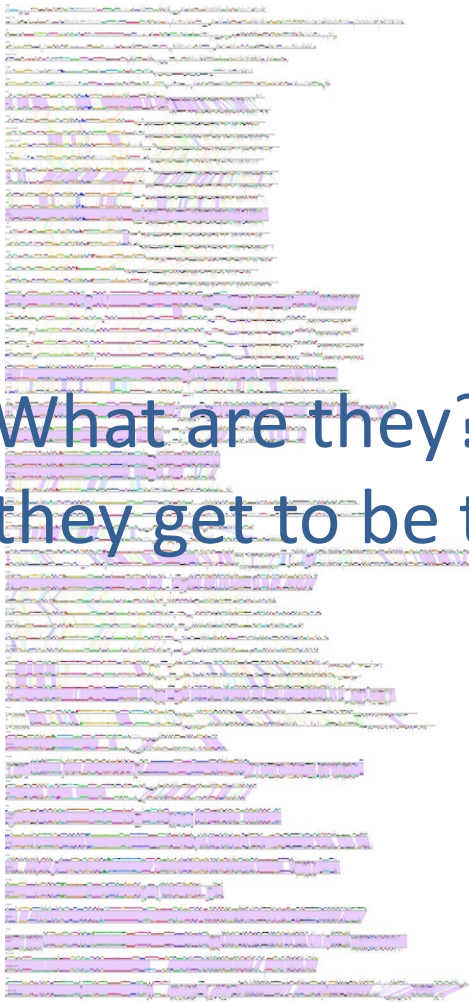


Predicting Genes in Actinobacteriophages

2022 Genomics Workshop Training

SEA-PHAGES Cohort 15

Deborah Jacobs-Sera



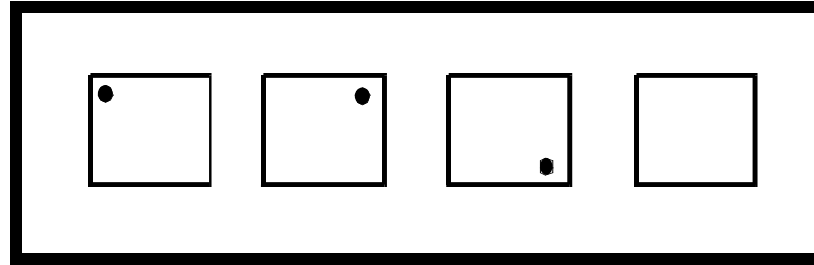
What are they?
How did they get to be that way?

It is all about finding the patterns...

Since the beginning of time, woman (being human) has tried to make order and sense out of her surroundings. Gene annotation and analysis is just a primal instinct to make order.

Young children, as they prepare to enter school, are tested to see if they are ready by recognizing patterns, a form of making order.

1. Where will the dot appear in the 4th box?



Remember, everything you need to know, you learned in kindergarten....

Make-Believe or Putative



Remember, you are working in the putative gene world. All gene **predictions** are made with the best evidence to date. Most of that evidence is computational (bioinformatic), not experimental. Tomorrow's data may give us better evidence, but your prediction today is the best it can be ... today! Make good predictions following a consistent approach. Let these predictions lead to experimentation that can provide the evidence to improve future predictions.

How many phage genome sequences are in GenBank?

~~13390~~

No longer countable

How many actinobacteriophage genomes are sequenced?

4153

How many As, Cs, Ts, and Gs are in a mycobacteriophage genome?

On average: ~70,000 base-pairs

Range: ~40,000 to ~165,000 bp

What is the universal format for a sequence?

FASTA

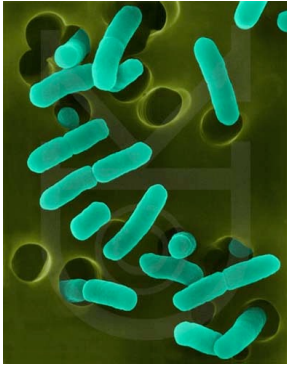
How do you make sense of the nucleotide sequence?

Convert to genes

How do you convert ATCGs to genes?

Codons

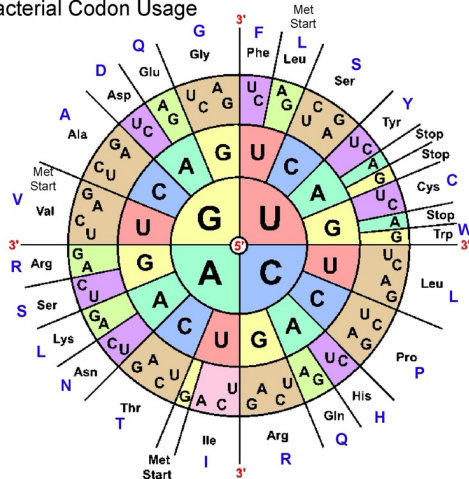
Code for Amino Acids, Starts, Stops



www.cen.ulaval.ca

- Phages use the Bacterial and Plant Plastid code (NCBI: Table 11)
- 3 starts
 - ATG (methionine)
 - GTG (valine)
 - TTG (leucine)

Bacterial Codon Usage



- 3 stops (TAA, TAG, TGA)
- Space in-between: Open Reading Frame -- ORF

ATGGACCTCTCGCCC

ATG GAC CTC TCG CCC

TGG ACC TCT CGC

GGA CCT CTC GCC

If there are 3 choices (frames) in the forward direction,
how many are in the reverse direction?

>Echild complete sequence, 53159 bp including
10bp overhang (CGGTCGGTTA), Cluster A2

TGCGGCCGCCCATCCTGTACGGGTTTCCAAGTCGATCGGAGTCCCGAGC
CGGCGCAGGAGCGCCTCACCCAGCCTCTGTGCGCCCCCAGGACGCAAGAT
CCCCGCTCACGCGGGTAGTTGTATGGGCTAATCGGCAAACGGCCTCTGAG
GCCGCGAGACCAATGTCACACCAGGTGGTGGATGTTATTGACGCACGCGT
CCGTTAAGAGGACATGGCCTAGGTATGGCTACCCAACTTAGATTCAAAA
CCAGTCCCCTGGCCCCGTCGTCGGTGTTGCCGCTCCTGGCGGGCGGGGG
CCAGGTCCACCACCGCAGGGAGGACCGCCATGAAGATCATCCGCTCGCTC
GCCGGGGCGCTCGCACTCATTGCGATCGCCGCCCGAGGGCGGCTGGGATGC
GGAAATCTACGAGCCGTGGGATGAGGACGAATACCTCCTATAGTGATCTA
CGCCACTTGCTCGGTGGGTGTCAAGTGATACTCATGTATCTAGTTATTGA
GGCCTAAAGGCCCGAATAAGAGCCGCACAGGCGGCTCTCTAAGAGCGCC
CACTAGGGCGCTCGAAGTAATACCGGCCTTGAGGGCCGGTTATCTGACCC
GGCAACCGCCGGGTCTTCTGCCGCGCCAGTGGCGCGGCTCATAGAAGGG
GTGAGGCAACCGTGTACGGCACTCGCTCGAGTGCCTACTGGGCCTCGCAG
CCGGGGAAGTTCGACGTTCTGAACCTGCGGATGACGTTCCCGAGCACGTC

