Ria Lamba

Bioinformatics Project Draft

# Identification of Spider Latrotoxins Using Diamond and NCBI-BLAST

**Video Tutorial:** [https://youtu.be/dEykHyQTfZ0](https://youtu.be/dEykHyQTfZ0)

**Datasets Used in Tutorial:**

[https://drive.google.com/drive/folders/1XrOYhCSlVW7oy5RVuTRS6YywRgAHPh1O?usp=sharing](https://drive.google.com/drive/folders/1XrOYhCSlVW7oy5RVuTRS6YywRgAHPh1O?usp=sharing)

Ria Lamba

Bioinformatics Project Draft

**Background**

In this tutorial, we will be using Diamond and NCBI-BLAST on the command line to help us identify potential latrotoxin sequences in the Mediterranean Black Widow. We will be running these jobs on the HPCC UCR cluster, but the commands can be used on other servers. Diamond and NCBI-BLAST are both programs that run various BLAST searches. While both programs give the same BLAST output, they each have different search algorithms and capabilities, which is why one may need to be used over the other for specific situations. While Diamond is faster than NCBI-BLAST, it can only do two types of BLAST searches: BLASTx and BLASTp, whereas NCBI-BLAST can also do BLASTn, tBLASTx and tBLASTn. Both programs can be used to generate BLAST databases and BLAST outputs in different formats.

BLAST stands for Basic Local Alignment Search Tool ("BLAST: Basic Local Alignment Search Tool.", n.d.). It is used to align and compare nucleotide or protein sequences to one another and to sequence databases. There are different types of blast searches that can be done for query sequences versus subject sequences.

| Blast Type | Query | Subject |
|---|---|---|
| blastp | Protein | Protein |
| blastx | Translated Nucleotide | Protein |
| blastn | Nucleotide | Nucleotide |
| tblastx | Translated Nucleotide | Translated Nucleotide |
| tblastn | Protein | Translated Nucleotide |

BLAST helps scientists get insight on sequence similarities that can be used to understand relationships between sequences among species and to identify gene families. Gene

Ria Lamba

Bioinformatics Project Draft

families are groups of similar genes formed by gene duplication. While one can do BLAST searches on the NCBI website, it is not an efficient way to BLAST multiple sequences, especially if you have a file with hundreds or thousands of sequences. Therefore, it is important to learn to run BLAST searches on the command line. If you connect to a high-performance computing cluster, you will have access to programs without having to download them on your own computer. You can also run bigger jobs with ease by submitting scripts with your commands, that then get sent to a queuing system to be processed. This tutorial shows how to submit BLAST jobs on a computing cluster using BASH scripts.

**Dataset**

Latrotoxins are a type of neurotoxin found in spiders of the genus *Latrodectus*. These spiders are incredibly poisonous and a bite from them can cause the illness latrodectism if enough venom is injected into the victim (Bonnet, 2004). The Mediterranean Black Widow (*Latrodectus tredecimguttatus*) is one of these toxic spiders, commonly found throughout the Mediterranean region as the name suggests. In this tutorial, we will perform BLAST and Diamond searches for latrotoxin-like sequences in the Mediterranean Black Widow, using the spider's venom gland RNA sequences downloaded from NCBI's Sequence Set Browser (bio sample SAMN02318955). We will also be using a FASTA file containing the first 320 amino acids of several previously described latrotoxins. This FASTA file was made by compiling the sequences of several latrotoxins collected and found on NCBI, followed by trimming after the first 320 amino acids. This FASTA file will be used to find potential latrotoxins in the Mediterranean Black Widow.

Ria Lamba

Bioinformatics Project Draft

First, we will do a BLASTx search using Diamond with the venom gland RNA sequences as the query against nr, the non-redundant protein database from NCBI. The BLASTx results will identify sequences that resemble the venom RNA sequences. Second, we will use NCBI-BLAST makeblastdb to make a BLASTable database out of the venom gland RNA sequences that will allow us to specifically search for latrotoxins; we will run a tBLASTn search with NCBI-BLAST using our FASTA file containing the first 320 amino acids of known latrotoxins as a query against these venom gland RNA sequences. The tBLASTn search will be repeated with the output format 0, a format that gives a visual of the alignments between each query and subject. The BLASTx and tBLASTn results give us different information that can be useful for identifying sequences of interest. The BLASTx results include protein sequences from nr with descriptive titles, that will help give us context as to what each venom gland sequence's function and role may be. The tBLASTn results show us the venom gland sequences that match to known latrotoxin sequences. We then will know which venom gland sequences are likely to be toxins.

