



## Finding data

### BLAST

Answer key

#### Contents

1. BLAST basics
  - E- value **new!**
2. How to use the tool and interpret its output?
3. Questions and practice exercises
  - Gene ID lookup **new!**

### 1. BLAST basics

The Basic Local Alignment Search Tool, or BLAST, is a way to compare sequences. It performs pairwise alignment between the *query* and the *targets* in the database to find a specific or similar sequences. Alignments with the best-matching sequences are shown and scored. BLAST similarity searches can be used to find homologous genes<sup>1</sup>, identify potential functions of novel genes, predict the size of PCR products, identify transcript or peptide evidence (from RNAseq or mass spectrometry experiments), and assign VectorBase and/or GenBank accession to genes among others.

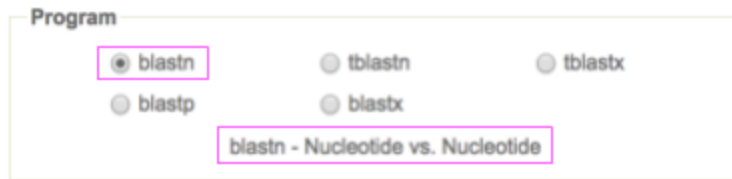
BLAST is actually a family of programs and these are:

Program	Query	Database
<b>blastn</b>	nucleotide	nucleotide
<b>blastp</b>	peptide	peptide
<b>blastx</b>	nucleotide (translated to peptide)	peptide
<b>tblastn</b>	peptide	nucleotide (translated to peptide)
<b>tblastx</b>	nucleotide (translated to peptide)	nucleotide (translated to peptide)

---

<sup>1</sup> The task of homolog finding is better accomplished using HMMER, <https://www.vectorbase.org/tutorials/tools-and-resources-tutorials/hmmer>

BLAST will *not* run if you choose an inappropriate program for the selected database, *e.g.*, blastn for a protein query. Hover with your mouse over each program and short description will describe you the type of query and database used in each one.



BLAST provides statistical evaluations that serve as a guide for the evaluation of the alignment scores.

BLAST algorithm is heuristic, thus BLASTn results are not always optimal due to the degeneracy of codons. Its strongly recommended to use BLASTp for homology searches, because protein searches are more specific than DNA searches. Translated BLAST, *i.e.*, tblastn, tblastx<sup>2</sup> and blastx, translates in all six reading frames. These programs are useful when the reading frame is unknown or the query is unknown. The image below shows the six potential reading frames of a sequence.

ATG → Met  
 TGA → Stop  
 GAC → Asp  
 5' - ATGACGAGAGAGCAGCCATTTTAG - 3'  
 3' - TACTGCTCTCTCGTCGGTAAAATC - 5'  
 Leu ← ATC  
 Stop ← AAT  
 Lys ← AAA

In the forward strand there are three reading frames, and its translations can start at position one, two or three and read in codons (in groups of three nucleotides at the time). Similarly happens in the reverse strand. Which one is the real reading frame? Possibly the one with no stop codons in the open reading frame (ORF). All six protein sequences, from the six reading frames, are used as a query against the target database.

BLASTx takes a DNA query, translates it into a protein sequence and then searches a protein database. tBLASTn translates an entire nucleotide database, in all six reading frames, to create a temporary protein database to search with the protein query. tBLASTx takes a translated DNA query to search a translated DNA database.

<sup>2</sup> tblastx is computationally intensive and should be used only as a last resource. Searching with large genomic queries is not recommended.

To perform a BLAST you can use a single or multiple query sequences, likewise, your target can be single, multiple or all databases, from one or all species. What you cannot do is to search using a mixture of nucleotide and peptide queries and/or databases. The allowed searches take place independently and the results will be reported for each database separately.

The screenshot shows the 'Datasets' section on the left and the 'Nucleotide' section on the right. In the 'Datasets' section, 'All Datasets' is unchecked. Under the 'Aedes' species dropdown, 'All Aedes' is unchecked, 'Aedes aegypti' is checked, and 'Aedes albopictus' is unchecked. The 'Anopheles' species dropdown is also visible. In the 'Nucleotide' section, 'Contigs Liverpool Strain' and 'Contigs Aag2 cell line str' are unchecked. 'ESTs Publicly submitted' is checked. 'Scaffolds Liverpool strain and TRF.' and 'Scaffolds Aag2 cell line s and TRF.' are checked. 'Transcripts Liverpool stra' is unchecked.

