

Agriculture & Agri-Food Canada Phage Genomics Workshop

Andrew M. Kropinski
Department of Pathobiology
University of Guelph, Canada
Email: Phage.Canada@gmail.com



Acknowledgement

We thank Brian Anderson from DNASTAR Inc. for access to the latest version of Lasergene Suite software



Biography

- ☐ Been working on phages since mid-1960s
- ☐ Held academic and government research positions
- ☐ Currently advise students & faculty at the University of Guelph – Adjunct Professor
- ☐ Past Chair, Bacterial and Archaeal Viruses Subcommittee of ICTV
- ☐ Genome Advisor to NCBI
- ☐ Sequenced: >150 phages
- ☐ Default-setting bioanalyst - Online Analysis Tools (<http://molbiol-tools.ca>)

If you're not a programmer...You're not a Bioinformatician! John H.E. Nash (PHAC)

- ☐ Apologies: I am a Windows and Mac person not a Unix/Linux user



Overview

You can:

1. have a commercial company sequence your phage, and assemble its genome
2. submit this sequence to GenBank and receive an "Accession Number"
3. annotate its genome using an online pipeline which will find that most of the genes specify "hypothetical proteins" OR
4. Devote some time to (a) preparing your genome for annotation and (b) carefully scrutinizing the annotations i.e. follow what you learn in this workshop

Overview 2

Laboratory output:

- ☐ High numbers of bacterial and phage genomes
 - Bioinformatician
 - Computing resources
 - Python programs
- ☐ Low number of bacterial and phage genomes
 - Windows, Mac and internet resources as described in this workshop

Workshop Outline

- ☐ Phage genome assembly – emphasis on Illumina paired-end data
- ☐ Autoannotation coupled with manual proofreading
- ☐ Phage taxonomy

Objectives

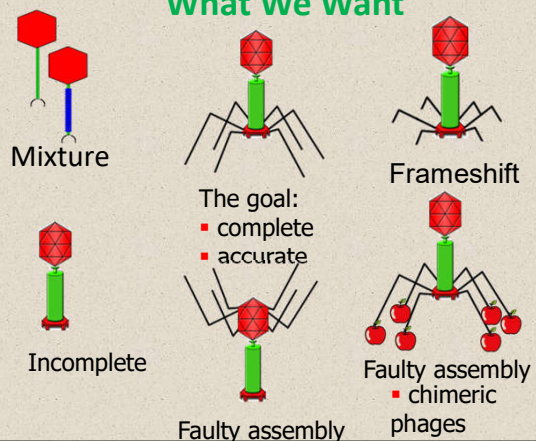
By the end of this workshop participants will

- ☐ have a deeper understanding of the steps involved in sequencing, assembling, and annotating phage genomes
- ☐ understand how phages are classified
- ☐ have an authoritative list of Internet resources and recommended software (commercial and free) for genome analysis.

PART 1

Genome Assembly

What We Want



Outline

Paired-end Illumina sequence data



Assemble
DNASTAR SeqMan NGen18

Primary assembly - contigs



Trim
DNASTAR SeqMan Ultra
Reassemble

Reassembly



Manipulation – RC, Cut/Paste, Splinting
SeqBuilder Pro

Genome for Annotation

Why emphasis on DNASTAR?

- ☐ Developed by geneticist Fred Blattner and computing science student John Schroeder (1984)
- ☐ I have been using their software since 1995
- ☐ Company responds readily to enquiries
- ☐ Software packages are **really** updated annually
- ☐ Available for Mac and PCs
- ☐ Intuitive software
- ☐ Excellent tutorials/videos

Non-commercial Alternatives

- ☐ SPAdes Genome Assembler
 - works with Ion Torrent, PacBio, Oxford Nanopore, and Illumina paired-end
 - Linux, macOS
 - URL: <https://github.com/ablab/spades>
 - Reference: Prjibelski A, Antipov D, Meleshko D, Lapidus A, Korobeynikov A. Using SPAdes De Novo Assembler. Curr Protoc Bioinformatics. 2020 Jun;70(1): e102. doi: 10.1002/cpbi.102. PMID: 32559359.
- ☐ **N.B.** I do not recommend Nanopore for phage sequencing

Setup for sequence assembly & analysis

- ☐ Create a directory with the name of the phage under analysis
- ☐ Create four subdirectories:
 - Data
 - Assembly
 - Reassembly
 - Annotation

Genome assembly using SeqMan NGen18

- ☐ Step-by-Step



SeqMan NGen

Welcome! What do you want to do?