Six-Frame Translations

```

C R F W S V R N P G V R G V S R P F R N G V * P A V V F A S T C * F P K W E T
A D F G L Y G T R G F R G F P E M G S D L R F S P A L L V D S R N M G T
Q I L V C T E P L G C G F A V S P K W G L T C G F R Q H L I P E M G T
1 TGCAGATTTGGTCTGTACGGAAACCGGGGGTTCCCGGTTCCCGGAAATGGGCTGCACCTGCGGTTTTCCGCCAGACTTGTGATCCCGGAAATGGG
.....
ACGCTAAAACAGACATGCCTTTGGGCCCCAAAGCCCAAGGGGTTTACCACAGACTGGACGCCAAAGGGCTGCTGAACAACTAAGGGCTTTACCC
.....
O L N O D T R F G Q L T E R N G R F P T Q G A T K A L V Q A N G F H
A S K P R Y F V P R P N R P K G S I P D B R R N E G A S T S E R F P H
S C I K T Q V S G P P K A T E G F H P R R V Q P K R W C K N I P
E V I M P P V P K D P S V R A R R R N K S A T R A T L S A D H D V V
R K S S C H L Y L K I L L C V L V A I S L R R R G L R C L R I M M W V
G S H H A T C T * R S F C A C S S Q * V C D A G Y V V C G S * C G
101 AGGAAGTCATCCACCTGTACCTAAAGATCCCTCTGTGCGTCTCGTGCCAAATAGTCTGCCAGCGGGGTACGGTGTCTCGGGATCATGATGGGT
.....
TCCTTCACTAGTACGGTGGACATGGATTTCTAGGAAGACACCGCAGGACGGCTTATCAGACGCTGGCGGATGCAACAGACGGCTACTACTACACCA
.....
S S T M M G G T G L S G E T R A R R L L D A V R A V N D A S * S T A
L F D D H W R Y R F I R R H T S T A I L L R R R P S R O R R I H
P L * * A V Q V * L D K Q A H E D C Y T Q S A P * * T T Q P D H I P
A P D L P D G V A W H P L T V R W W N D I W A S P M A P E Y T D S
L Q T C R M V L R G I R * R C A G G M T F G R R R W P P R S T Q T R
G S R L A G W C C V A S V D G A L V E * H L G V A D G P G V H L C
201 GGCTCCAGACTTGGCGGATGGTGGGGGGCATCCGTGACGGTGGCTGGTGGATGACATTTGGGCGTGGCGGATGGCGCGGAGTACACAGACTCC
.....
CCGAGGCTGAACGGCTACAGCAACGACCGTAGGCACTGCCACGGACCACTTACTGTAACCCGGAGGGTACCAGGGCCCTATGGGTCTGAGG
.....
A G S K G S P T A H C G N V T R Q H F S M Q A D G I A G S Y V S E
H S W V Q R I T N R P M R Q R H A P P I V N P R R R H G R L N P Q C V R
P E L S A P H H Q T A D T S P A S T S H C K P T A S P G P T C L S
D I N G L F R V A M L Y N D F W T A D N A K A R A E A Q V R L E K A
I S T G C S V W R C S I T I F G L P I T R R R V R R L R F G W R R
Y Q R V V P C G D A L * R F L D C R * R E G A C G G S C S V G E G
301 GATATCAACGGTGTTCCTGTGGCGATGCTTATAAGATTTTGGACTGCCATACCCAGGCGCGTGGGGAGGTCAAGTTCGTTGGGAAGC
.....
CTATAGTGGCCACAGGCAACCGGTACAGATATGCTAAAACCTGACGGCTATTCGGTTCGCGGCAACGCTCCAGTCCAACTNCCTTTCC
.....
S I L P N N R T A I S * L S K Q V A S L A F A R A S A * T R N S F
I D V P Q E T H R H E I V I K P S G I V R L R T R L S L N P Q L L
P Y * R T T G H P S A R Y R N K S Q R Y R S P A H P P E P T P S P
D T D Y G T N P L A R R R L E W O I E A T E D S K A K G S K R R K
P T L I M G R I R W L V A V W S G R L R R P R I R R L R G R S G G
R H * L W D E S V G S P F G V A D * G V E A E
401 CGACACTGATTTGGGACGATCCGTGGCGTGGCGCGTTGGAGTGGGAGATGAGGGGACCGAGGATTCGAAAGCTAAGGGCTCGAAGCGGGGAA
.....
GGCTGGACATACCCTCGTACGCTGCCGACGCGGACCTGATCCGCTATCCGCGCTCCCTACGCTTCCCGCCTTCCGCGCTTCCGCGCTT
.....
A S V S * P V F G N A R R K S H C T S A V S S E F A L P D F R R
G V S I I P R I R Q S T A T Q L P L N L R G L I R L S L P R L P P
R C Q N H S S D T P E D G N P T A S Q P S R P N S P * P T S A S
S E A A P V C P P E P G D D P R L K L V T * R P Y G C F A G P C C
L R L R R C A H R S L V T I L V * S L * P D G L M A V I Q V P A V
V * G C A G V P T G A W * R S S F E A C D L T A L L W L F C R S L L W
501 GTCTGAGGCTGGCGGTTGGCCACGGGCGTGGTACGATCCCTGTTGAAGCTTGTGACCTGACGGCGTTTGTGGGCTGTTCCTGTTG
.....
CAGACTCCGACGGCCACCGGGTGGCTCGGACCACTGCTAGGAGCAAACCTCGAACACTGGACTGCGGGAATACCGACAAAACGTCAGGGAGCAC
.....
D S A A G T H G G S G P S S G R K F S T V Q R G * P Q K A P G Q Q
L R L S R R H A W R L R T V I R T Q L K H G S P R I A T K C T G A T
T Q P Q A P T G V P A Q H R D E N S A Q S R V A K H S N Q L D R S
G F S V S Y V G S A G V * L H * G S D G V R P W L I V G A G R T S R
D L A F P T L G P Q V C D F I E D R M V F G P G S L S G Q A R L
I * R F L R W V R R C V T S L R I G W C S A L A H C R G R P H V S
601 GGATTCAGCTTTCCACTGTTGGTCCGGAGGTGTGTACTTATGAGGATCGGATGGTGTTCGGCGCTGGCTCATTTGGGGGCGGGCGACCTGCTC
.....
CCTAATCGCAAAGGATGCAACCCAGGGCTCCACACACTGAAGTAACTCCTAGCTTACCAGCGGGAGCCGATTAACAGCCCGCTCCGGCTGCAGAG

```

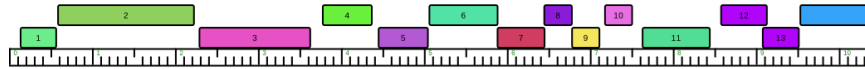
Ovechkin_Draft



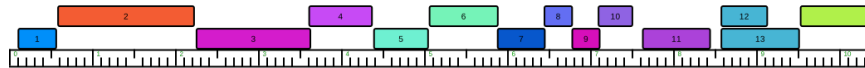
Ovechkin



Ovechkin_Draft



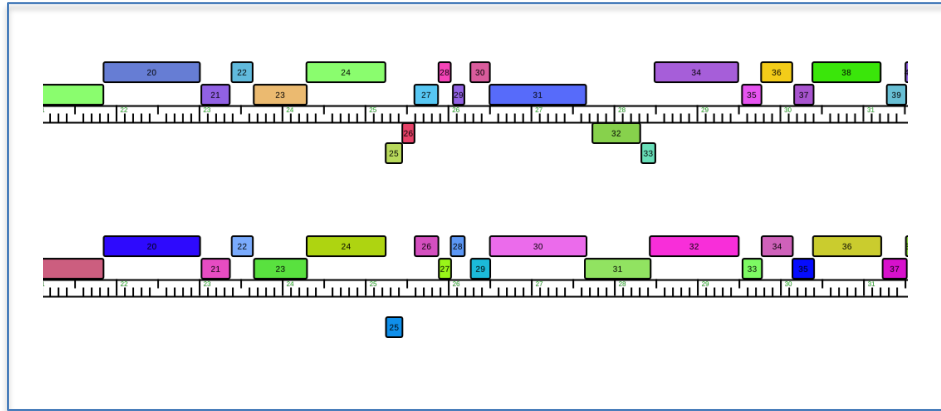
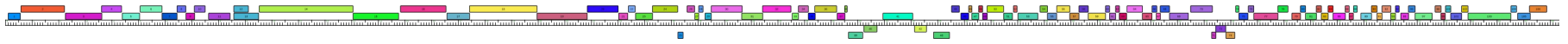
Ovechkin



