

COMMUNITY SCIENCE PROJECT WORKFLOW (with example)

Introduction

How does exchange of genetic material between bacteriophages and bacteria drive bacterial evolution? This is the question we're attempting to answer with the Community Science Project. You can help us gather data to address this question.

We will start with a phage called Fen4701_41 with Accession Number: **NC_027641**.

One of the predicted proteins in the database for Fen4704_41 is the **Major capsid protein VP1**. This phage has a small number of genes that are predicted to encode proteins, but our previous work has shown that only VP1 has homologs in bacteria, so that's the only one we'll investigate here.

```
major capsid protein VP1"
/protein_id="YP_009160373.1"
/db_xref="GI:906476351"
/db_xref="GeneID:25102737"
/translation="MGLFETNQAPRVPRAKFDLSHDRKLTMMNGTLIPVLNKEVLPGD
VWRCSP TVFLRFLALLAPIMHKVNVKLEVFVFCAYRLVWPNFPDFVSGGESGTAAPVFP
TFNTTQC NVTNTTILQDGSLLDYLGFPTYPQSGTYNAA YNGNTFSILHPRVYTLIYDT
WYRNQNVETTVIGTTFGADG PQNAGDYPIFSLRSRQWERDYFTSCQPSTQRGTQVVI
PISSTITSPSGNTNIPIVRVSGGAANTVSHAIGSKAADGSLMDQTS LAGVYLDPNGL
IASNTSLTVNALRASITLQTYKERMQRGGGRYIEYLWNVFGVASDDARLQRP ELYTL
REPVTVSEVLTTSTGCTGQTTPGTLAGAAMSIGSGSGSVYAAKEHGCIMALMSVIPR
SAYQNGIGPEFWRQVNTDFYVPDFAHLGEQTVIQNVYYDMTNA AQAASNGGVVFGYQS
RFAEYKYANDQVSGQFKSTLGFHWHLGRQFNSTSVPALNL AFCQCVPRKDIFAVTTDVD
HIVVDIYWNISVLRPMPYFGTPSTLNAVS"
```

Part 1: Forward BLAST – looking for bacterial homologs of bacteriophage proteins

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](#)

NEWS

BLAST+ 2.4.0 released
A new version (2.4.0) of the BLAST+ executables is now available.
Thu, 02 Jun 2016 14:00:00 EST [More BLAST news...](#)

Web BLAST

blastx
translated nucleotide > protein

tblastn
protein > translated nucleotide

1. BLAST the Major capsid protein VP1 with Accession ID YP_009160373 in the protein database
2. Use <blastp>, setting the parameters to search the <nr database> but only against <bacteria>. This way your top hits will be bacteria, not other bacteriophages.

Use the default settings for blastp.

The screenshot shows the NCBI BLAST search interface. At the top, there are tabs for different BLAST programs: blastn, **blastp**, blastx, tblastn, and tblastx. A tooltip for the blastp tab says "Click to close this tab; Option-click to close all tabs except this one". The main heading is "BLASTP programs search protein databases using a protein query. m".

The "Enter accession number(s), gi(s), or FASTA sequence(s)" section has a text input field with a red box around it containing the text "Paste Accession number (YP_009160373) here". To the right of this field is a "Clear" button and a "Query subrange" section with "From" and "To" input fields.

Below the text input field is the "Or, upload file" section with a "Choose File" button and "no file selected" text. There is also a "Job Title" input field with a placeholder "Enter a descriptive title for your BLAST search".

The "Choose Search Set" section contains several options:

- Database:** A dropdown menu showing "Non-redundant protein sequences (nr)".
- Organism:** A dropdown menu showing "bacteria (taxid:2)".
- Exclude:** A checkbox labeled "Exclude" with a "+" button.
- Entrez Query:** A text input field with a placeholder "Enter an Entrez query to limit search".

The "Program Selection" section contains a list of algorithms:

- ☒ blastp (protein-protein BLAST)
- ☐ PSI-BLAST (Position-Specific Iterated BLAST)
- ☐ PHI-BLAST (Pattern Hit Initiated BLAST)
- ☐ DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

At the bottom of the "Program Selection" section is a link "Choose a BLAST algorithm".