New Assembly (indicated by a red arrow)

Manage Cloud Assemblies

Variant Analysis / Resequencing	ABI / Sanger
RNA-Seq / Transcriptomics	De novo assembly
De Novo Assembly and Finishing (indicated by a red arrow)	Genome finishing - refinement
Metagenomics	NGS-Based
Variant Call Format (VCF) Files	De novo assembly (indicated by a red arrow)
Combine / Reanalyze Existing Assemblies	Genome finishing - initial error correction
	Genome finishing - refinement
	Combined reference-guided/de novo assembly
	PacBio / Nanopore
	De novo assembly (beta)
	De novo assembly and polishing (beta)
	NGS polishing of draft genome (beta)

Move from screen to screen with NEXT

Adding Illumina data

Input sequences

Read technology: ☒ Paired-end data

Experiment setup:

Sequence File: Pair Distance:

M-Grey_S6_L001_R2_001.fastq.gz

Input sequences

Read technology: ☒ Paired-end data

Experiment setup:

Sequence File	Pair Distance
M-Grey_S6_L001_R1_001.fastq.gz	500
M-Grey_S6_L001_R2_001.fastq.gz	500

N.B. Compressed fastq data

Changing number of reads

Preassembly options

Read Filtering

☒ Maximum total reads:

☐ Remove reads smaller than:

Preassembly options

Read Filtering

☒ Maximum total reads:

☐ Remove reads smaller than:

Assembly options and names

Assembly options

Coverage Calculation

☐ Unknown genome length

☒ Estimated genome length bp

Assembly options

Coverage Calculation

☐ Unknown genome length

☒ Estimated genome length bp

Project name:

Project folder:

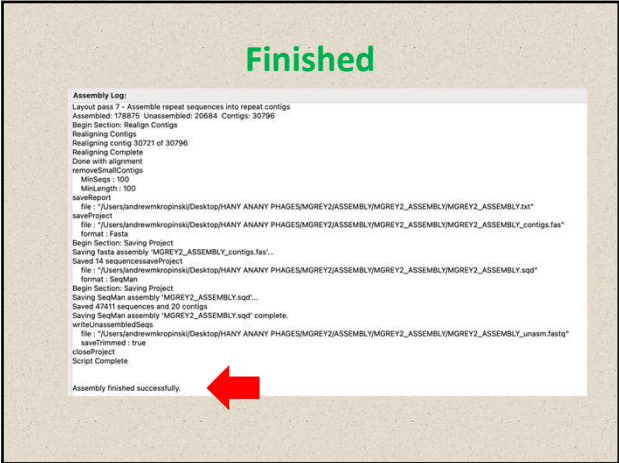
Assembly output:

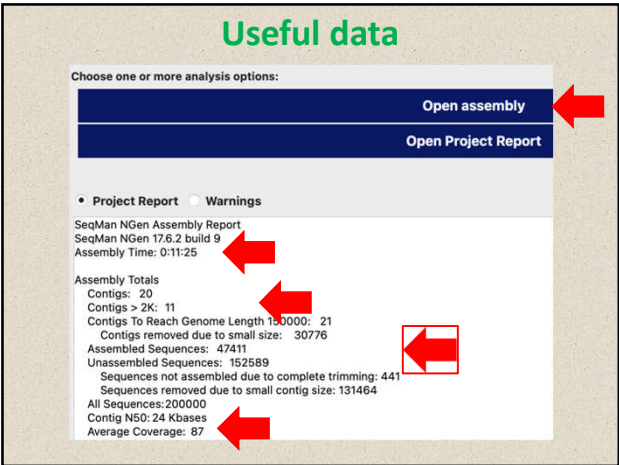
MGREY2_ASSEMBLY.sqd

MGREY2_ASSEMBLY_contigs.fas

MGREY2_ASSEMBLY.script





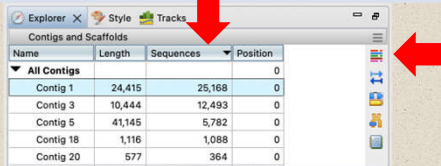


Examining & editing data in SeqMan Ultra

MGREY2_ASSEMBLY.sqd

/Users/andrewmkropinski/Desktop/HANY ANANY PHAGES/MGREY2/ASSEMBLY/MGREY2_ASSEMBLY

Contig N50: 24.0 kb
Largest contig size: 41,145 bp
Number of contigs: 20
Reads assembled: 47,411

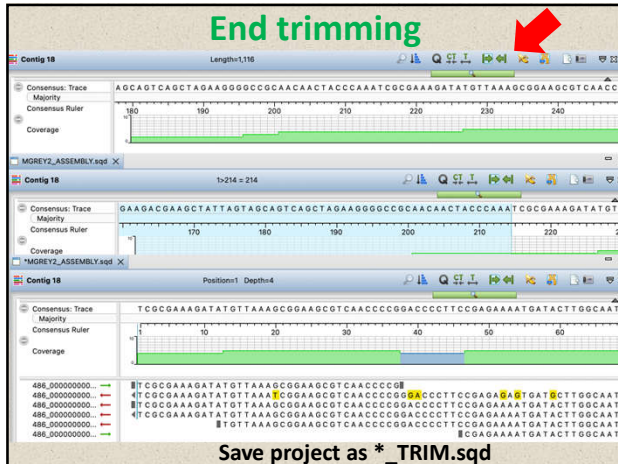


Name	Length	Sequences	Position
All Contigs			0
Contig 1	24,415	25,168	0
Contig 3	10,444	12,493	0
Contig 5	41,145	5,782	0
Contig 18	1,116	1,088	0
Contig 20	577	364	0

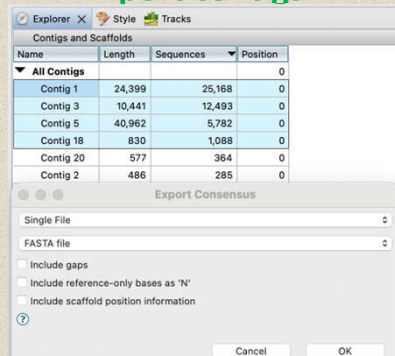
Click here to order by number of sequences

Click here to see a specific contig

End trimming



Export contigs



Explorer X Style Tracks

Contigs and Scaffolds

Name	Length	Sequences	Position
All Contigs			0
Contig 1	24,399	25,168	0
Contig 3	10,441	12,493	0
Contig 5	40,962	5,782	0
Contig 18	830	1,088	0
Contig 20	577	364	0
Contig 2	486	285	0

Export Consensus

Single File ☐

FASTA file ☐

☐ Include gaps

☐ Include reference-only bases as 'N'

☐ Include scaffold position information

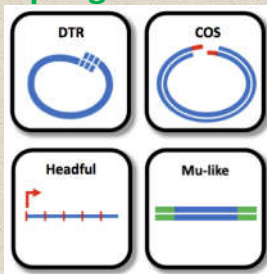
Cancel OK

☐ Select and export consensus of desired contigs as a single file in fasta format (*.fas)

9

Genome orientation - termini

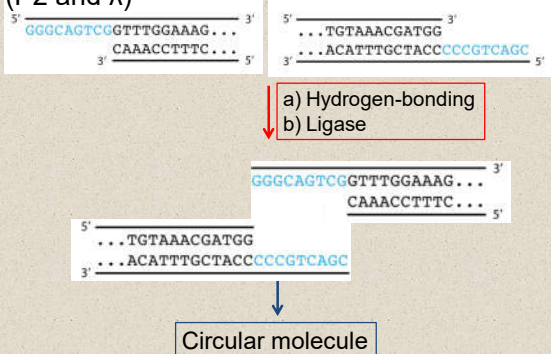
Bacteriophage Genome Termini



- ☐ Phage termini are problematic
- ☐ Bottom line – if your phage is related to an existing virus its genome should be colinear
- ☐ If a choice exists use RefSeq/ICTV homologs
- ☐ **Next six slides are largely for information only**