Ovechkin_Draft



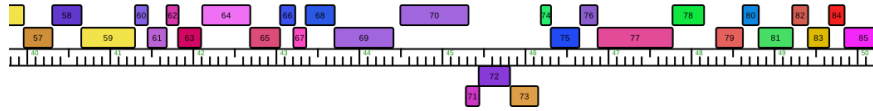
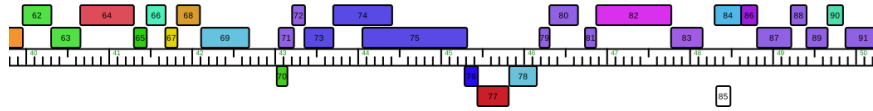
Ovechkin

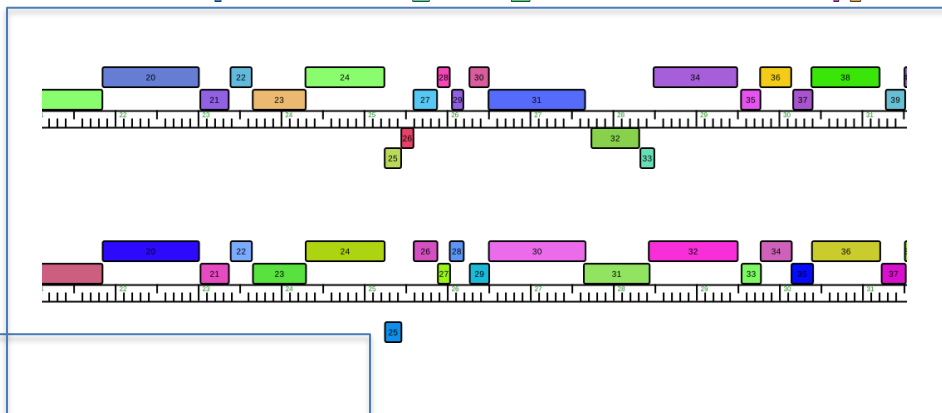


Ovechkin_Draft

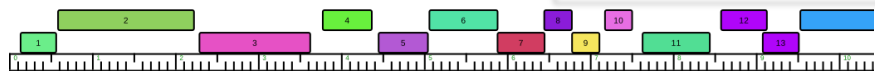


Ovechkin

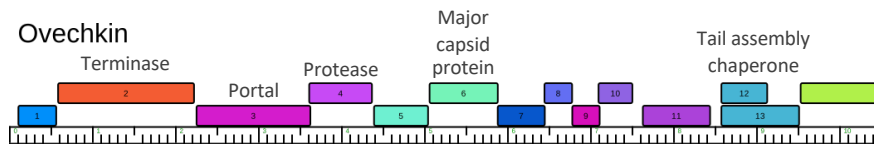




Ovechkin_Draft



Ovechkin



Gene Evaluations

For each feature we have 3 questions.

- Is it a gene?
- What is its start?
- What is its function?

For each question, what is the supporting data?

Gene Evaluations

- We are always looking for the **supporting data**.
- DNA master is our genome editor
- We (it) use 2 programs, Glimmer and GeneMark, to identify coding potential.
- We use Phamerator output for a visual representation of gene and nucleotide similarity.
- We use the Guiding Principles to remind us all of the parameters.
- As we evaluate, we can:
 - Add a gene
 - Delete a gene
 - Change a gene start

Supporting Data #1: Coding Potential

Glimmer and GeneMark

- Use Hidden Markov Models to identify coding potential
- Use **a sample** of the genome
- Identify longest ORFS in that sample
- Calculate patterns in the nucleotides:
 - 2 at a time, 4 at a time

```

GLIMMER (ver. 3.02; iterated) predictions:
orfID      start      end frame  score
-----
>Sheen complete sequence, 52927 bp including 10 bp 3' overhang (CGGGCGGTAA), Cluster A7
orf00001   732       1166    +3    11.17
orf00002   1259      1576    +2    14.04
orf00004   1566      2318    +3    11.01
orf00006   2347      3570    +1    10.85
orf00007   3599      3877    +2     1.93
orf00008   3889      4512    +1    10.39
orf00009   4509      5477    +3     5.52
orf00011   5731      7155    +1    12.91
orf00012   5772      5635    -1     2.63
orf00013   7152      7595    +3    10.63
orf00014   7592      8332    +2     5.98
orf00016   8359     10059    +1    11.16
orf00018   10056     11552    +3    15.98
orf00020   11549     12562    +2    11.29
orf00021   12621     13130    +3    13.44
orf00022   13160     14149    +2    18.56
orf00023   14229     14390    +3     9.03
orf00025   14394     14768    +3    11.22
orf00026   14765     14920    +2     2.58
orf00028   14917     15300    +1    11.86
orf00029   15303     15647    +3    10.29
orf00030   15660     16109    +3     7.67
orf00032   16124     16708    +2    15.52
orf00033   16821     17186    +3    12.50
orf00035   17354     17614    +2     5.52
orf00037   17618     20998    +2    11.78
orf00038   21003     22982    +3    15.20
orf00041   22979     24781    +2    16.51
orf00042   24798     25265    +3     6.64
orf00043   25298     25588    +2     6.78
orf00044   25593     27047    +3    13.74
orf00045   27051     27377    +3     7.90
orf00047   28925     27417    -3     7.82
orf00048   29214     29071    -1    14.71
orf00049   29802     29311    -1     3.17
orf00050   29936     29799    -3     9.82
orf00051   30417     30229    -1    13.15

```

GLIMMER

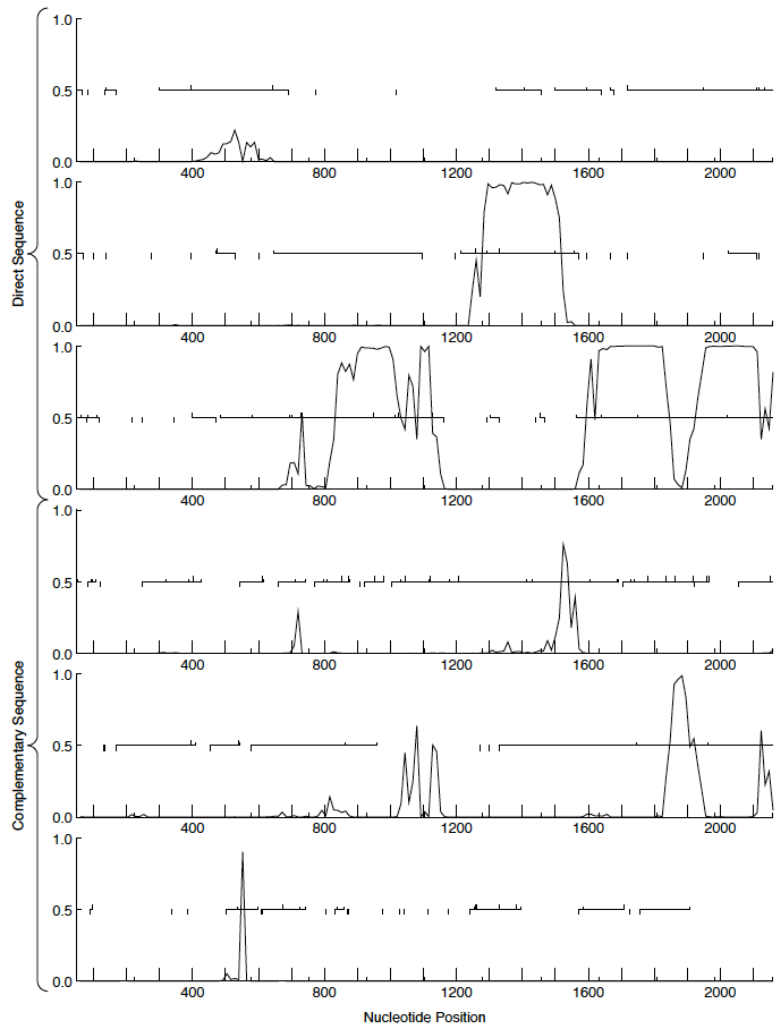


http://www.ncbi.nlm.nih.gov/genomes/MICROBES/glimmer_3.cgi

Microbial Genome Annotation Tools

GLIMMER is a system for finding genes in microbial DNA, especially the genomes of bacteria, archaea, and viruses. GLIMMER (Gene Locator and Interpolated Markov Model[ER]) uses interpolated Markov models to identify coding regions.

