



MAG5

Metagenomic Assembly Analysis

By David Podorefsky
Na Liu
Ketaki Malaviya
Nikola Malinov
Julian Elmurib
Robert Davison

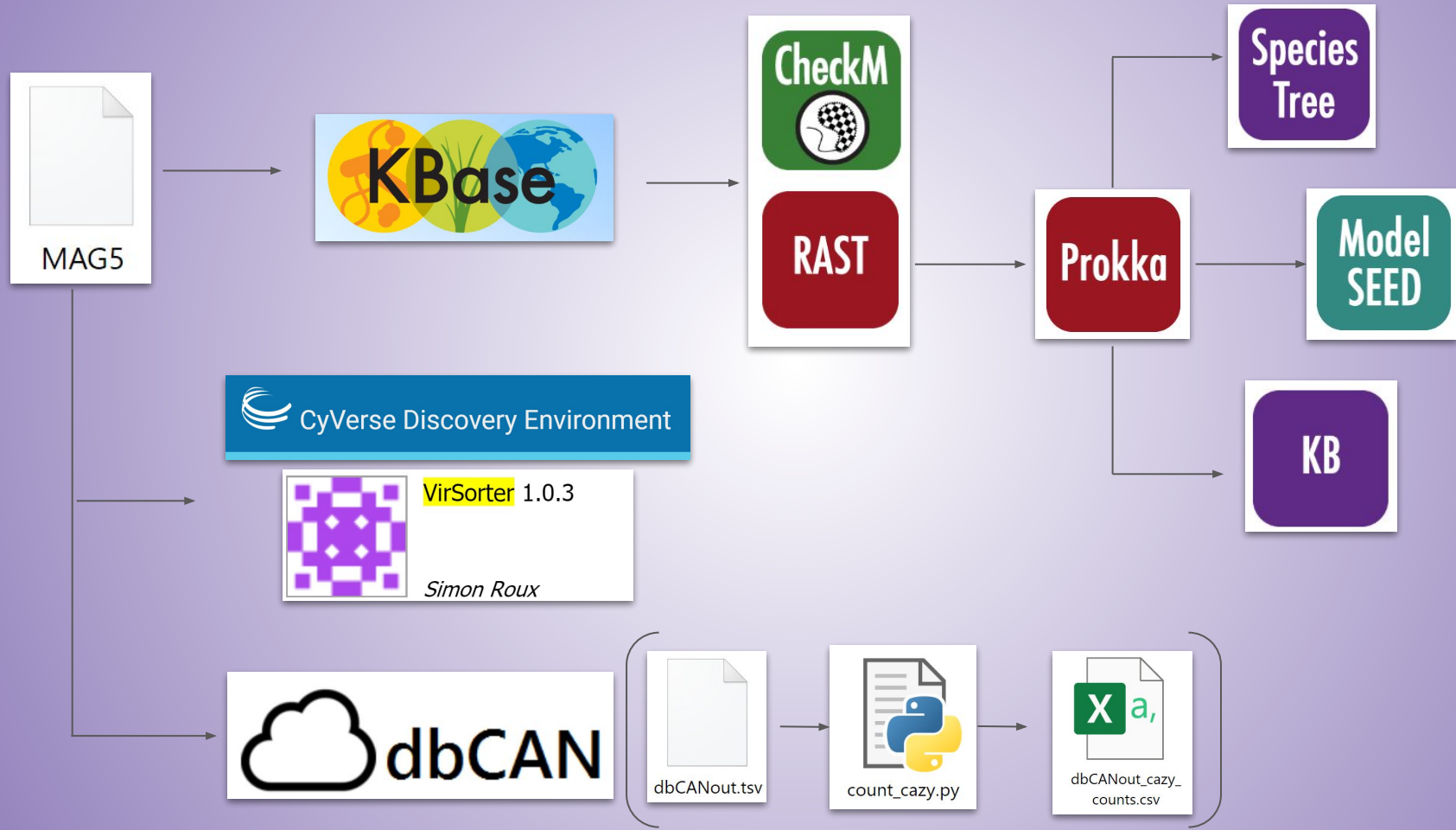
Introduction

- Possible untapped **biotech potential** in goat rumen
 - Biofuels, Antibiotics, etc.
- Identify and research potential applications of MAG5



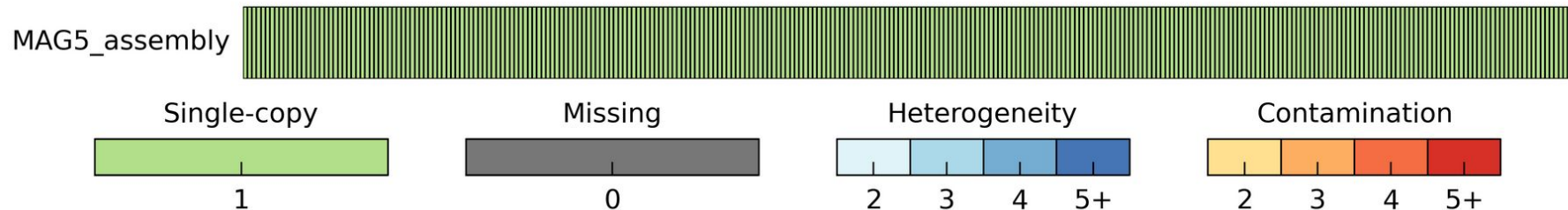
San Clemente Island Goats

Methods



Genome Quality

Figure 1.
Genome Quality Plot



Note: Plot of phylogenetic marker genes. Green bars represent a single copy of a marker, while a grey bar indicates a missing clade-specific marker.