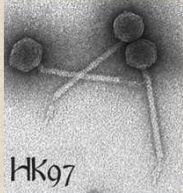
Termini of *Caudoviricetes*

A. Cohesive (sticky) ends – 5'-extensions (P2 and λ)



Termini of Caudoviricetes 2

B. Cohesive (sticky) ends – 3'-extensions (HK97, D3, many mycobacteriophages)

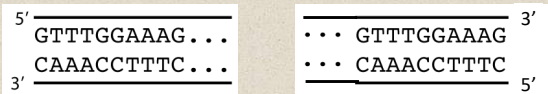


Site-specific packaging

(<http://www.pitt.edu/~duda/HK97.html>)

Termini of Caudoviricetes 3

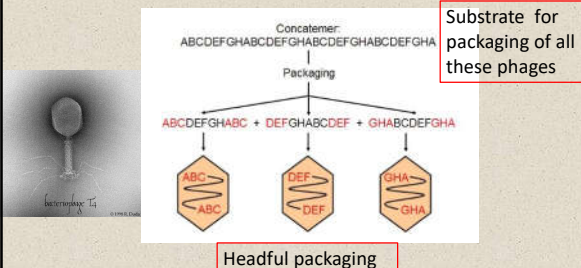
C. Terminal redundancy (TR; direct repeats (DR)) – *Autographivirinae* (T7, SP6, ϕ KMV), T5 phage, *Listeria* phage A511



- *Escherichia* phage T7 – 160 bp
- *Listeria* phage A511 – 3,125 bp
- *Escherichia* phage T5 – 10,219 bp
- *Bacillus* phage SPO1 – 13,185 bp

Termini of Caudoviricetes 4

D. Terminal redundancy (TR) with Circular Permutation (CP) – *Escherichia* phages T4 and P1



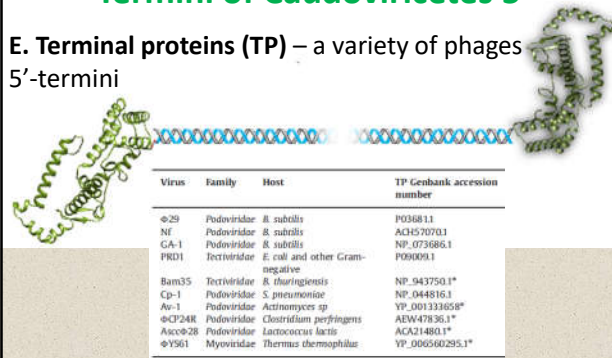
Substrate for
packaging of all
these phages

Headful packaging

N.B. No member of the class Caudoviricetes possess
a circular genome

Termini of Caudoviricetes 5

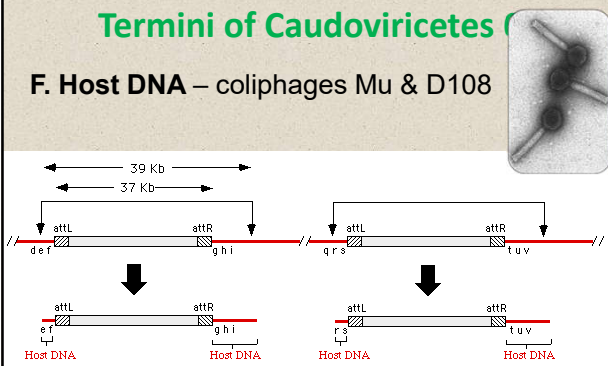
E. Terminal proteins (TP) – a variety of phages 5'-termini



(<http://www.sciencedirect.com/science/article/pii/S0042682214003742>)