- Delcher AL, Harmon D, Kasif S, White O, Salzberg SL. Improved microbial gene identification with GLIMMER. *Nucleic Acids Research* 27:23 (1999), 4636-4641.
- Salzberg S, Delcher A, Kasif S, White O. Microbial gene identification using interpolated Markov models. *Nucleic Acids Research* 26:2 (1998), 544-548.



GeneMark Output (trained on *M. tuberculosis*)

DNA Master

ORF Analysis for Jefe_Draft

bp: 1140 | 173 345 517 689 861 1033 1205 1377 1549 1721 1893 2065 2237 2409 2581 2753 2925 3097 3269 3441 3613 3785 3957 4129 4301 4473 4645 4817 5000

ORF 827 - 1126

Overview Features References Sequence Documentation

Sort By Index Tag Name 5' End 3' End Length

Tag	Name	5' End	3' End	Length
DNAM_1	1	1	645	645
DNAM_2	2	642	830	189
DNAM_3	3	839	1126	288
DNAM_4	4	1123	2745	1623
DNAM_5	5	2820	3053	234
DNAM_6	6	3050	3457	408
DNAM_7	7	3457	3621	165
DNAM_8	8	3682	4422	741
DNAM_9	9	4422	4748	327
DNAM_10	10	4745	5122	378
DNAM_11	11	5132	5404	273
DNAM_12	12	5401	5811	411
DNAM_13	13	5811	6179	369
DNAM_14	14	6176	6406	231
DNAM_15	15	6403	6747	345
DNAM_16	16	6840	7385	546
DNAM_17	17	7451	7621	171
DNAM_18	18	7624	8016	393
DNAM_19	19	8013	8201	189
DNAM_20	20	8427	10094	1668
DNAM_21	21	10091	10637	447
DNAM_22	22	10518	10616	99

Description Sequence Product Regions Blast Context

Name 3 **GeneID**

Type CDS **GI**

5' End 839 Locus Tag DNAM_3

3' End 1126 Regions 1

Length 288

Direction Forward

Translation Table **Standard Code**

EC Number

Product gp3

Function

Notes
Original Glimmer call @bp 839 has strength 5.10. GeneMark calls start at 878

Choose ORF start

Start: 15 ORF Start : 839 Cdn1 Cdn2 Cdn3 Length SD Scoring Matrix Kibler6 Explore

Selected: 1 ORF Stop : 1126 5' End 100.0 0.0 100.0 3 Spacing Weight Matrix Karlin Medium Document

ORF Length : 288 3' End 63.8 57.7 75.8 447

Sta	Raw SD	Genomic	Spacer	Final	Sequence of the Region	Start	Start	ORF
#	Score	Z Value	Distance	Score	Upstream of the Start	Codon	Position	Length
1	-2.624	2.661	18	-4.925	CTGAGGATCGGCATGTGAT	GTG	677	450
2	-6.213	0.918	10	-6.908	AGGATCGGCATGTGATGTG	GTG	680	447
3	-3.349	2.309	10	-4.043	GCATGTGATGTGTCACCG	GTG	689	438
4	-3.225	2.369	12	-4.061	TGACGTGGCGAGGCTCTGAG	TTG	784	393
5	-4.299	1.848	9	-5.074	GGACGCCACGACAGGGGAC	GTG	758	369
6	-6.188	0.931	11	-6.945	CGTGGCTCCCGCCCGCGCA	GTG	794	333
7	-2.187	2.873	10	-2.882	AGGACAGTGGAGTGGGCGC	ATG	827	300
8	-5.308	1.358	11	-6.065	GGATGGGCCATGACCCGTTG	ATG	839	288
9	-6.359	0.848	11	-7.116	CGCATGACCCGTTGATGTCC	GTG	845	282
10	-6.193	0.928	14	-7.539	TGATCCCTGACCCCGTGCACA	TTG	878	249
11	-4.918	1.547	11	-5.675	GATCTACTACACTAGAGGAC	GTG	944	183
12	-3.942	2.021	10	-4.637	CTACAGTACAGGCGGCGCC	GTG	950	177
13	-2.071	2.930	16	-3.867	CCTGGAGGAGCCATGAGCGG	ATG	1034	93
14	-6.940	0.566	12	-7.775	CCTGGCGATGCTCTCTCTCG	ATG	1097	30
15	-2.624	2.661	8	-3.846	TCTCTCGGATGAGGACCGG	ATG	1109	18

91 Features Live Six-frame map of starts and stops 53125

Accession Favorite Genome Recent Genome Recent File ATG ME1 OD

Supporting Data #2:

- DNA Master – our genome editor
- Draft Annotation
 - Glimmer
 - GeneMark
- Refinement of Draft Annotation
 - Blast
 - NCBI
 - PhagesDB
 - HHPred Suite
 - Comparative Data
 - Phamerator
 - Starterator

Blast Comparisons

BLAST® Basic Local Alignment Search Tool

Home Recent Results Saved Strategies Help

NCBI BLAST/ blastp suite

blastn blastp **blastx** tblastn tblastx

Enter Query Sequence

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

From

To

Or, upload file No file selected.

Job Title

Enter a descriptive title for your BLAST search

Align two or more sequences

Choose Search Set

Database

Organism Exclude

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

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

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

Enter an Entrez query to limit search

Program Selection

Algorithm

blastp (protein-protein BLAST)

PSI-BLAST (Position-Specific Iterated BLAST)

PHI-BLAST (Pattern Hit Initiated BLAST)

DELTA-BLAST (Domain Enhanced Lookup Time Accelerated BLAST)

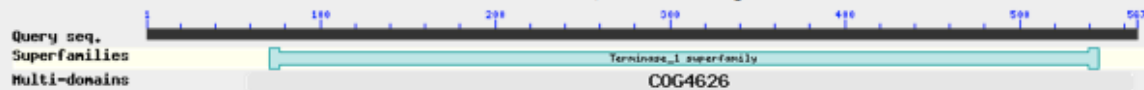
Choose a BLAST algorithm

BLAST Search database Non-redundant protein sequences (nr) using Blastp (protein-protein BLAST)

Show results in a new window

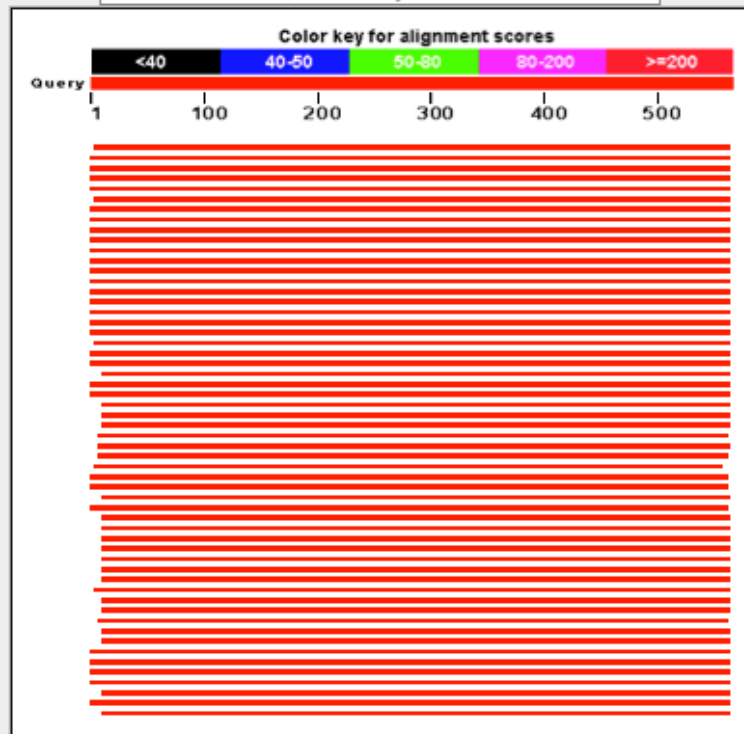
[Algorithm parameters](#)

Putative conserved domains have been detected, click on the image below for detailed results.