Table 1.
Quality Check

Bin Name	Marker Lineage	# Genomes	# Markers	# Marker Sets	0	1	2	3	4	5+	Completeness	Contamination
MAG5_assembly	p__Firmicutes	100	295	158	1	294	0	0	0	0	99.367	0.0

Note: The KBase app Check M v1.0.8 was used to assess the quality of the MAG5 genome.

Clusters of Orthologous Genes (COGs)

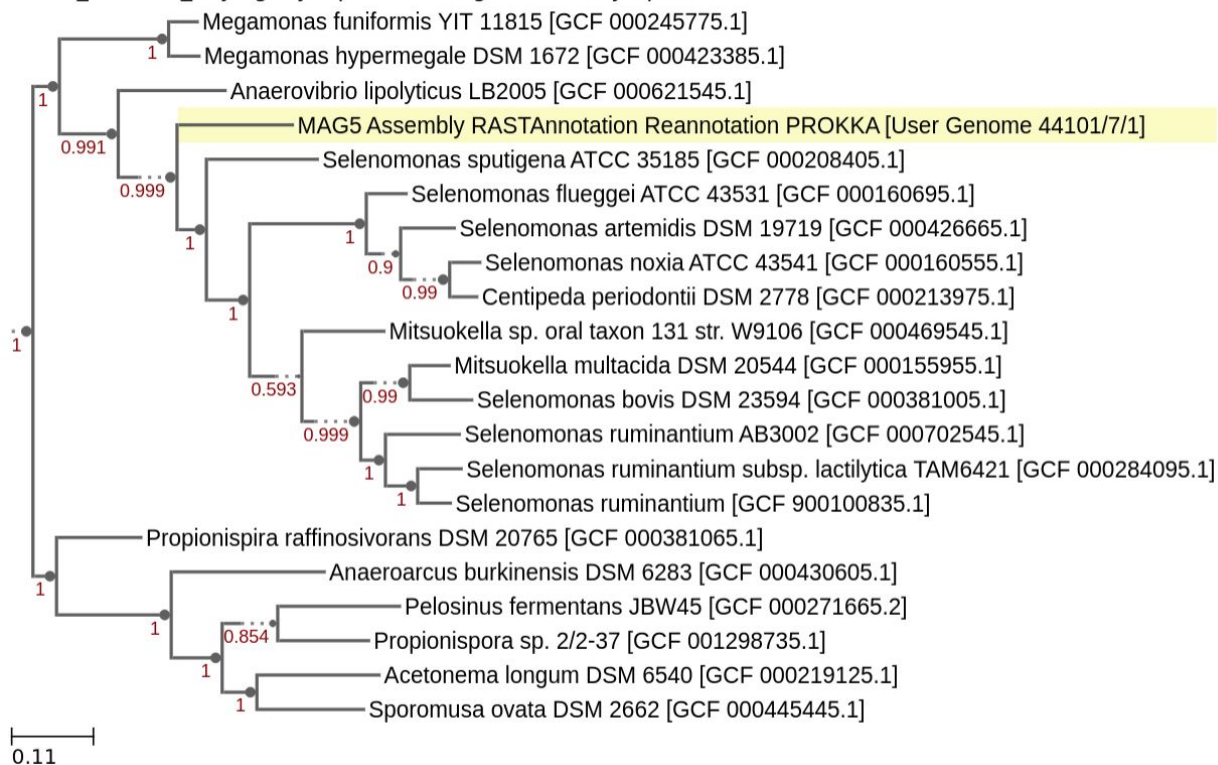
- The **Insert Genome Into Species Tree 2.1.10** App uses **49** COGs when building the species tree

- These COGs are predominantly ***ribosomal subunits***
- Ribosomes are used for ***phylogenetic analysis*** they're present in all organisms
- The COGs used for ***Energy acquisition*** are much broader and cannot be compared from species to species to produce the phylogenetic tree

COG0012	COG0012	Predicted GTPase, probable translation factor [Translation ribosomal structure and biogenesis].
COG0013	AlaS	Alanyl-tRNA synthetase [Translation, ribosomal structure and biogenesis].
COG0016	PheS	Phenylalanyl-tRNA synthetase alpha subunit [Translation, ribosomal structure and biogenesis].
COG0018	ArgS	Arginyl-tRNA synthetase [Translation, ribosomal structure and biogenesis].
COG0030	KsgA	Dimethyladenosine transferase (rRNA methylation) [Translation, ribosomal structure and biogenesis].

Phylogenetic Tree

MAG5_PROKKA_Phylogeny: Species Tree generated by Species Tree Builder



Phylogeny

- RNAmmer used to identify ribosomal RNA
 - 5s sub-unit identified
- The obtained sequence was blasted

```
>rRNA_NODE_12_length_459105_cov_16.432523_254-369_DIR+ /molecule=5s_rRNA /score=72.6
CGGTGGTGATGCCTACGTGGTTCCACCTGTTCCCATTCGGAACACAGTAGTTAAGCACGT
AAAGGCCGAAAGTACTTGGCTGGAGACGGCCCGGAGGATAGGAAGCTGCCGGTTA
MAG5_rRNA_seq.fa (END)
```

Sequences producing significant alignments:

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

[Alignments](#) [Download](#) [GenBank](#) [Graphics](#) [Distance tree of results](#)