Termini of Caudoviricetes 6

F. Host DNA – coliphages Mu & D108



(<http://www.sci.sdsu.edu/~smaloy/MicrobialGenetics/topics/transposons/Mu.html>)
http://2012.igem.org/Team:Marburg_SYNMIKRO/Project

Some general rules on genome termini

- ☐ T4-like phages begin with rIIA gene on complementary strand
- ☐ Many other phages begin with TerS/L

Termini of Caudoviricetes

- ☐ PhageTerm - Garneau JR et al. 2017. Sci Rep. **7(1)**:8292. doi: 10.1038/s41598-017-07910-5. PMID: 28811656.
- ☐ accessible via Galaxy Pasteur (<https://galaxy.pasteur.fr/>)
- ☐ video "How to run PhageTerm tool in Galaxy" <https://www.youtube.com/watch?v=9y2gfUSLkkg>
- ☐ for terminal repeats you can use the magnification bar in DNASTAR

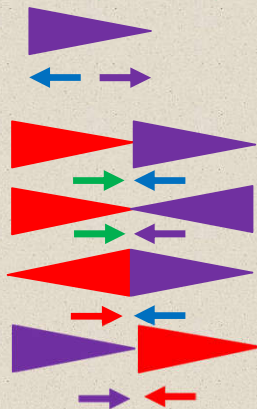
BREAK

What to do about gaps

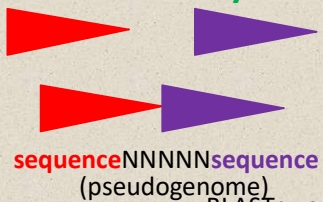
- ☐ Occasionally assembly result in two or more contigs which will not collapse into one.
- ☐ Gap closure techniques:
 - Primer walking and Sanger sequencing
 - requires specialize sequencing abilities
 - PCR and Sanger sequencing
 - Splinting and reference-guided assembly

PCR and Sanger sequencing

- ☐ Two fragments
- ☐ Possibilities:
- ☐ Sequence amplicon



Splinting and reference-guided assembly



BLASTn versus
Caudoviricetes (taxid:
2731619); exclude -
Uncultured/envirom

Teseptimavirus T7 genome assembly, chromosome: 1
Sequence ID: [OZ035731.1](#) Length: 39778 Number of Matches: 2

BLASTn results & splint

```

Query 421  AGACACCTAAAGTTGGTGGTATCTTTAAGAAGCCTAAGAACAAGGCACAGCGAGAAGGCC 480
Sbjct 18476 AGACACCTAAAGTTGGTGGTATCTTTAAGAAGCCTAAGAACAAGGCACAGCGAGAAGGCC 18535

Query 481  GTGAGCCTTGCGAACTTGAT 500
Sbjct 18536 GTGAGCCTTGCGAACTTGAT 18555

Query 511  GGTGGGTCCCACCAAGTACACCGATAAGGGTGCTCCTGTGGTGACGATGAGGTACTCG 570
Sbjct 18656 GGTGGGTCCCACCAAGTACACCGATAAGGGTGCTCCTGTGGTGACGATGAGGTACTCG 18715
  
```

- ☐ gap in homologous genome of 99 bp (18556-18655)
- ☐ splint – homologous regions (200 bp) plus missing sequence
- ☐ 18356 - 18855

Splint

Teseptimavirus T7 genome assembly, chromosome: 1
GenBank: OZ035731.1
FASTA Graphics

Teseptimavirus T7 genome assembly, chromosome: 1
GenBank: OZ035731.1
FASTA Graphics

Change region shown
☐ Whole sequence
☒ Selected region
from: begin to: end
Update View

Change region shown
☐ Whole sequence
☒ Selected region
from: 18356 to: 18855
Update View

Teseptimavirus T7 genome assembly, chromosome: 1
GenBank: OZ035731.1
FASTA Graphics

Go to

LOCUS OZ035731 500 bp DNA linear PH0 25-APR-2024
DEFINITION Teseptimavirus T7 genome assembly, chromosome: 1.
ACCESSION OZ035731 REGION: 18356..18855

What next

- ☐ use splint to generate a single contig
- ☐ use this contig as the template in a reference-guided assembly with your Illumina sequence data.
- ☐ you will be able to export the consensus for further manipulations and analyses.