Distribution of 100 Blast Hits on the Query Sequence

Mouse over to see the define, click to show alignments



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"/>	gp12 [Mycobacterium phage Timshel]	1047	1047	99%	0.0	90%	AEJ92326.1
<input type="checkbox"/>	terminase [Mycobacterium phage Obama12]	1045	1045	99%	0.0	88%	YP_009007203.1
<input type="checkbox"/>	terminase [Mycobacterium phage Dhanush]	1045	1045	99%	0.0	88%	AGK87127.1
<input type="checkbox"/>	gp11 [Mycobacterium phage Flux]	1044	1044	99%	0.0	87%	AFL47832.1
<input type="checkbox"/>	terminase [Mycobacterium phage Nyxis]	1044	1044	99%	0.0	88%	YP_009005850.1
<input type="checkbox"/>	terminase [Mycobacterium phage HINdeR]	1043	1043	99%	0.0	89%	YP_008051863.1
<input type="checkbox"/>	gp10 [Mycobacterium phage Arturo]	1043	1043	99%	0.0	88%	AFU20464.1
<input type="checkbox"/>	terminase [Mycobacterium phage Kampy]	1043	1043	99%	0.0	87%	YP_009031871.1
<input type="checkbox"/>	gp11 [Mycobacterium phage Shaka]	1043	1043	99%	0.0	87%	AEF57321.1
<input type="checkbox"/>	gp11 [Mycobacterium phage Sabertooth]	1042	1042	99%	0.0	87%	AFU20551.1
<input type="checkbox"/>	terminase [Mycobacterium phage BellusTerra]	1040	1040	99%	0.0	87%	YP_009005569.1
<input type="checkbox"/>	gp11 [Mycobacterium phage ICleared]	1040	1040	99%	0.0	87%	AFL46618.1
<input type="checkbox"/>	gp11 [Mycobacterium phage Peaches]	1040	1040	99%	0.0	87%	YP_003358714.1
<input type="checkbox"/>	gp11 [Mycobacterium phage Wile]	1038	1038	99%	0.0	87%	YP_009014094.1
<input type="checkbox"/>	terminase [Mycobacterium phage QuinnKiro]	1023	1023	99%	0.0	86%	AIS73685.1

Download GenPept Graphics

gp12 [Mycobacterium phage Timshel]

Sequence ID: [gb|AEJ92326.1](#) Length: 564 Number of Matches: 1

Range 1: 4 to 564 GenPept Graphics Next Match Previous Match

Score	Expect	Method	Identities	Positives	Gaps
1047 bits(2707)	0.0	Compositional matrix adjust.	504/562(90%)	531/562(94%)	1/562(0%)
Query 5	YLSTPELLPQPPHKGIPVWLCHEDGSWALPKYTLGWGLNWLAEYVRSFAGGGPIPTLE	64			
Sbjct 4	YLNPGLLPQPPHKGIPVWQVHEDGSWALPARTLQWGLNWLAEYVRSFAGGGVFIPTLE	63			
Query 65	QARFILWYAVDENGVIYAYREGVLRMRKMGWKDPLCAAIALVELCGPVAFSHWDEKGNPV	124			
Sbjct 64	QARFILWYAVDE G YAYREG LRRMKGWKDPLCAAIALVELCGPVAFSHW G+PV	123			
Query 125	GKRRHAANITIAAVSQDQTKNTFSLFPPVMSKQMKTEYGLDQVNFVIYTEDGGRIEAATS	184			
Sbjct 124	GK RHAANIT+AAVSDQDQTKNTFS+FPVMISK+MK +YGLDQVNFVIY+E+GGRIEAATS	183			
Query 185	SPASMEGNRPRTLVIENETQWGWGPDGNVNDGVAADDVIEGNVSKIPGARLKAICNAHIP	244			
Sbjct 184	SPASMEGNRPRTLVIENETQWGWGPDGNVNDG AMDDVIEGNVSKIPGARLKAICNAHIP	243			
Query 245	GNDTVAEKAYDHWQDILSGKAVDTGLMYDALEAPADTPVSEIPSEKEDPEGYEAGIAQLM	304			
Sbjct 244	GNDTVAEKAYDHWQDILSGKAVDTG+MYDALEAPADTPVSEIPSEKEDPEGYE GIAQLM	303			
Query 305	DGLEVARGDSYWLPLEELGSLVNTNRNPVTESSRRKFLNQVNAHEDSWIAPSEWDRLA+TD	364			
Sbjct 304	EGLIARGDSYWLPLE+EI+GSLVNTNRNPVTESSRRKFLNQVNAHEDSWIAP+ WDRLA+TD	363			
Query 365	KALALQKDDRI TLGFDGSKSDWTALVACRVS DGMFLIKGNWPNEDYPHEVPRFEDVAV	424			
Sbjct 364	L+K+DRITLGFDFGSKSDWTALVACRVS DGMFL+ WNP DYPH+EVPR++VDVAV	423			
Query 425	VRSAFQRVDVVGFRADVKEFEAYVDQWRDFKRKLNATPGNPVAFDMRGQTKRFALDC	484			
Sbjct 424	VRSAFQRVDVVGFRADVKEFEAYVDQWRDFKRKLNATPGNPVAFDMRGQTKRFALDC	483			
Query 485	ERFVDAVIEHELHHDGNPVLRQHVLNARRHPTTYDAISIRKESKDSKKIDAACAVLAF	544			
Sbjct 484	ERFLDAVIEKELHHDGNPVLRQHVLNARRHPTTYDAISIRKESKDSKKIDAACAVLAF	543			
Query 545	GARQDYQMSKKHRS GAKAVIIR 566				
Sbjct 544	G+RDY MSKKHR G AVI+R 564				
	GSRQDYMSKKHRRGGG-AVIVR 564				

Download GenPept Graphics

terminase [Mycobacterium phage Obama12]

Sequence ID: [ref|YP_009007203.1](#) Length: 565 Number of Matches: 1[See 3 more title\(s\)](#)