Target single or multiple datasets in a single species

The screenshot shows the 'Datasets' section on the left and the 'Nucleotide' and 'Peptide' sections on the right. In the 'Datasets' section, 'All Datasets' is checked. Under the 'Aedes' species dropdown, 'All Datasets' is checked. Under the 'Anopheles' species dropdown, 'All Datasets' is checked. Under the 'Glossina' species dropdown, 'All Datasets' is checked. Under the 'Ixodes' species dropdown, 'All Datasets' is checked. In the 'Nucleotide' section, 'Contigs', 'ESTs', 'Scaffolds', 'Transcripts', 'Chromosomes', and 'Assembled transcriptome' are all unchecked. In the 'Peptide' section, 'Peptides' and 'Assembled proteome' are both checked.

Target single or multiple datasets in multiple species

There are different default options for each BLAST program. The following is a brief description of each one.



The image shows two side-by-side screenshots of BLAST options panels. The left panel has a title 'Options' and contains two settings: 'Maximum E-Value' with a dropdown menu showing '10' and 'Word Size' with a dropdown menu showing '11'. The right panel also has a title 'Options' and contains two settings: 'Complexity Masking' with a dropdown menu showing 'On' and 'Results per Query per Database' with a dropdown menu showing '10'.

#### *E-value:*

Is the number of unrelated sequences in a similarity search of a sequence database that are expected to achieve a local alignment score as high or higher than the one obtained between the query sequence and the matching database sequence. Like with a p-value, the lower the better. Recently, new options for higher e-values were added. *Previously*, our highest e-value allowed was 10. **Current options include 100, 1000, 10000 and 100000.** You can test this new feature when using protein queries against scaffolds in tBLASTn for the search of (new) highly derived/rapidly evolving genes.

#### *Word size:*

The query will be divided up into smaller chunks which will then be compared to the databases(s) separately. Default is 3 for DNA and 11 for protein. A larger words is more efficient proving an exact match.

#### *Scoring Matrix:*

This option is available for proteins searches, *i.e.*, all programs except blastn, and are used to evaluate the quality of the alignment with a score, in which the higher the score the better. By default BLOSUM62 is selected (and its calculated from comparisons of sequences were 62% of the similar ones are grouped), but it can be changed depending on the type of sequences you are searching with:

- *BLOSUM* is based on local multiple sequence alignment (MSA) of distantly related proteins. Higher numbers mean smaller evolutionary distance (*i.e.*, higher sequence identity). BLOSUM62 is used for closer sequences than BLOSUM45
- *PAM* is based on global MSA of closely related proteins. Higher numbers mean larger evolutionary distance. PAM70 is used for more distant sequences than PAM30

#### *Complexity Masking:*

It removes segments of the query sequence that have low compositional complexity. Filtering leaves the more biologically interesting regions of the query sequence available for specific matching against database sequences. Filtering low-complexity regions from the query sequence helps to reduce the number of false positives.

#### *Number of results:*

Limits the output of BLAST to a maximum of the specified numbers of hits and can potentially speed up a search. Default is 10.

## 2. How to use the tool and interpret its output?

### *Quick start*

*Optional:* Login to VectorBase. Your BLAST (ClustalW and Hmmer) jobs will be saved and viewable in your user page.



Go to BLAST (from the Tools navigation menu or from the organism pages)

Paste or upload the query sequence(s)

Perform Search of target database(s)

Are matching sequences found with reasonable alignments and significant E values for alignment scores? Depending on your query 100% identical matches or sequences with mismatches but clearly related are of interest.

### *Output interpretation*

The alignments should be examined for a small value of E, absence of low-complexity regions that falsely give high alignment scores, and quality of the alignment (*i.e.*, presence of long stretches of aligned regions that do not depend on gaps to produce an alignment). The results are divided into three sections, shown below as a, b, c and d.

**Results**  
**Job 180005**

Description Test  
Submitted Thursday, April 27th, 2017 09:59:07 -0400  
Compute Time 11 seconds

**a.** CLEAR RESULTS  
EXPAND ALL

Checked Hits

**b.** Quick align Pass to ClustalW Download  
include query

Organism	Database	HSPs
Aedes albopictus	(Peptides) C6/36 cell line strain peptide sequences, NCBI-101 geneset.	10
Aedes albopictus	(Peptides) Foshan strain peptide sequences, AaLoF1.2 geneset.	10

**c.**

**d.** Show Query/Hit Numbers

Hit	Gene Name	Description	Query	Aln Length	E-value	Score	Identity	Query Hit	DB Sequence Hit
<input type="checkbox"/>	AALF017696-PA	long wavelength sensitive opsin	AaGPRop1	320	0.0	1453	75.9%	>	>
<input type="checkbox"/>	AALF009534-PA	Rhodopsin	AaGPRop1	320	0.0	1485	79.6%	>	>
<input type="checkbox"/>	AALF009531-PA	long wavelength sensitive opsin	AaGPRop1	323	5e-179	1302	66.7%	>	>

#### Section a:

- 'Clear Results' keeps input sequence and parameters, unlike Reset, which also deletes the query.
- 'Expand All (/Sort all)': Allows sorting across all organism hits
- 'View Raw Results' shows data in plain format on a new tab.
- 'Download Results' downloads raw results in a text file.