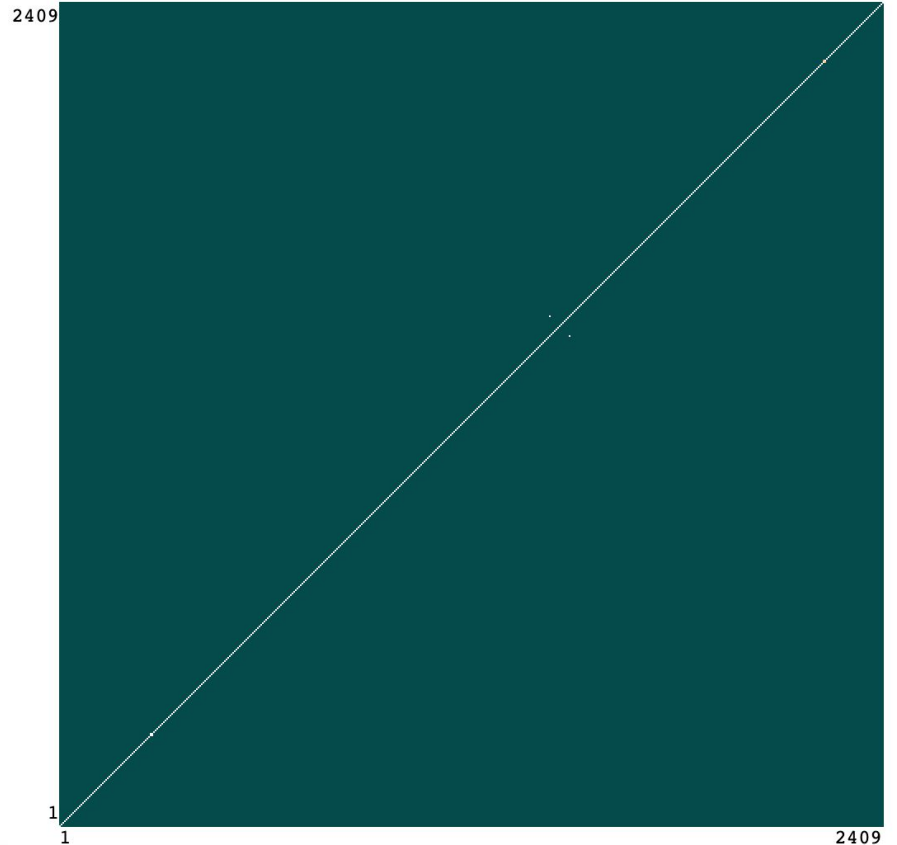
	Description	Max Score	Total Score	Query Cover	E value	Per. Ident	Accession
<input type="checkbox"/>	Veillonella rodentium strain NCTC12018 genome assembly, chromosome: 1	137	551	89%	6e-29	90.38%	LT906470.1
<input type="checkbox"/>	Christensenella sp. Marseille-P2438 genome assembly, chromosome: contig00001	137	137	97%	6e-29	88.50%	LT700187.1
<input type="checkbox"/>	Uncultured bacterium BD contig02516 genomic sequence	106	106	77%	2e-19	87.91%	KJ799706.1
<input type="checkbox"/>	Selenomonas sp. oral taxon 126 strain W7667 genome	102	371	96%	2e-18	83.04%	CP016201.1
<input type="checkbox"/>	Selenomonas sp. oral taxon 920 strain W5150 chromosome, complete genome	97.1	383	96%	1e-16	82.14%	CP017042.1
<input type="checkbox"/>	Selenomonas sp. oral taxon 136 strain F0591, complete genome	97.1	269	96%	1e-16	82.14%	CP014239.1
<input type="checkbox"/>	Eubacterium siraeum V10Sc8a strain V10Sc8a 5S ribosomal RNA, complete sequence	97.1	97.1	81%	1e-16		

Questions/comm

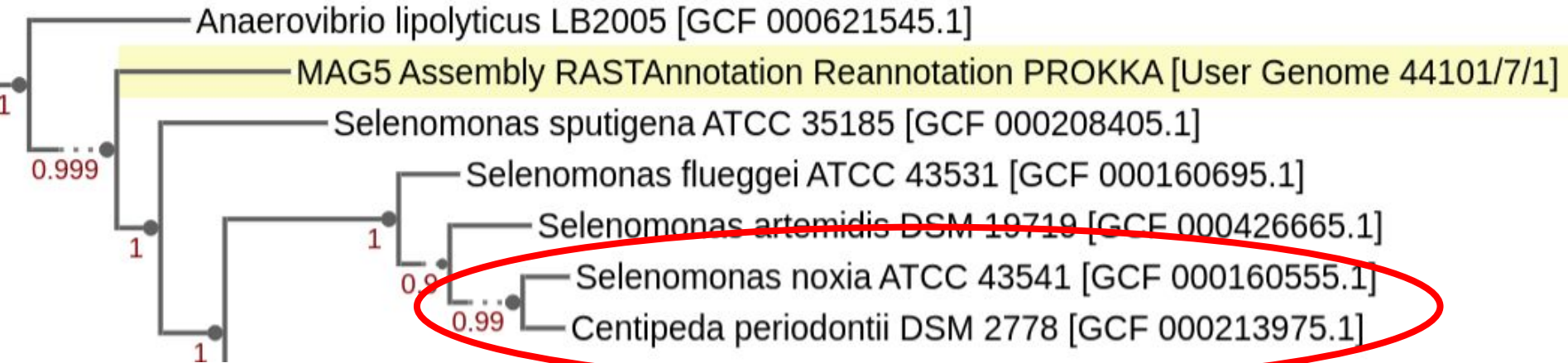
Synteny Plot Comparison

An identical match
would look like this:

MAG5 compared to MAG5



Example of two known close relatives



Selenomonas Noxia ATCC

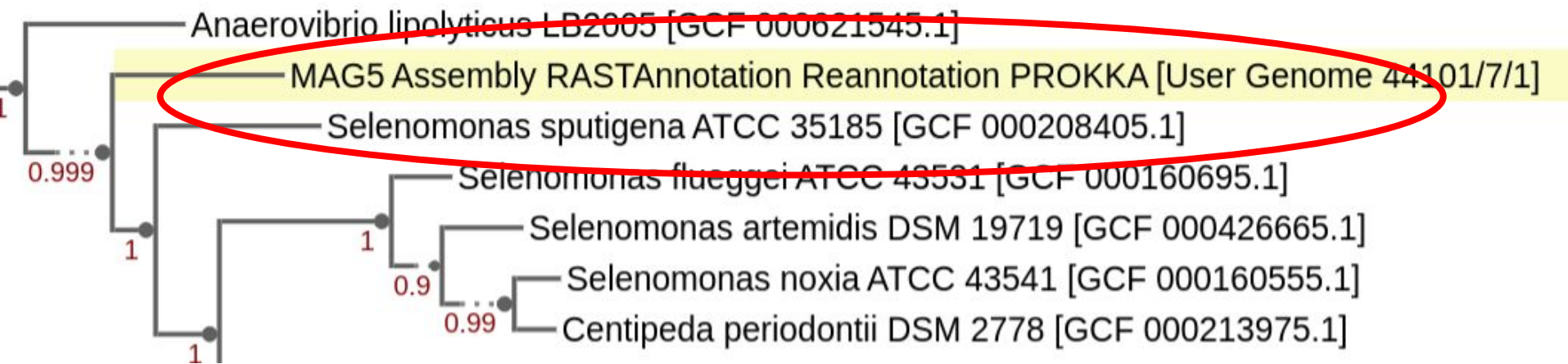
Centipeda Periodontii DSM

2453

1

1

1876



Our MAG5 compared to its closest relative *Selenomonas* *Sputigena*

MAG5:

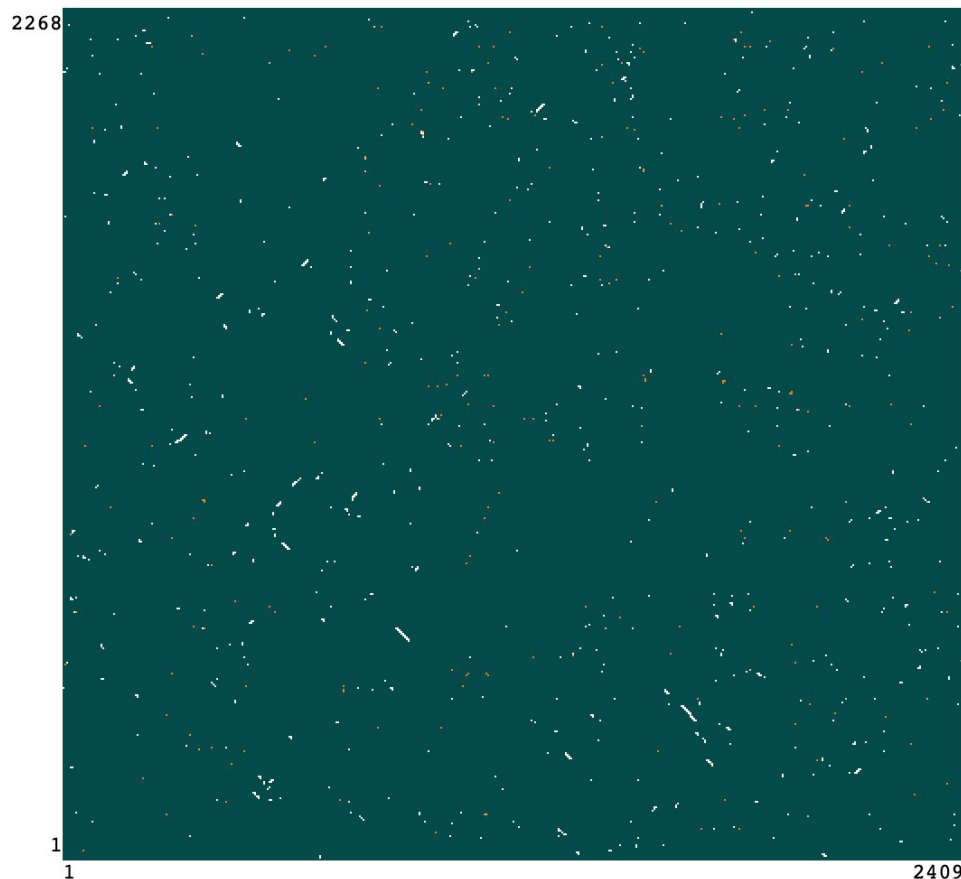
2409 genes

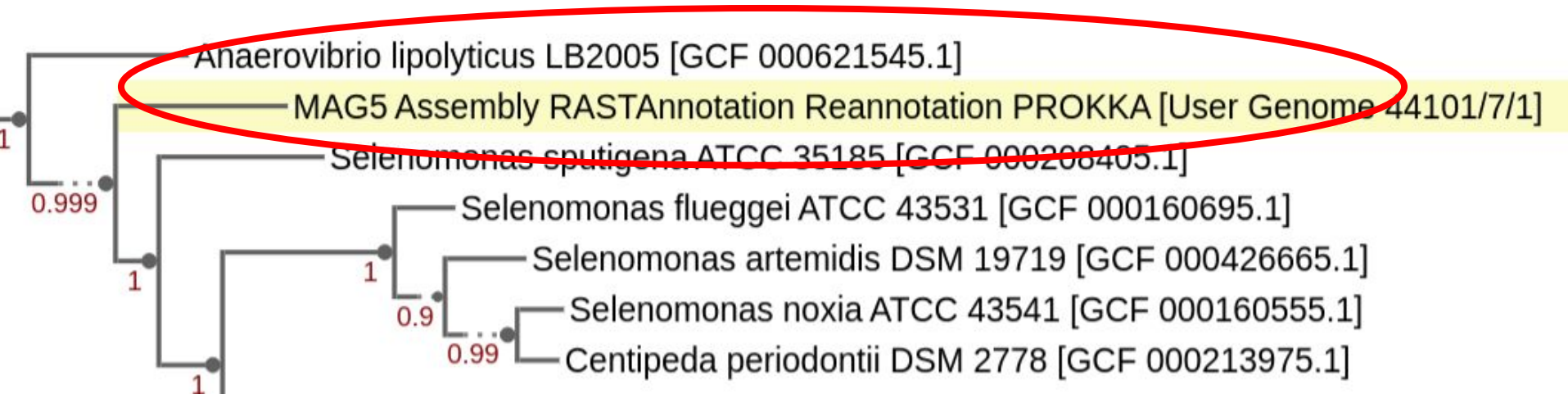
1675 hits

Selenomonas Sputigena:

2268 genes

1641 hits



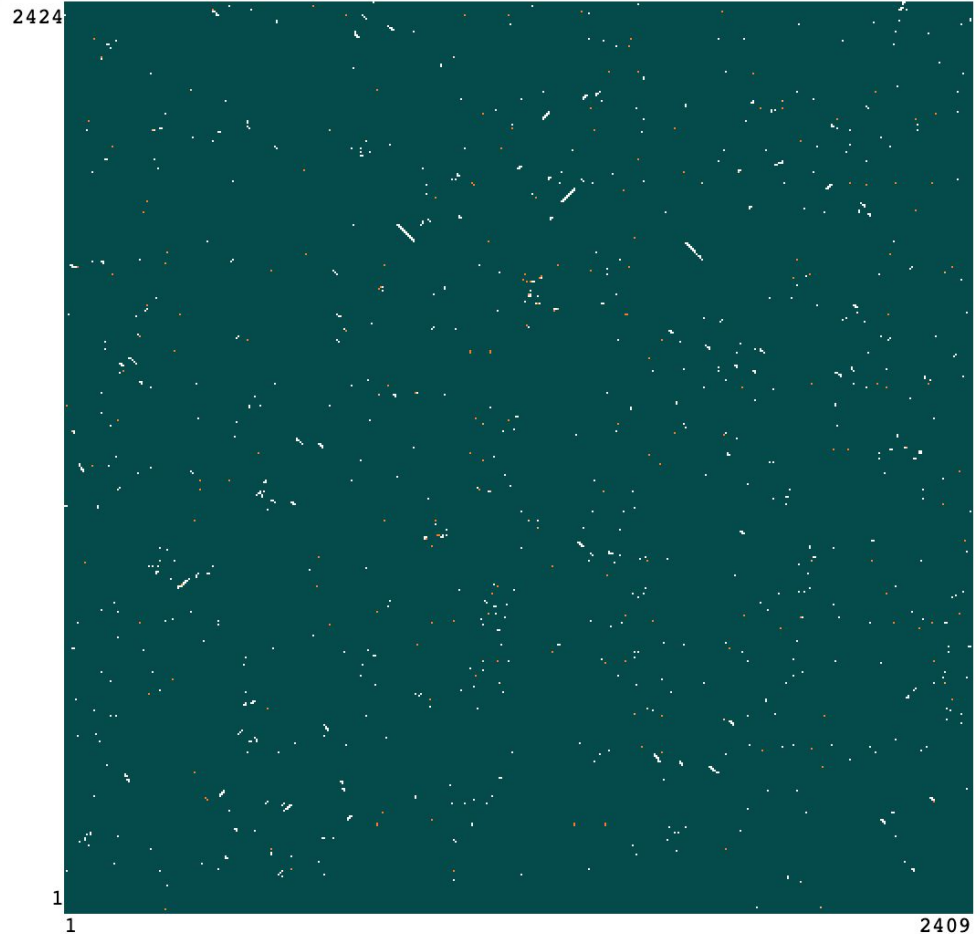


MAG5 compared to a
second close relative
slightly further on the
tree.