3. The results of this BLAST will take you to the "Hits" page. Record the **top 10 hits** in an Excel spreadsheet (**provided on the website**) to enter into the Community Science Project database. Make sure to record the hits based on the following cut-off values:

Sequences producing significant alignments:

Select: [All](#) [None](#) Selected:0

Alignments Download GenPept Graphics Distance tree of results Multiple alignment							
	Description	Max score	Total score	Query cover	E value	Ident	Accession
<input type="checkbox"/>	hypothetical protein [Parabacteroides distasonis]	340	340	98%	1e-107	37%	WP_036631168.1
<input type="checkbox"/>	Capsid protein (F protein) [Parabacteroides merdae]	337	337	98%	4e-106	37%	CUP01217.1
<input type="checkbox"/>	capsid protein VP1 [Parabacteroides merdae CAG.48]	333	333	98%	6e-105	37%	CDD13573.1
<input type="checkbox"/>	hypothetical protein [Parabacteroides distasonis]	330	330	98%	1e-103	38%	WP_005867318.1
<input type="checkbox"/>	hypothetical protein [Elizabethkingia anophelis]	315	315	98%	8e-98	35%	WP_051631734.1
<input type="checkbox"/>	hypothetical protein [Elizabethkingia anophelis]	309	309	97%	1e-95	35%	WP_024568106.1
<input type="checkbox"/>	hypothetical protein [Parabacteroides distasonis]	245	245	72%	5e-73	37%	WP_036630554.1
<input type="checkbox"/>	unnamed protein product [uncultured bacterium]	248	248	99%	4e-72	31%	CDL66973.1
<input type="checkbox"/>	unnamed protein product [uncultured bacterium]	246	246	96%	5e-72	32%	CDL65960.1
<input type="checkbox"/>	Capsid protein (F protein) [Chlamydia trachomatis]	246	246	99%	2e-71	32%	CRH84958.1

We're looking for things that are relatively closely related, therefore use **> 70% query cover** and **< 1e-50 E-value** as the cutoffs for this. If there are fewer than 10 that match these criteria, just capture these in the spreadsheet. If there are cogent reasons to change these cutoffs depending on the system you're looking at, capture this information as well.

Part 2: Reverse BLAST – looking for bacteriophage homologs of bacterial proteins

Now take the top hit from the Forward BLAST result, which in this case is:
hypothetical protein [Parabacteroides distasonis]
Accession number: WP_036631168.

1. Perform another BLAST using <blastp> but this time we will be looking for phage proteins. There is no "Bacteriophage" setting, so we'll use the <nr> database and restrict our results to the <virus> sector of the database.

The screenshot shows the NCBI BLAST web interface. The 'Enter Query Sequence' section has a text input field containing 'WP_036631168' and a 'Query subrange' section with 'From' and 'To' fields. Below this is an 'Or, upload file' section with a 'Browse...' button and a 'Job Title' field. The 'Choose Search Set' section includes a 'Database' dropdown set to 'Non-redundant protein sequences (nr)', an 'Organism' dropdown set to 'Viruses (taxid:10239)', and an 'Exclude' section with checkboxes for 'Models (XM/XP)' and 'Uncultured/environmental sample sequences'. The 'Program Selection' section has a radio button selected for 'blastp (protein-protein BLAST)'. At the bottom, there is a 'BLAST' button and a checkbox for 'Show results in a new window'.

2. The results of the REVERSE BLAST will take you to the “Hits” page and again you will record the top 10 hits (**spreadsheet provided on the website**). Use the same cutoffs (normally, > 70% query cover and < 1e-50 E-value) as the you did for the forward BLAST.

Part 3: Learning about the protein universe in which the protein of interest resides

Now check to see if your Top Hit from the REVERSE Blast matches the initial phage protein you had searched for.

In our search, the Reverse BLAST top hit is: major capsid protein VP1 [Parabacteroides phage YZ-2015a] with Accession Number: **YP_009218533**

This result matches with the initial phage protein we searched for -

major capsid protein VP1 [Microviridae Fen4707_41] with Accession Number:
YP_009160373

This result conveys that the two proteins are each other's closest homologs. This isn't always the case and suggests that there are other homologs that could be explored. This information can be recorded on the Reverse Blast spreadsheet provided on the website.

Part 4: Submitting your work to Genome Solver

If you want to take part in the Community Science Project, the information in your two spreadsheets should be submitted to the Community Science Project pages on the Genome Solver Qubeshub website (<https://qubeshub.org/community/groups/genomesolver>).