#### Section b: It only works when hits are selected using the tick box next to them

- 'Quick align' in the same page a small box shows a Clustal alignment.
- 'Pass to ClustalW' a new tab is open with the sequences in FASTA format on VectorBase ClustalW.
- 'Download' text file with sequences in FASTA format.

#### Section c:

- 'HSPs' or high-scoring segment pair, is a high-scoring word or alignment between the query and the hits. In protein searches, the HSP is usually three amino acids long. The word will be enlarged if sequence positions neighboring the word also match and provide a higher-scoring alignments.

#### Section d:

Results are presented in a table format, with each column providing different lines of evidence to select the best hit(s).

Use only high confidence hits, i.e., genomic regions identified with high similarity with your query, for further analysis. The confidence in a hit can be judged by the:

E-value    The smaller the better

Score  
Identity  
Length }    The bigger the better

which should be considered jointly.

$$\text{Identity} = \frac{\text{\# of identical letters}}{\text{alignment length (\# matching letters + \# gaps)}}$$

When looking for homologous sequences, only hits with close to 100% identity and approximately the same length as the query should be retained.

For example, hit A is more indicative of common ancestry than hit B, even though the E-value is similar in the two hits.

Hit	E-value	Identity	Alignment Length
A	0.005	50%	100% (whole query)
B	0.006	80%	10% of the query

For example, a 50% identity of a match that extends for the whole protein sequence is more indicative of common ancestry than a 80% identity over a region = 10% of the length of one of the two proteins, without significant similarity in the other 90% region.

In summary. In the BLAST output there are exact matches and sequences with mismatches. For all the hits that are not exact matches, there are calculations that generate statistics that help evaluate which hits are significant.

### 3. Questions and practice exercises

#### Question 1. 1

What is BLAST?

	True	False
A program to construct sequence alignments		X
The BLAST algorithm is a way to perform a sequence similarity search	X	

A program to do small and big scale data mining queries such as the ones done with VectorBase Search or Biomart		X
---	--	---

### Question 1. 2

What are the five types of VectorBase BLAST programs?

BLASTn, tBLASTn, tBLASTX, BLASTp and BLASTx
---

### Question 1.3

When planning a PCR experiment, it is necessary to know the size of the amplicon products for the gDNA and cDNA, to interpret the obtained results in the agarose gel. This requires a knowledge of the gene structure to localize where the primers will bind.

Using the primer sequences (in sample file provided in the tutorial page), run a BLASTn job for *Ixodes scapularis* 'Scaffolds' and 'Transcripts'.

## BLAST

### Basic Local Alignment Search Tool

```
>Serp2Ae1-6_F
.....
TTACGCTCCCGACGTTATTC

>Serp2Ae1-6_R
.....
TTCGAGGGATCAAACAGGTC
```

In the BLAST results click on the "Transcripts" database. Are the best hits binding to the gene in the table below?

**Table 2.** PCR primer sequences used in this study.

Gene	Primer	5'→3'	Amplicon size	Genbank /Vectorbase
Serp1n 2	Serp2Ae1-6_F	TTACGCTCCCGACGTTATTC	651	NW_002630218.1 IscW_ISCW018607
	Serp2Ae1-6_R	TTCGAGGGATCAAACAGGTC		