MAG5:
2409 genes
1712 hits

Anaerovibrio Lipolyticus:
2424 genes
1739 hits

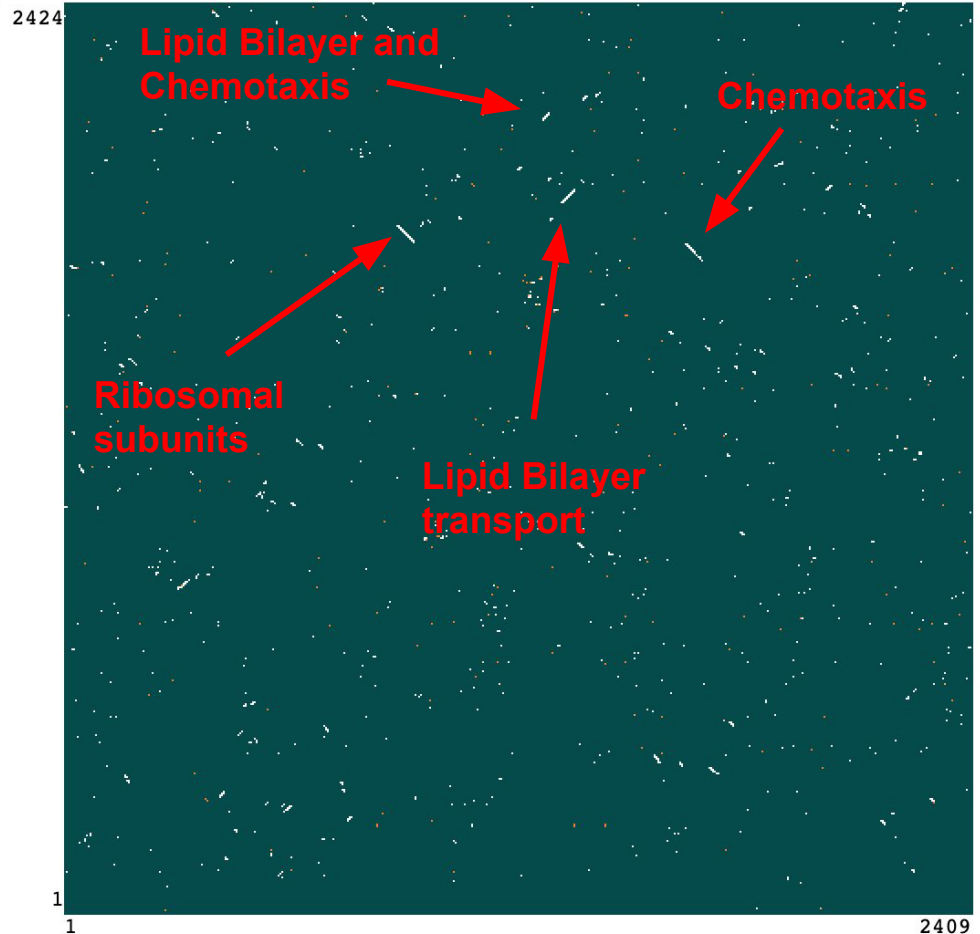
More Hits!



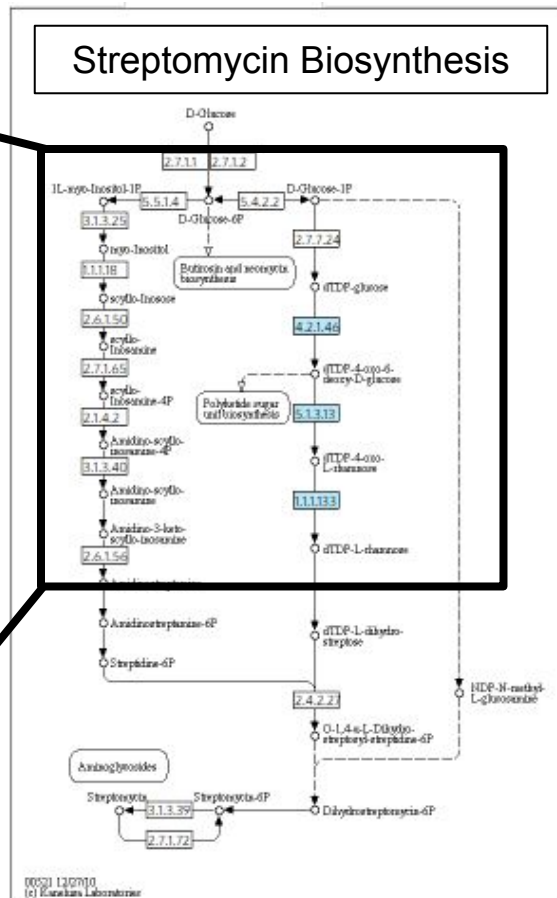
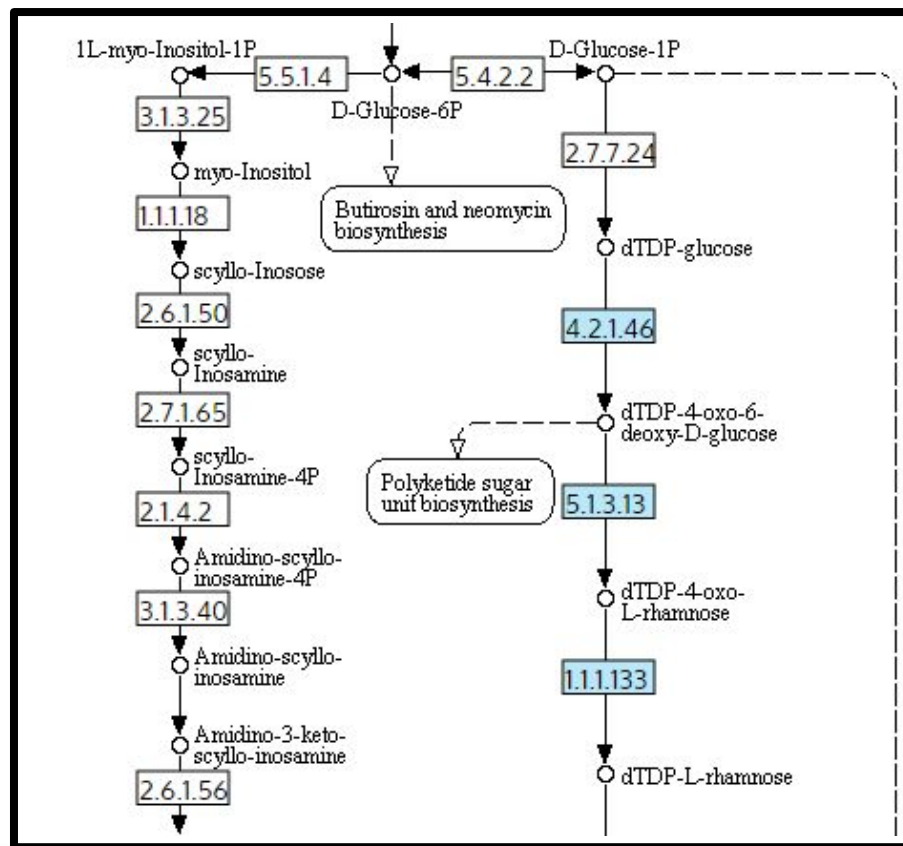
MAG5 compared to a second close relative slightly further on the tree.