Nonaligned contigs

- ☐ what are they and should I worry?
 - Unaligned fragments of your phage
 - Host DNA
 - Prophage DNA
 - Second phage

Is my phage related to anything?

Standard Nucleotide BLAST

blastn blastp blastx tblastn tblastx

BLASTN programs search nucleotide databases using a nucleotide query. more...

Enter Query Sequence

Enter accession number(s), gi(s), or FASTA sequence(s) [Clear](#) Query subrange [From](#) [To](#)

Or, upload file [Browse...](#) [Jeff fasta](#) [?](#)

Choose Search Set

Database ☒ Standard databases (nr etc.) ☐ rRNA/ITS databases ☐ Genomic + transcript databases ☐ Betacoronavirus ☐ Experimental databases

Core nucleotide database (core_nt) [?](#)

Organism [Optional](#) Caudoviricetes (taxid:2731619) ☐ exclude [Add organism](#)

Program Selection

Optimize for ☐ Highly similar sequences (megablast) ☐ More dissimilar sequences (discontiguous megablast) ☒ Somewhat similar sequences (blastn) [Choose a BLAST algorithm](#)

[BLAST](#) Search database core_nt using [blastn](#) (optimize for somewhat similar sequences) ☒ Show results in a new window

Is my phage related to anything 2?

Pseudomonas phage vB_Pae_LESphi2, complete genome

GenBank: OQ594955.1

[FASTA](#) [Graphics](#)

[Go to:](#) [☺](#)

LOCUS OQ594955 42123 bp DNA linear PHG 08-MAR-2024

DEFINITION Pseudomonas phage vB_Pae_LESphi2, complete genome.

ACCESSION OQ594955

VERSION OQ594955.1

KEYWORDS .

SOURCE Pseudomonas phage vB_Pae_LESphi2

ORGANISM Pseudomonas phage vB_Pae_LESphi2

Viruses; Duplodnaviria; Heunggongvirae; Uroviricota; Caudoviricetes.

Realigning your genome 2

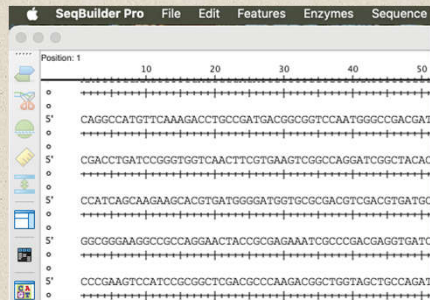
- ☐ May require RC
- ☐ May require cutting and reassembly using SeqMan Ultra
- ☐ Requires SeqBuilder Pro



ORIGIN 1 tcttctgcct tatgtttgat ttctatgtg tatcagtagc ttaagactga tcgcttcac

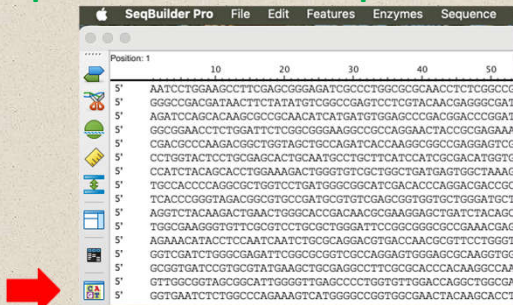
- ☐ Identify 5'-end of homologous phage
- ☐ Cutting & reassembly will result in your phage genome being colinear with this phage

SeqBuilder Pro – Default format



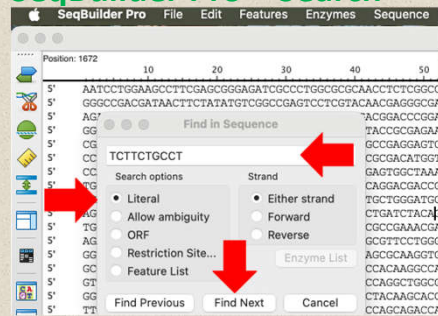
☐ Shows restriction sites automatically

SeqBuilder Pro – EditSeq format



☐ My preferred format

SeqBuilder Pro – Search



- ☐ Identify 5'-end and cut yields two fragments which I label **LEFT END** and **RIGHT END**
- ☐ Save in fasta format
- ☐ Reassemble using SeqMan Ultra

Is sequence ready for annotation?