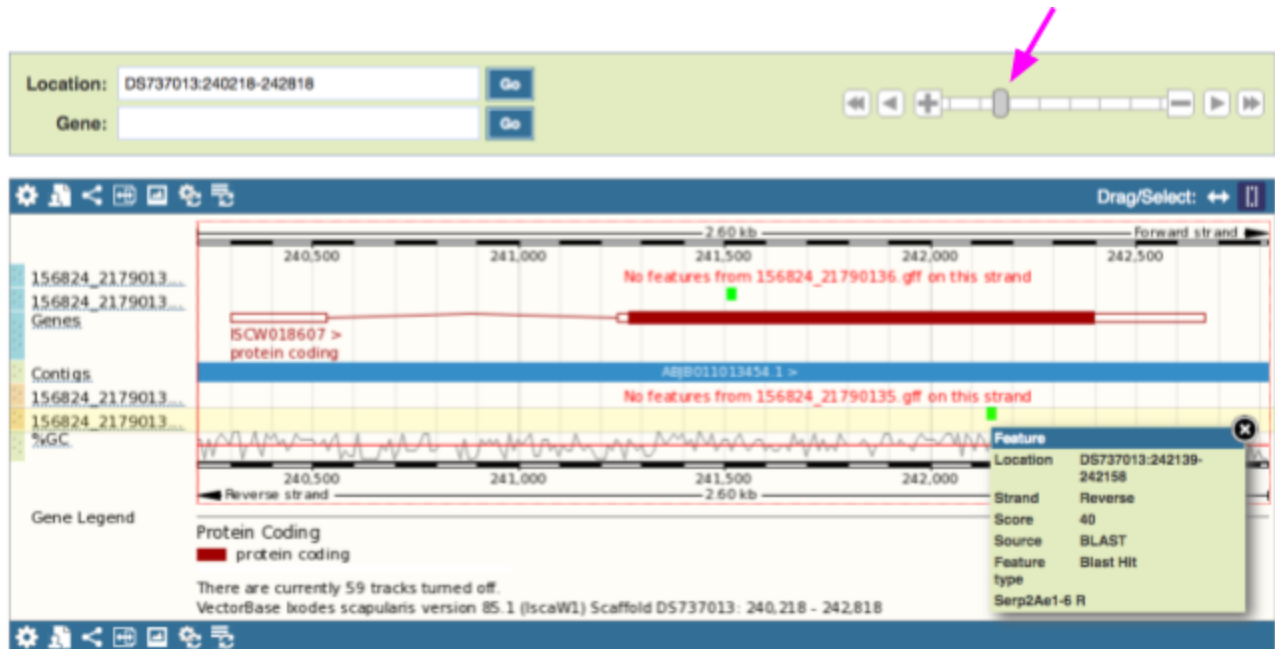
Yes <u>X</u> No <u>      </u>
-------------------------------



In the BLAST results click on the “Scaffold” database. Click on the first hit. In the popup window click on “Browse Genome”. What is the green bar in the bottom display?

	True	False
The gene sequence		X
The scaffold sequence		X
The primer (query) sequence	X	

Repeat previous step but for the second hit. Move the zoom slide bar to 2,600 bp. Click on the green squares to obtain popup windows with details.



#### Question 1. 4

Anopheline mosquitoes are of prime medical importance as vectors of human *Plasmodium*, a parasite that cause malaria, yet *Anopheles* phylogeny is poorly known and there are many species complexes. In 2001, Krzywinski and colleagues published a study where they analyze sequences of two protein-coding single copy genes from Anophelines, *glucose-6-phosphate dehydrogenase* (*G<sub>6</sub>pd*) and *white*. These genes sequences were used in 16 species to determine their phylogenetic relationships.

Syst Biol. 2001 Aug;50(4):540-56.

**Toward understanding Anophelinae (Diptera, Culicidae) phylogeny: insights from nuclear single-copy genes and the weight of evidence.**

Krzywinski J<sup>1</sup>, Wilkerson RC, Besansky NJ.

In this exercise, you are going to use BLASTp to identify their VectorBase gene IDs in *An. stephensi*, using the GenBank accession numbers provided in the paper 'Material and Methods', AF317808 and AF318208 for *G<sub>6</sub>pd* and *white*, respectively. There is a file with the GenBank sequences in the tutorial page. **Hint:** there is no need to run BLAST twice, you can run a single job with both sequences

## BLAST

### Basic Local Alignment Search Tool

```
>g6pd_Astephensi
LWWLFRDNLLPSDTKFIGYARSKLSVAELKEKCRQY
KGWNRIIVEKPFGRDAESSNVLSVHLAKLFTEDQLY

>white_Astephensi
VQMLAIAPAKEAECRDMIKKICDSFAVSPiAREVLE
LDQDGVMNINGSIFLFLTNMTFQNVFAVINVFAEI
SCASSISMAISVGPPVVIPFLI
```

	<i>Anopheles stephensi</i>	
	SDA-500 strain	Indian strain
<i>white</i>	ASTE001213	ASTEI04427
<i>G<sub>6</sub>pd</i>	ASTE009404	ASTEI04175

Do not close this job results, you need them for the next question.

Are *An. stephensi* *G<sub>6</sub>pd* and *white* sequences from GenBank the same length when compared to VectorBase gene models? **Hint:** look at query and hit graphics (and/or numbers).