MAG5:
2409 genes
1712 hits

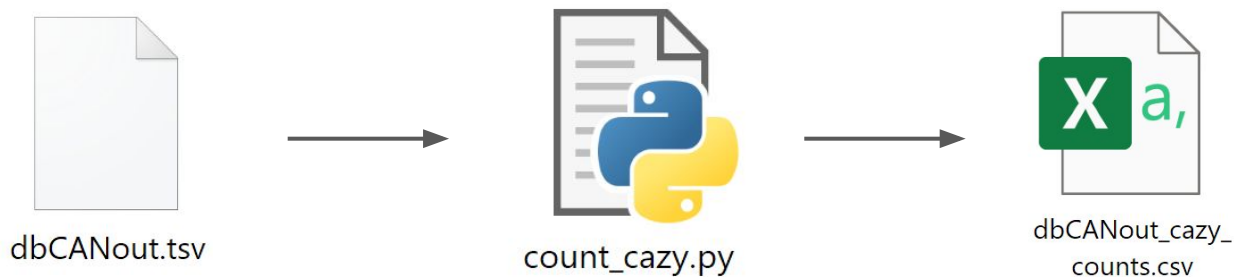
Anaerovibrio Lipolyticus:
2424 genes
1739 hits



Weak Evidence Suggests No Potential for Antibiotic Production



CAZyme Pipeline



1. Determines CAZymes with 3 hidden Markov model (HMM) variations
2. Utilize python script to keep those present in all three databases - credit to Tom
3. Count and realize enzymes

Hidden Markov Model

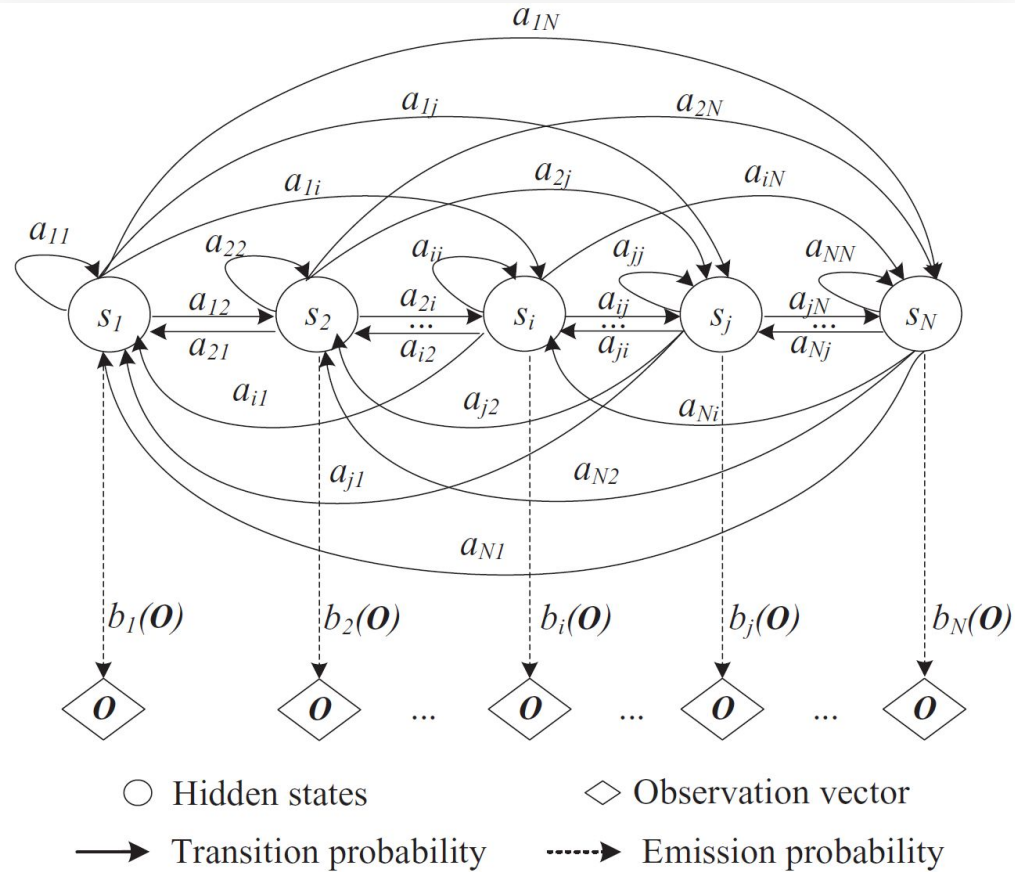


Fig. 5. The relationship of hidden states and observations in continuous HMM.



