> <a href="#">Cupriavidus sp. USMAHM13 chromosome 2, complete sequence</a>	52.8	52.8	5%	0.011	100%	

Questions/comments

# Interpreting Results of **BLAST**

- Score:
- Max score:
- Total score:
- Query coverage
- E Value:
- Max Identity:

# Score

- Score: jumlah keselarasan semua segmen dari sekuens database yang cocok dengan sekuen nukleotida kita.
- Nilai skor menunjukkan keakuratan nilai penjajaran sekuens berupa nukleotida yang tidak diketahui dengan sekuens nukleotida yang terdapat di dalam GenBank.
- Semakin tinggi nilai skor yang diperoleh maka semakin tinggi tingkat homologi kedua sekuens
- Max score: Score of single best aligned sequence
- Total score: Sum of scores of all aligned sequences

# Query coverage

- *Query coverage* adalah persentasi dari panjang nukleotida yang sesuai dengan sekuen database yang terdapat pada BLAST

# Max identity

- *Max identity* adalah nilai tertinggi dari persentasi identitas atau kecocokan antara sekuen nukleotida dengan sekuen database yang tersejajarkan
- **Hagstrom *et al* (2000)** menyatakan bahwa bakteri yang mempunyai nilai max identity 16S rRNA lebih besar dari 97% adalah **spesies yang sama**.
- Sedangkan persamaan sekuen antara 93%-97% dapat mewakili identitas pada tingkat **genus** tetapi berbeda pada tingkat spesies

# What is an E-value?

- **E-value**
  - The chance that the match could be random
  - The lower the E-value, the more significant the match
    - $E = 10^{-4}$  is considered the **cutoff point**
    - $E = 0$  means that the two sequences are statistically **identical**

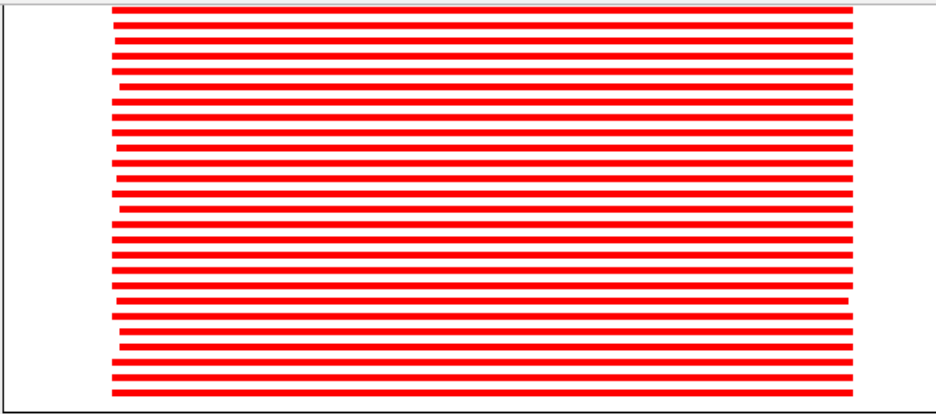


# E-value

- Nilai E-value merupakan nilai dugaan yang memberikan ukuran statistic yang signifikan terhadap kedua sekuen.
- Nilai E-value yang semakin tinggi menunjukkan tingkat homologi antara sekuen semakin rendah, sedangkan nilai E-value yang semakin rendah menunjukkan tingkat homologi antar sekuens semakin tinggi.
- Nilai E-value bernilai 0 (nol) menunjukkan bahwa kedua sekuen tersebut identik

# Accession Number

- Accession number (bioinformatics), a unique identifier given to a biological polymer sequence (DNA, protein) when it is submitted to a sequence database.



Descriptions

Sequences producing significant alignments:

Select: All None Selected:0

Alignments Download GenBank Graphics Distance tree of results

Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/> <a href="#">Streptomyces aurantiacogriseus</a> gene for 16S rRNA, partial sequence, strain: NBRC 13668	2584	2584	98%	0.0	99%	<a href="#">AB184451.2</a>
<input type="checkbox"/> <a href="#">Streptomyces sp. B7</a> 16S ribosomal RNA gene, partial sequence	2579	2579	98%	0.0	99%	<a href="#">AF332543.1</a>
<input type="checkbox"/> <a href="#">Streptomyces luridus</a> strain YIM 131011 16S ribosomal RNA gene, partial sequence	2555	2555	96%	0.0	99%	<a href="#">KX502956.1</a>
<input type="checkbox"/> <a href="#">Streptomyces sp. W26</a> 16S ribosomal RNA gene, partial sequence	2553	2553	99%	0.0	99%	<a href="#">EU596428.1</a>
<input type="checkbox"/> <a href="#">Streptomyces sp. CPE275</a> 16S ribosomal RNA gene, partial sequence	2536	2536	98%	0.0	99%	<a href="#">JN969009.1</a>
<input type="checkbox"/> <a href="#">Streptomyces setonensis</a> strain HBUM82838 16S ribosomal RNA gene, partial sequence	2536	2536	98%	0.0	99%	<a href="#">EU841570.1</a>
<input type="checkbox"/> <a href="#">Streptomyces cyaneogriseus</a> strain 179 16S ribosomal RNA gene, partial sequence	2534	2534	98%	0.0	99%	<a href="#">EF063457.1</a>
<input type="checkbox"/> <a href="#">Streptomyces kurssanovii</a> strain ICTA165 16S ribosomal RNA gene, partial sequence	2531	2531	98%	0.0	99%	<a href="#">EU841570.1</a>
<input type="checkbox"/> <a href="#">Streptomyces sp. YH4</a> 16S ribosomal RNA gene, partial sequence	2531	2531	98%	0.0	99%	<a href="#">EU841570.1</a>

Questions/comments

## TUGAS PRAKTIKUM ---- KUMPULKAN MINGGU DEPAN

1. Buat Lah langkah langkah dalam analisis sekuens dengan teknik BLAST
2. Interpretasikanlah hasil BLAST berdasarkan E-Value, Max score, Total Score, Max Identified and Query Coverage

THANK  
YOU



607132.wordpress.com

Noviani's Blog