Range 1: 1 to 565 GenPept Graphics Next Match Previous Match

Score	Expect	Method	Identities	Positives	Gaps
1045 bits(2703)	0.0	Compositional matrix adjust.	496/566(88%)	537/566(94%)	1/566(0%)
Query 1	MALEYLSTPELLPQPPHKGIPVWLCHEDGSWALPKYTLGWGLNWLAEYVRSFAGGGPI	60			
Sbjct 1	M+L PLLPQPPHKGIPVW EDGSW LP+ TLGWG+LNWLA+YVRSFAGGGPF+ MSLANHHFVLLPQPPHKGIPVWQVREDGSHLPERTLQWGLNWLAKYVRSFAGGGPFL	60			
Query 61	PTLEQARFILWYAVDENGVIYAYREGVLRMRKMGWKDPLCAAIALVELCGPVAFSHWDEK	120			
Sbjct 61	PTLEQARFILWYAVDE GVIYAYREGVLRMRKMGWKDPLCAAIAL ELCGPVAFSHW + PTLEQARFILWYAVDERGVIYAYREGVLRMRKMGWKDPLCAAIALAELCGPVAFSHWDL	120			
Query 121	GNPVGKRRHAANITIAAVSQDQTKNTFSLFPPVMSKQMKTEYGLDQVNFVIYTEDGGRIE	180			
	GNPVGK RHAANITIAAVSQDQTKNTFSLFPPVMSK+KT+YGLDQVNFVIY+E GGRIE				

The Mycobacteriophage Database

at PhagesDB.org

Home Phages Data BLAST Publications Resources Software Social Other DBs About

Search PhagesDB.org



Recently Added Phages

Savanna

Bryanna

RadDad

Dyanna

Daishi

Recently Modified Phages

HINder

Ferb

Biaggio

Cgwhichard

Filch

Recently Finished Phages

Cambiare (G)

AlanGrant (B4)

Corofin (B3)

Baee (B5)

Local Protein BLAST

Go to [Nucleotide BLAST](#)

This tool will run a local BLAST search against our protein databases. This includes all proteins from the most recent Phamerator update. Proteins marked "Draft" are from auto-annotated files.

Choose program to use and database to search

Program Database

Enter sequence below in **FASTA** format

```
AATSSPASMEGNRPTLVIENTQWVGVPDGNVNDGVAMDDVIEGNVSKIPGARKLAICN
AHIIPGNDTVAEKAYDHWQDILSGKAVDTGLMYDALEAPADTPVSEIPSEKEDPEGYEAGI
AQLMDGLEVARGDSYWLPLEEILGSVLNTRNPVTESSRRKFLNQVNAHEDSWIAPSEWDRL
AVTDKALALQKDDRI TLGFDGSKSDDWTALVACRVSDGMLFLKSWNPEDYPHEEVPRED
VDAVVRSAFORVDVVGFRADVKEFEAYVDQWGRDFKRKLINATPGNPPVAFDMRGQTKRF
ALDCERFVDAVIEHELHHDGNPVLRQHVNLNARRHPTTYDAISIRKESKDSSKKIDAAVCA
VLAFGARQDYQMSKHHRSKAVIIRZ
```

Or load it from disk

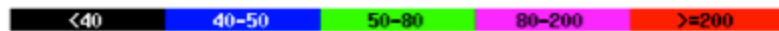
No file selected.

Set subsequence: From To

Distribution of 100 Blast Hits on the Query Sequence

Mouse-over to show define and scores. Click to show alignments

Color Key for Alignment Scores



Sequences producing significant alignments:	Score	E
	(bits)	Value
Sheen-Draft_14, function unknown, 566	<u>1171</u>	0.0
Timshel_12, Terminase., 564	<u>1053</u>	0.0
HINder_11, terminase, 564	<u>1048</u>	0.0
Obamal2-draft_10, function unknown, 565	<u>1046</u>	0.0
LHTSCC_11, Terminase, large subunit, 565	<u>1046</u>	0.0
Arturo_10, terminase, 565	<u>1046</u>	0.0
Phighter1804-Draft_10, function unknown, 565	<u>1046</u>	0.0
Eagle_10, Terminase, 565	<u>1046</u>	0.0
Camperdownii_Draft_10, function unknown, 565	<u>1046</u>	0.0
Flux_11, terminase, 565	<u>1045</u>	0.0
Dhanush_11, terminase, 565	<u>1045</u>	0.0
Wander_Draft_10, function unknown, 565	<u>1044</u>	0.0
TygerBlood-Draft_10, function unknown, 565	<u>1044</u>	0.0
TiroTheta9_11, terminase, 565	<u>1044</u>	0.0
TinaFeyge-Draft_gp10, function unknown, 565	<u>1044</u>	0.0
Shaka_11, terminase, 565	<u>1044</u>	0.0
Sabertooth_11, Terminase, 565	<u>1044</u>	0.0
Millski_Draft_10, function unknown, 565	<u>1044</u>	0.0
Melvin MELVIN_11, Terminase, 565	<u>1044</u>	0.0
MeeZee_11, Terminase, 565	<u>1044</u>	0.0
Medusa_11, terminase, 565	<u>1044</u>	0.0
Maverick-DRAFT_11, function unknown, 565	<u>1044</u>	0.0
Lemur_Draft_10, function unknown, 557	<u>1044</u>	0.0
Kratark_Draft_10, function unknown, 565	<u>1044</u>	0.0
KFPoly_Draft_10, function unknown, 565	<u>1044</u>	0.0
Kampy-Draft_10, function unknown, 565	<u>1044</u>	0.0
Holli-draft_10, function unknown, 565	<u>1044</u>	0.0
HamSlice-Draft_10, function unknown, 565	<u>1044</u>	0.0
Gadost_Draft_10, function unknown, 565	<u>1044</u>	0.0
Funston_Draft_DRAFT_10, function unknown, 565	<u>1044</u>	0.0
Eris_Draft_10, function unknown, 565	<u>1044</u>	0.0

>Timshel_12, Terminase., 564
Length = 564

Score = 1053 bits (2724), Expect = 0.0
Identities = 504/562 (89%), Positives = 531/562 (94%), Gaps = 1/562 (0%)

Query: 5 YLSTEP LLPQPPHKIGPVWLCHEDGSWALPKYTLGWGVLNWLAEYVRS PAGGGPFIPTLE 64
YL+ PLLPQPPHKIGPVW HEDGSWALP TLGWGVLNWLAEYVRS PAGGGPFIPTLE
Sbjct: 4 YLNPGLLPQPPHKIGPVWQVHEDGSWALPARTL GWGVLNWLAEYVRS PAGGGVFIPTLE 63

Query: 65 QARFILWYAVDENG VYAYREGVLRMRMGWGKDPLCAAIALVELCGPVAFSHWDEKGNPV 124
QARFILWYAVDE G YAYREG LRRMRMGWGKDPLCAAIALVELCGPVAFSHWD G+PV
Sbjct: 64 QARFILWYAVDERGNYAYREGCLRRMRMGWGKDPLCAAIALVELCGPVAFSHWDLGSPV 123

Query: 125 GKRRHAAWITIAAVSQDQTKNTFSLFPVMISKQMKTEYGLDVKNFVIYTEDGGRIEAATS 184
GK RHAAWIT+AAVSQDQTKNTFS+FPVMISK+MK +YGLDVKNFVIY+E+GGRIEAATS
Sbjct: 124 GKPRHAAWITVAAVSQDQTKNTFSMPFVMISKMKKVDYGLDVKNFVIYSEEGGRIEAATS 183

Query: 185 SPASMEGNRPTLV IENETQWVGVPDGNVNDGVAMDDVIEGNVSKIPGARKLAICNAHIP 244
SPASMEGNRPTLV IENETQWVGVPDGNVNDG AMDDVIEGNVSKIPGARKLAICNAHIP
Sbjct: 184 SPASMEGNRPTLV IENETQWVGVPDGNVNDGPAMDDVIEGNVSKIPGARKLAICNAHIP 243

Query: 245 GNDTVAEKAYDHWQDILSGKAVDTGLMYDALEAPADTPVSEIPSEKEDPEGYEAGIAQLM 304
GNDTVAEKAYDHWQDILSGKAVDTG+MYDALEAPADTPVSEIPSEKEDPEGYE GIAQLM
Sbjct: 244 GNDTVAEKAYDHWQDILSGKAVDTGIMYDALEAPADTPVSEIPSEKEDPEGYERKIAQLM 303