Ria Lamba

Bioinformatics Project Draft

# Instruction Manual Table of Contents

Ria Lamba

Bioinformatics Project Draft

# Instruction Manual

 First open a terminal. In this tutorial, we will be using MobaXTerm on a Windows computer. However, you can access the cluster using any terminal you want. To log on to the cluster, type the server name into the search bar (for example, cluster.hpcc.ucr.edu). You will then be prompted to type in your username and then your password.

If you are using another terminal such as git-bash, type the following to log in:
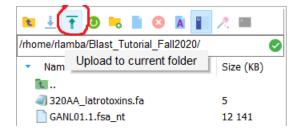
ssh -X rlamba@cluster.hpcc.ucr.edu

If it is your first time logging on to the cluster, you will be prompted to enter your password.

**Make** a directory for this tutorial and **cd** into it:

mkdir Blast_Tutorial_Fall2020
cd Blast_Tutorial_Fall2020
To upload files from your computer to the directory on the cluster, move over to the graphical interface on the left and click upload to current folder. Then select the files on your computer that you would like to upload.



Alternatively, transfer files using the command line using **scp**:

scp -r <path_to_directory-in-pc> <username>@<host_name>:<path-to-file-on-cluster>

**scp** stands for "secure copy". In <path_to_directory> type the path that the file you want to upload is located at in your computer. Then add your username on the cluster in <username>. Enter the name of the cluster you are connected to in <host_name>. In <path-to-file-in-cluster> type the path in the cluster that you want to upload the file to.

## BLASTx Search With Diamond

On the cluster we are using, bioinformatics software we need are already installed and can be accessed via modules. To load DIAMOND on the cluster, type the following:

module load diamond

Next, we want to open a file to write our script, we will use the text editor **nano**:

Ria Lamba

Bioinformatics Project Draft

nano tredecimguttatus_diamond_blastx.sh

Your script may vary depending on the cluster and queuing system you are using. Here is an example script:

```
#!/bin/bash -l

#SBATCH --nodes=1
#SBATCH --ntasks=6
#SBATCH --mem-per-cpu=6G
#SBATCH --time=1-00:00:00
#SBATCH --output=/rhome/rlamba/Blast_Tutorial_Fall2020/tredecimguttatus
_DIAMOND_BLASTx_runlog.out

module load diamond

diamond blastx \
-q /rhome/rlamba/Blast_Tutorial_Fall2020/GANL01.1.fsa_nt \
-d /rhome/rlamba/shared/DIAMOND_v0.9.24_nr_db_Feb2020.dmnd \
-o /rhome/rlamba/Blast_Tutorial_Fall2020/tredecimguttatus_nr_database.tsv \
-k 3 -f 6 qseqid sseqid evalue qframe qlen slen qstart qend sstart send length pident stitle
```

The first line in the script is called the sha-bang line. It is crucial that every script starts with this line, otherwise it cannot run. It is the path to the Bash Interpreter.

The next five lines that start with #SBATCH specify what resources are being requested for the job. In this case, we are asking for 1 node, 6 GB of CPU, and 1 day to be allocated for this job. This does not necessarily mean the job will take up all 6 GB or 1 whole day. It just means this is the maximum that will be given for the job. If you request too little resources, your job will immediately stop running once it runs out of resources. If you request a large amount of resources, your job may remain pending in the queue for a while. This is why one must be reasonable when requesting resources for a job.

**Diamond Parameters Used In Script:**

| | |
|---|---|
| -q | Path to the file we are using for the query |
| -d | Path to the database we are blasting our query against (the subject) |
| -o | Path to where we want our output file to be when created + the name we want for the file. |
| -k | This parameter is optional. It specifies how many matches you want for each query sequence. If you do not use the -k parameter, you will get the default number of matches (25). |