▼ Database

(Peptides) Indian strain predicted peptide sequences, AsteI2.2 geneset. 20

(Peptides) SDA-500 strain predicted peptide sequences, AsteS1.3 geneset. 20

Quick align Pass to ClustalW Download

☒ include query

Gene Name Description Query Aln Length E-value Score Identity Query Hit DB Sequence Hit

white\_Asteph... 243 1e-175 1751 100% > >

g6pd\_Asteph... 153 5e-110 1131 100% > >

Show Query/Hit Numbers

Show Query/Hit Graphics

Query	Query	Hit	Hit
Start	End	Start	End
1	243	352	594
1	153	122	274

No they are not, the ones from Genbank are shorter/incomplete gene sequences.

### Question 1.5

A file with the VectorBase gene sequences, AALB010162 and AALB006905, is provided in the tutorial page. Run a separate BLASTn job for each query, changing to '50' the 'Results per query per database' and select Assembled transcriptome (RNAseq) as the target database. Is there transcript evidence for *An. albimanus*, *G6pd* and *white* genes? **Hint:** Remember to check alignment length, e-value, score and identity.

There is transcript evidence, with statistical confidence, only for *G6pd*.

What the blue and green colors mean in the bars of the query and hit graphics?

The direction of the transcript, forward is green and reverse is blue.

What is the difference between ‘assembled transcriptomes’ and ‘transcripts’?

Transcripts are *in silico* predictions from the genome sequences. The ‘assembled transcriptomes’ is experimental evidence, in the case of *An. albimanus* from ESTs and RNAseq experiments.

### Question 1.6

Similarity or identification searches can also be done using VectorBase data (<https://www.vectorbase.org/downloads>). Commercial parties have developed enhanced BLAST applications, but its use is not free.

To study the effect of Dengue-2 virus infection in the salivary glands of *Ae. aegypti* mosquitoes, scientist performed an experiment to identify differential expression in this organ. Complete the experimental strategy shown below with VectorBase data.

- Treatment groups: mosquitoes feed with (1) blood and (2) blood + DEN-2
- After a 10-day extrinsic incubation period, the salivary glands were dissected
- Proteins were extracted from the sample tissue and run on a two-dimensional gel electrophoresis
- Differentially expressed proteins (spots) were cut from the gel and used for mass spectrometry
- Protein identification was performed using the commercial software x. Using BLAST the protein spectra were compared against

**Aedes-aegypti-Liverpool\_PEPTIDES\_AegL3.3** (target database)

downloaded from VectorBase.

- Results are reported in the paper using as stable accession

**VectorBase gene IDs**

- For genes with unknown symbol or description GenBank non-redundant protein sequences (nr) was used for a second run of similarity searches.

### Question 1.7 Gene ID lookup **new!**

In conjunction with inputting sequence(s) through either the text area or file upload in Blast, a new feature called **Gene ID lookup** in Blast is introduced. In lieu of inputting sequences, one can input Gene ID(s) through either Textarea or File upload. To distinguish between Nucleotide (cds type) and Amino-Acid (aa type) sequences, the first line of the input should be either **#GeneID\_cds** (for Nucleotide sequences) or **#GeneID\_aa** (for Amino-Acid

sequences): please refer to below figures. Note that one should select correct Program to run Blast depending on \_cds or \_aa. After inputting Gene IDs and submitting Blast job, sequences corresponding to Gene ID input(s) are retrieved and Blast will be automatically executed based on the sequences.

### Example #1) \_cds input

**BLAST**  
Basic Local Alignment Search Tool

#GeneID\_cds  
AGAP005203  
AGAP005203-RA  
AGAP005203-RF  
AGAP005203-RJ  
GAPW01000096.1

Upload FASTA File  
 No file selected.

Program

☒ blastn      ☐ tblastn      ☐ tblastx  
☐ blastp      ☐ blastx

blastn - Nucleotide vs. Nucleotide

**Example #2) \_aa input**

**BLAST**

**Basic Local Alignment Search Tool**

#GeneID\_aa  
AGAP005203  
AGAP005203-PA  
AGAP005203-PD  
AGAP005203-PH  
DEK\_2005\_000013

**Upload FASTA File**

No file selected.

**Program**

☐ blastn      ☐ tblastn      ☐ tblastx  
☒ blastp      ☐ blastx

blastp - Peptide vs. Peptide

If you need help with any question and its answer contact us at [info@vectorbase.org](mailto:info@vectorbase.org). Because VectorBase data, tools and resources are updated every two months (6 release cycles per year), answers to these exercises will change too.