Query: 305 DGLEVARGDSYWLPLEEILGSVLNTRNPVTSERRKFLNQVNAHEDSWIAPSEWDRDLAVTD 364
+GLE+ARGDSYWLPL+EI+GSVLNTRNPVTSERRKFLNQVNAHEDSWIAP+ WDRLA+TD
Sbjct: 304 EGLEIARGDSYWLPLDEIMGSVLNTRNPVTSERRKFLNQVNAHEDSWIAPAWDRDLALTD 363

Query: 365 KALALQKDDRITLGF DGSKSDWTALVACRVSDGMLFLIKSWNPEDYPHEVPRDVEDAV 424
L+K+DRITLGF DGSKSDWTALVACRVSDGMLFL+ WNP DYPH+EVPR++VDVAV
Sbjct: 364 PLFLKKNDRITLGF DGSKSDWTALVACRVSDGMLFLLDKWNPN DYPHDVEDVPRDVEDAV 423

Query: 425 VRSFAFQRYDVVGF RADVKEFEAYVDQWGRDFKRKLNATPGNPVAFDMRGQTKRFALDC 484
VRSFAFQRYDVVGF RADVKEFEAYVDQWGRDFKRKLNATPGNPVAFDMRGQTKRFALDC
Sbjct: 424 VRSFAFQRYDVVGF RADVKEFEAYVDQWGRDFKRKLNATPGNPVAFDMRGQTKRFALDC 483

Query: 485 ERFVDAVIEHELHHDGNPVL RQHVLNARRHPTTYDAISIRKESKDSKKIDAAVCAVLAF 544
ERF+DAVIE EL HD NPVL RQHVLNARRHPTT+DAISIRKESKDSKKIDAAVCAVLAF
Sbjct: 484 ERFDAVIEKELWHDQNPVL RQHVLNARRHPTTFDAISIRKESKDSKKIDAAVCAVLAF 543

Query: 545 GARQDYQMSK KHRSGAKAVIR 566
G+RQDY MSKKHR G AVI+R
Sbjct: 544 GSRQDYLM SKKHRGGG-AVIVR 564

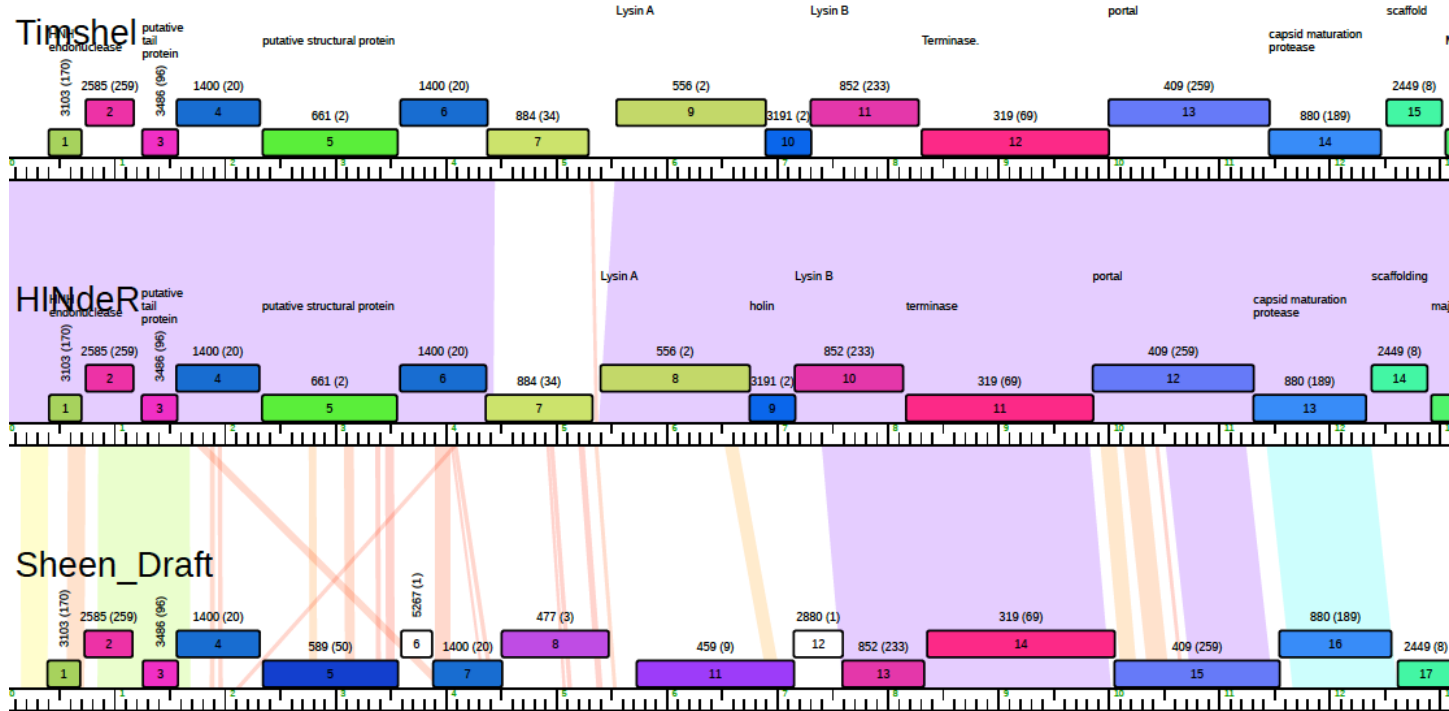
>HINder_11, terminase, 564
Length = 564

Score = 1048 bits (2709), Expect = 0.0
Identities = 501/562 (89%), Positives = 530/562 (94%), Gaps = 1/562 (0%)

Query: 5 YLSTEP LLPQPPHKIGPVWLCHEDGSWALPKYTLGWGVLNWLAEYVRS PAGGGPFIPTLE 64
YL+ PLLPQPPHKIGPVW HEDGSWALP TLGWGVLNWLAEYVRS PAGGGPFIPTLE
Sbjct: 4 YLNPGLLPQPPHKIGPVWQVHEDGSWALPARTL GWGVLNWLAEYVRS PAGGGPFIPTLE 63

Query: 65 QARFILWYAVDENG VYAYREGVLRMRMGWGKDPLCAAIALVELCGPVAFSHWDEKGNPV 124
QARFILWYAVDE G YAYREG LRRMRMGWGKDPLCAAIALVELCGPVAFSHWD G+PV
Sbjct: 64 QARFILWYAVDERGNYAYREGCLRRMRMGWGKDPLCAAIALVELCGPVAFSHWDLGSPV 123

Phamerator map



Starterator data

Note: In the above figure, yellow indicates the location of called starts comprised solely of computational predictions (i.e. auto-annotations by Glimmer/GeneMark), green indicates the location of called starts with at least 1 manual gene annotation.

Pham 4391 Report

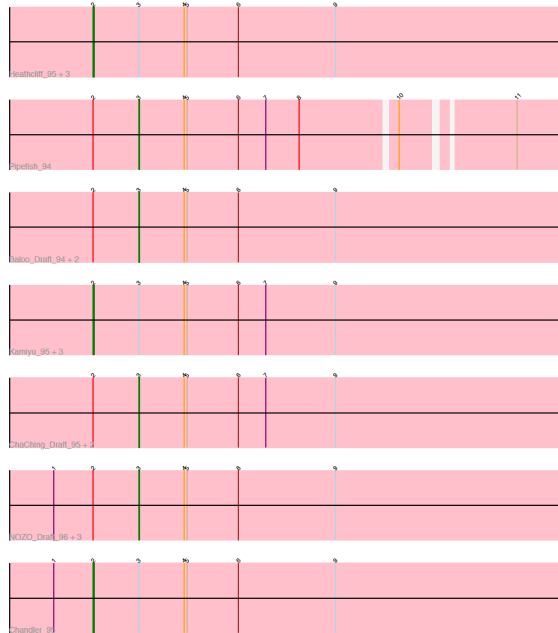
This analysis was run 11/27/16.