Ria Lamba

Bioinformatics Project Draft

| -f | This is the output format line, where you specify what output format you would like (in this case, we have 6 to specify output format 6) and what information you want in your output. |
|---|---|
| qseqid | Query sequence id |
| sseqid | Subject sequence id |
| evalue | E-value (Value that gives an idea of how good the match is. The lower the e-value, the better the match is.) |
| qframe | The frame that the query is translated in. |
| qlen | How long the alignment is for the query sequences (will be in bp since query is nucleotide sequences) |
| slen | How long the alignment is for the subject sequence (Will be in amino acids since the subject is protein sequences) |
| qstart | When the alignment starts for the query sequences. |
| qend | When the alignment ends for the query sequence. |
| sstart | When the alignment starts for the subject sequences. |
| send | When the alignment ends for the subject sequence. |
| length | Length of the alignment |
| pident | Percent of identical matches |
| stitle | The titles of the subject sequences. Note that these will be more descriptive when doing blast searches against databases like nr, which has been put together by ncbi. |

Once you have finished your script, submit it to the queuing system:

sbatch tredecimguttatus_diamond_blastx.sh

Your job should be finished running in a few hours. You can check the status of jobs you have submitted by typing:

squeue -u rlamba

When you use **squeue -u** to check your job status, you should see something like this if it is running correctly:

JOBID PARTITION    NAME    USER ST      TIME  NODES NODELIST(REASON)
       2537826    intel tredecim  rlamba  R    2:40:12      1 i23

If you do not see any jobs running, that means that your job was submitted but there is an issue that caused it to not run. Check the *runlog.out* file to see what may have gone wrong and check your script for errors such as extra spaces.

## Making a BLAST Database With NCBI-BLAST

To load NCBI-BLAST, type:

module load NCBI-BLAST/2.9.0+

Ria Lamba

Bioinformatics Project Draft

To make a database of the *tredecimguttatus* sequences, type the following:

makeblastdb -dbtype nucl -in /rhome/rlamba/Blast_Tutorial_Fall2020/GANL01.1.fsa_nt -out /rhome/rlamba/Blast_Tutorial_Fall2020/Tredecimguttatus_database_nucdb

**Makeblastdb Parameters Used In Script:**

| dbtype | What type of database are you making. In this case, we type nucl for nucleotide. If we were making a protein database, we would type prot. |
|--------|-------------------------------------------------------------------------------------------------------------------------------------------|
| -in    | Stands for input. Path to the file you want to make a database out of.                                                                     |
| -out   | Stands for output. Path to where you want the database to be when it is created + the name you want for your database.                     |

The database should be done in approximately a second. It is a very fast process. Note that when you check your output files with **ls** you will see three files:

-Tredecimguttatus_database_nucdb.nsq
-Tredecimguttatus_database_nucdb.nin
-Tredecimguttatus_database_nucdb.nhr

These are all components of the database. To do a BLAST search against the database, simply refer to the database as *Tredecimguttatus_database_nucdb*.

**tBLASTn With NCBI-BLAST**

Open another text file using **nano**:

nano Tredecimguttatus_Latrotoxin_tblastn.sh

Here is an example script:

#!/bin/bash -l

#SBATCH --nodes=1
#SBATCH --ntasks=6
#SBATCH --mem-per-cpu=6G
#SBATCH --time=1-00:00:00
#SBATCH --output=/rhome/rlamba/Blast_Tutorial_Fall2020/tredecimguttatus _Latrotoxin_tblastn_runlog.out

module load NCBI-BLAST/2.9.0+
tblastn \
-query /rhome/rlamba/Blast_Tutorial_Fall2020/320AA_latrotoxins.fa \
-out /rhome/rlamba/Blast_Tutorial_Fall2020/tredecimguttatus _Latrotoxin_tblastn.tsv \

Ria Lamba

Bioinformatics Project Draft


-db /rhome/rlamba/Blast_Tutorial_Fall2020/Tredecimguttatus_database_nucdb \
-outfmt '0 qseqid sseqid evalue qframe qlen slen qstart qend sstart send length pident stitle' \
-max_target_seqs 5