❑ Questions:

1. Is it full length?
2. Is it error free?

❑ Quick and dirty approaches:

1. BLASTn
2. BLASTx

BLASTn – complete example?

LOCUS PP039555 39435 bp DNA linear PHG 18-SEP-2024
 DEFINITION Enterobacter phage vB_Ecl_MII_004, complete genome.
 ACCESSION PP039555
 VERSION PP039555.1
 KEYWORDS .
 SOURCE Enterobacter phage vB_Ecl_MII_004
 ORGANISM Enterobacter phage vB_Ecl_MII_004
 Viruses.
 REFERENCE 1 (bases 1 to 39435)

- ❑ BLASTn reveals it is related to *Escherichia* phage Peacock (MK903279; 39233 bp)
- ❑ Caudoviricetes; Autographiviridae; Studiervirinae; Kayfunavirus

BLASTn – incomplete?

LOCUS PP935704 15914 bp DNA linear PHG 03-AUG-2024
 DEFINITION Salmonella phage vB_SE126_2P, complete genome.
 ACCESSION PP935704
 VERSION PP935704.1
 KEYWORDS .
 SOURCE Salmonella phage vB_SE126_2P
 ORGANISM Salmonella phage vB_SE126_2P
 Viruses.
 REFERENCE 1 (bases 1 to 15914)

- ❑ BLASTn reveals it is incomplete

	Description	Scientific Name	Max Score	Total Score	Query Cover	E value	Per. Ident	Acc. Len	Accession
✓	Salmonella phage f18SE, complete genome	Salmonella phage f18SE	21603	22575	100%	0.0	91.90%	41968	NC_028908.1
✓	Salmonella phage f2SE, complete genome	Salmonella phage f2SE	21599	22572	100%	0.0	91.89%	41965	KJ051149.1
✓	Salmonella phage f3SE, complete genome	Salmonella phage f3SE	21594	22568	100%	0.0	91.89%	41967	KJ051147.1
✓	Salmonella phage fSE4S, complete genome	Salmonella phage fSE4S	21581	22554	100%	0.0	91.87%	41768	KT881477.1
✓	Salmonella phage fSE1C, complete genome	Salmonella phage fSE1C	21547	22520	100%	0.0	91.81%	41720	KT962832.1
✓	Salmonella phage vB_Sen9_PVP-SE2, complete genome	Salmonella phage vB_Sen9_PVP-SE2	18206	22565	100%	0.0	91.92%	42425	NC_073198.1

BLASTx - frameshifts

LOCUS PQ039555 39435 bp DNA linear PHG 18-SEP-2024
 DEFINITION Enterobacter phage vB_Ecl_MII_004, complete genome.
 ACCESSION PQ039555
 VERSION PQ039555.1

- ☐ BLASTn reveals it is related to *Escherichia* phage Peacock (MK903279; 39233 bp) [& complete]
- ☐ BLASTx versus *Escherichia* phage Peacock (taxid:2591100) – **huge number of frameshifts**

Alignment Scores ■ < 40 ■ 40 - 50 ■ 50 - 80 ■ 80 - 200 ■ >= 200 ?

Distribution of the top 119 Blast Hits on 47 subject sequences



What are the problems

- ☐ scientists with little knowledge and experience working with phages
- ☐ individuals not taking advantage of free expertise in the International Committee on Taxonomy of Viruses (ICTV)



- ☐ Nanopore versus Illumina sequencing technology

COMMENT ##Assembly-Data-START##
 Assembly Method :: Canu v. 1.7.1; Flye v. 2.9; Racon v. 1.4.13
 Coverage :: 24.22x
 Sequencing Technology :: Oxford Nanopore Technology
 ##Assembly-Data-END##
 FEATURES Location/Qualifiers

End of Part 1



Questions?