Pham number 4391 has 20 members, 6 are drafts.

Phages represented in each track:

- Track 1 : Heathcliff_95, Bernardo_92, Akoma_96, Audrey_95
- Track 2 : Pipefish_94
- Track 3 : Baloo_Draft_94, Phaerdrus_90, Mortcellus_Draft_96
- Track 4 : Kamiyu_95, Athena_97, Daisy_94, Corofin_95
- Track 5 : ChaChing_Draft_95, Yahalom_Draft_91, Phlyer_95
- Track 6 : NOZO_Draft_96, Composita_Draft_98, Gadjet_93, OrangeOswald_94
- Track 7 : Chandler_95

Pham 4391



Summary of Final Annotations (Info on gene starts based on numbers in diagram):

The start number called the most often in the published annotations is start number 2, it was called in 9 of the 14 non-draft genes in the pham.

Genes that call this "Most Annotated" start:

- Heathcliff_95, Kamiyu_95, Akoma_96, Bernardo_92, Athena_97, Chandler_95, Audrey_95, Daisy_94, Corofin_95,

Genes that have the "Most Annotated" start but do not call it:

- Pipefish_94, Baloo_Draft_94, Phaerdrus_90, Yahalom_Draft_91, NOZO_Draft_96, Composita_Draft_98, Phlyer_95, OrangeOswald_94, Gadjet_93, Mortcellus_Draft_96, ChaChing_Draft_95,

Genes that do not have the "Most Annotated" start:

•

Summary by start number:

- Start number 2 is called in: Heathcliff_95, Kamiyu_95, Akoma_96, Bernardo_92, Athena_97, Chandler_95, Audrey_95, Daisy_94, Corofin_95,

Percent with start 2 called: 45.0%

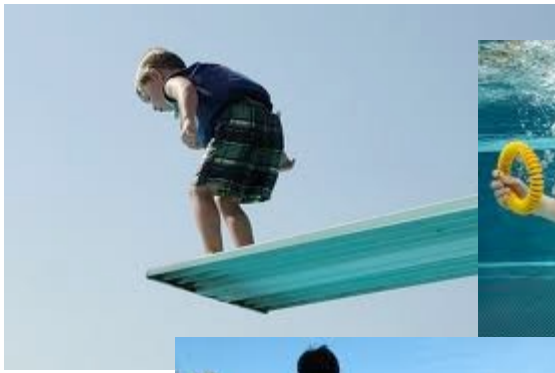
- Start number 3 is called in: Pipefish_94, Baloo_Draft_94, Phaerdrus_90, Yahalom_Draft_91, NOZO_Draft_96, Composita_Draft_98, Phlyer_95, OrangeOswald_94, Gadjet_93, Mortcellus_Draft_96, ChaChing_Draft_95,

Percent with start 3 called: 55.0%

GUIDING PRINCIPLES OF BACTERIOPHAGE GENOME ANNOTATION

- Found in “Phage Annotation, Genomics and Data Interpretation” Section of the Bioinformatics Guide
- 15 Key Directives
- Read for tomorrow

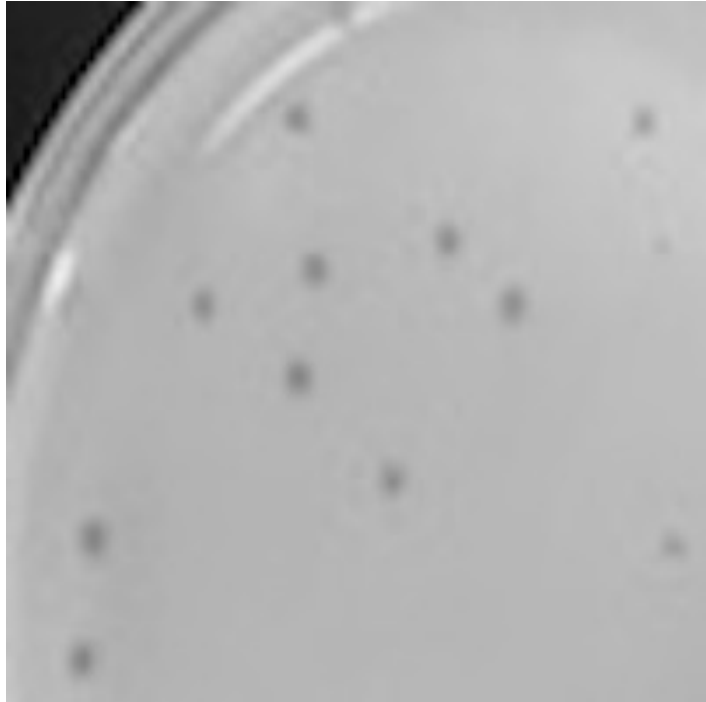
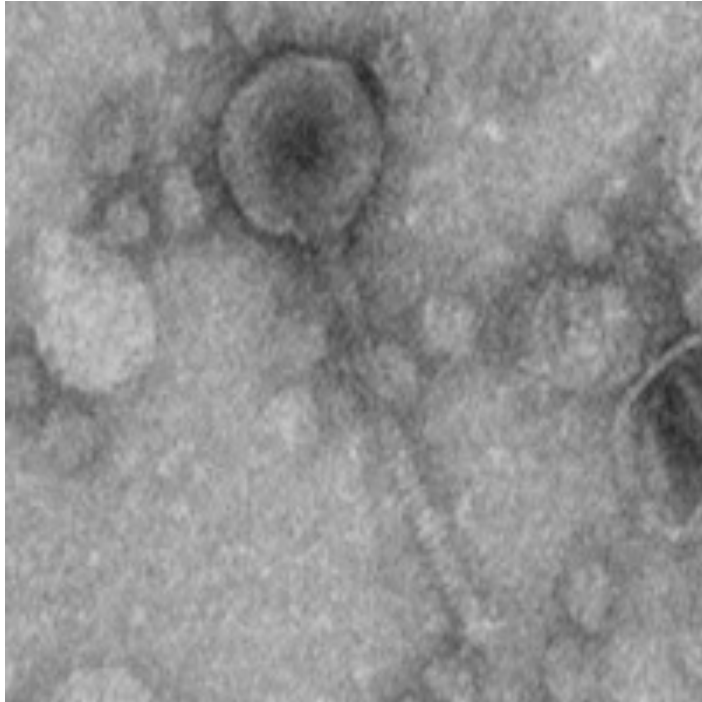
<https://seaphagesbioinformatics.helpdocsonline.com/guiding-principles>



Let's get started!

1. Gather Data
2. Auto-annotate in DNA Master
3. Gene Calling
4. Functional Assignments

Microbacterium Jefe



Tonight's Tasks:

Annotation Outline

- Outline Introduction
- Setting Up Your Computer
 - DNA Master
 - Installing DNA Master ✓
 - Updating DNA Master ✓
 - Setting DNA Master preferences ←
 - Websites to bookmark ←
- Surveying Your Genome
 - Retrieving your genome sequence ←
 - Comparing your genome's sequence ✓
 - Comparing your genome's genes ✓
 - Clustering your genome
- Gathering Data
 - Creating a coding potential graph ←
 - Making a Phamerator map
 - Guiding principles of phage annotation
- Automatically Annotating Your Genome
 - Auto-annotation ←

DNA Master
Current Build
2705

← Tonight's reading

Complete Genome Blastp in DNA Master ←

Tips for DNA Master Files

- Save .dnam5 file often
- Save .dnam5 file as a new name. (Then don't save the old named one.)
- If working in a virtual machine, be mindful about closing/shutting down
- Did I mention to save your files often using a new name?