**NCBI-BLAST Parameters Used In Script:**

| | |
|---|---|
| -query | Path to the file we are using for the query |
| -out | Path to where we want our output file to be when created + the name we want for the file. |
| -db | Path to the database we are blasting our query against (the subject) |
| -max_target_seqs | This parameter is optional. It specifies how many matches you want for each query sequence. If you do not use the -max_target_seqs parameter, you will get the default number of matches (25). |
| -outfmt | This is the output format line, where you specify what output format you would like (in this case, we have 6 to specify output format 6) and what information you want in your output. |
| qseqid | Query sequence id |
| sseqid | Subject sequence id |
| evalue | E-value (Value that gives an idea of how good the match is. The lower the e-value, the better the match is.) |
| qframe | The frame that the query is translated in. |
| qlen | How long the alignment is for the query sequences (will be in amino acids since query is protein sequences) |
| slen | How long the alignment is for the subject sequence (Will be in bp since the subject is nucleotide sequences) |
| qstart | When the alignment starts for the query sequences. |
| qend | When the alignment ends for the query sequence. |
| sstart | When the alignment starts for the subject sequences. |
| send | When the alignment ends for the subject sequence. |
| length | Length of the alignment |
| pident | Percent of identical matches |
| stitle | The titles of the subject sequences. Note that these will be more descriptive when doing blast searches against databases like nr, which has been put together by ncbi. |


**tBLASTn In Output Format 0:**

The BLASTx search and tBLASTn search we did previously were both in output format 6. We can easily change the output format. Different blast output formats show us the same information in different ways, giving us different insights. We will be repeating our last tBLASTn search, only this time we will ask the program to give us output in format 0 instead of format 6. While

10

Ria Lamba

Bioinformatics Project Draft

format 6 is a table with all the terms specified in the *-outfmt* line as headers, format 0 gives the same information specified, but also gives a visual of the alignments between each query and subject. See pages 13 and 14 to see screenshots of these two output formats.

Open another text file using **nano**:

nano Tredecimguttatus_Latrotoxin_tblastn.sh

Here is an example script:

```
#!/bin/bash -l

#SBATCH --nodes=1
#SBATCH --ntasks=6
#SBATCH --mem-per-cpu=6G
#SBATCH --time=1-00:00:00
#SBATCH --
output=/rhome/rlamba/Blast_Tutorial_Fall2020/tredecimguttatus_Latrotoxin_F0_tblastn_runlog.
out
module load NCBI-BLAST/2.9.0+
tblastn \
-query /rhome/rlamba/Blast_Tutorial_Fall2020/320AA_latrotoxins.fa \
-out /rhome/rlamba/Blast_Tutorial_Fall2020/tredecimguttatus_Latrotoxin_F0_tblastn.tsv \
-db /rhome/rlamba/Blast_Tutorial_Fall2020/Tredecimguttatus_database_nucdb \
-outfmt '0 qseqid sseqid evalue qframe qlen slen qstart qend sstart send length pident stitle' \
-max_target_seqs 5
```

Notice that the script is the same as the script for our tBLASTn in format 6, only at the *-outfmt* line, there is a 0 in place of the 6. An F0 has also been added to the *runlog.out* file and the output file names in order to be able to distinguish them from the files for the previous tBLASTn search that was in format 6.

## Downloading and Viewing Output:

If you are using MobaXTerm, you can download files from the cluster to your computer by clicking download selected files (blue downward arrow at top left).

Ria Lamba

Bioinformatics Project Draft

You can also transfer files on the command line by typing:

scp -r <path-to-file-in-cluster> <username>@<host_name>:<path_to_directory-in-pc>
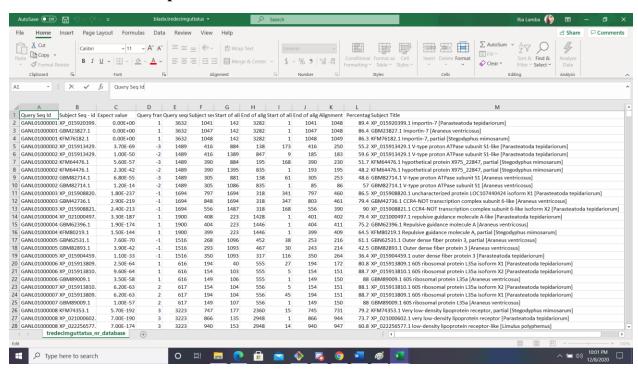
In <path-to-file-in-cluster> type the path that the file you want to download is located at in the cluster. Then add your username on the cluster in <username>. Enter the name of the cluster you are connected to in <host_name>. In <path_to_directory-in-pc> type the path to the directory you want the file to be downloaded to on your pc.

The output for these BLAST searches that are in format 6 are best viewed and analyzed in excel. You will notice several columns with various information. The order of the columns is the same as the order of the terms specified in the *-f/-outfmt* line. The output for the tblastn search in format 0 is best viewed and analyzed using notepad++ or a similar text editor. You will notice information for each query and subject match along with a visual of the sequence alignment between the two.

In this case we are interested in locating latrotoxins from these BLAST results.

-Open your BLASTx output with Excel.

**Diamond BLASTx Output:**
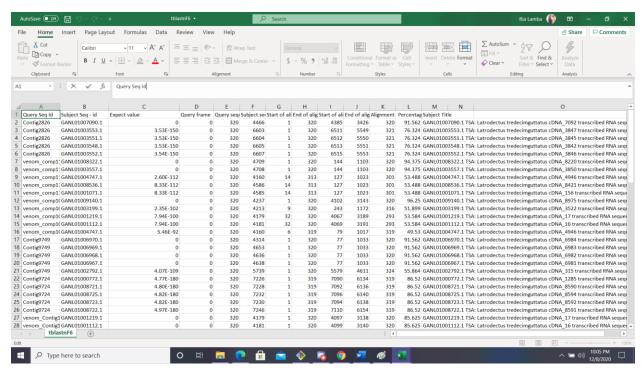


-Click find and select and type "latrotoxin"

-Click find all

Ria Lamba

Bioinformatics Project Draft

Excel will then show you all the subject titles that contain the word "latrotoxin". You can then see which of the query sequences in the venom gland matched to sequences in NCBI that are labeled as latrotoxins. This can give you an idea of which sequences in the venom gland may be latrotoxins that you should look at more closely.

-Open your tBLASTn format 6 output with Excel

## NCBI-BLAST tBLASTn Output In Format 6:



You can now see which sequences in the venom gland matched with our query file of latrotoxins. You can also take a look at the evalue to help you judge how good of a match there is between a certain query and subject. This also can give some insight into which sequences in the venom gland are latrotoxins.

If you want to view the alignments between the query and subject sequences, you can open and view your tBLASTn results in format 0.

Ria Lamba

Bioinformatics Project Draft

## NCBI-BLAST tBLASTn Output In Format 0:
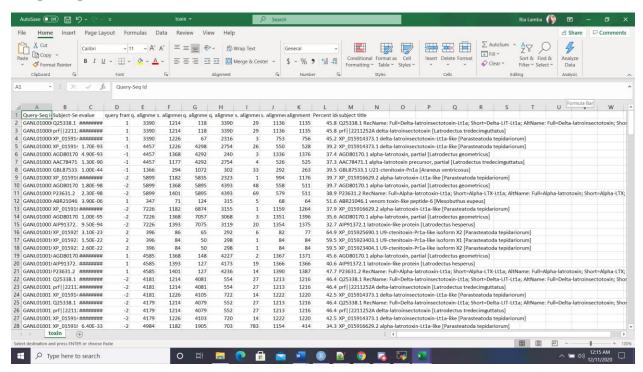
**grep Command to Pull Specific Rows:**

Another trick you can use to find matches of interest in your BLAST output is to use the **grep** command while you are still logged on to the cluster. In this example, we will use the command to find any lines in our output file with the word "toxin" and output them to a new file:

fgrep toxin tredecimguttatus_nr_database.tsv > toxin.tsv

The output file (toxin.tsv) will only contain rows that mention the word "toxin" in the line.

Ria Lamba

Bioinformatics Project Draft

# **<u>Works Cited</u>**

"BLAST: Basic Local Alignment Search Tool." *National Center for Biotechnology Information*,
    U.S. National Library of Medicine, blast.ncbi.nlm.nih.gov/Blast.cgi.


Bonnet, M.S. "The Toxicology of Latrodectus Tredecimguttatus: the Mediterranean Black
    Widow Spider." *Homeopathy*, No Longer Published by Elsevier, 9 Jan. 2004,
    www.sciencedirect.com/science/article/abs/pii/S1475491603001243.


"High-Performance Computing Center (HPCC)." *High-Performance Computing Center (HPCC)*
    *| HPCC @ UCR*, hpcc.ucr.edu/.


"Sequence Set Browser :: NCBI." *National Center for Biotechnology Information*, U.S. National
    Library of Medicine, www.ncbi.nlm.nih.gov/Traces/wgs/?view=TSA.


"Shebang." *Shebang - Linux Shell Scripting Tutorial - A Beginner's Handbook*,
    bash.cyberciti.biz/guide/Shebang.