dbCAN meta server: automated CAZyme annotation

Carbohydrate-active enzyme

[Home](#)[Annotate](#)[Download](#)[Help](#)[About us](#)

Cite us: [NAR/gky418](#) and [gky479](#)

DIAMOND:

Fast blast hits in CAZy database
(E-Value < 1e-102)

Diamond

HMMER



HMMER:

Domain boundaries according to
the dbCAN CAZyme domain
HMM database
(E-value < 1e-15, coverage > 0.35)

Hotpep

Hotpep:

Short conserved motifs in PPR library
(Frequency > 2.6, Hits > 6)

Download SignalP output **Download Prodigal predictions** **Download this table** (keep those with # of Tools >=2 will give you best result; and use dbCAN domain assignment is recommended) ?

Show 15 ▼ entries

Search:

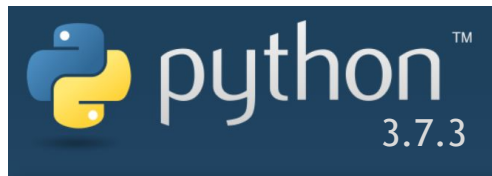
Gene ID	# of Tools	HMMER	DIAMOND	Hotpep	Signal Peptide
NODE_107_length_200687_c1 ov_16.378152_148		N	GH6	N	N
NODE_107_length_200687_c1 ov_16.378152_162		N	GT13	N	N
NODE_112_length_196477_c2 ov_17.935031_129		CE4 (77-234)+ GH153 (281-593)	CE4	N	Y (1-25)
NODE_112_length_196477_c3 ov_17.935031_130		GT2_Glyco_tranf_2_3 (55-275)	GT2	GT2	N
NODE_112_length_196477_c1 ov_17.935031_24		CE10 (72-238)	N	N	Y (1-35)
NODE_12_length_459105_co1 v_16.432523_146		N	GH13_30	N	N
NODE_12_length_459105_co3 v_16.432523_155		GT51 (69-244)	GT51	GT51	N
NODE_12_length_459105_co2 v_16.432523_185		GH23 (454-560)	GH23	N	Y (1-22)
NODE_12_length_459105_co1 v_16.432523_260		GT28 (207-356)	N	N	Y (1-21)
NODE_12_length_459105_co3 v_16.432523_295		GT2_Glyco_tranf_2_3 (45-268)	GT2	GT2	N
NODE_12_length_459105_co2 v_16.432523_298		GT26 (58-226)	N	GT26	N
NODE_12_length_459105_co1 v_16.432523_391		GT2_Glycos_transf_2 (2-165)	N	N	N
NODE_12_length_459105_co2 v_16.432523_468		GH73 (123-249)	N	GH73	Y (1-29)
NODE_207_length_144259_c2 ov_18.029109_10		GT9 (74-322)	GT9	N	N
NODE_207_length_144259_c1 ov_18.029109_2		GT2_Glycos_transf_2 (13-113)	N	N	N

Count CAZymes by Class

```

Gene ID  RM088  Motzger  DIAMOND  SignalS  #ofTools
NODE_112_length_196477_cov_17.935031_125  CE4(77-234)+GH153(281-593)  N  CE4  Y(1-25)  2
NODE_112_length_196477_cov_17.935031_130  GT2_Glyco_transf_2_3(55-275)  GT2  GT2  N  3
NODE_112_length_196477_cov_17.935031_24  CE10(72-238)  N  N  Y(1-35)  1
NODE_112_length_459105_cov_16.432523_165  GT51  GT51  N  3
NODE_112_length_459105_cov_16.432523_185  GH84  GH84  N  2
NODE_12_length_459105_cov_16.432523_260  GT28  GT28  N  N  Y(1-22)  2
NODE_12_length_459105_cov_16.432523_240  GT28  GT28  N  N  Y(1-21)  1
NODE_12_length_459105_cov_16.432523_295  GT2_Glyco_transf_2_3(45-268)  GT2  GT2  N  3
NODE_12_length_459105_cov_16.432523_298  GT28  GT28  N  N  Y(1-25)  2
NODE_12_length_459105_cov_16.432523_391  GT2_Glyco_transf_2_2(2-165)  N  N  N  1
NODE_12_length_459105_cov_16.432523_468  GH78(1123-249)  GH78  N  Y(1-29)  2
NODE_207_length_144259_cov_18.029109_10  GT9  GT9  N  2
NODE_207_length_144259_cov_18.029109_2  GT2_Glyco_transf_2(113-133)  N  N  N  1
NODE_207_length_144259_cov_18.029109_26  GT4  GT4  N  3
NODE_207_length_144259_cov_18.029109_27  GT4  GT4  N  1
NODE_207_length_144259_cov_18.029109_28  GT2_Glyco_transf_2(4-176)  GT2  N  2
NODE_207_length_144259_cov_18.029109_3  GT9(76-340)  N  GT9  N  2
NODE_207_length_144259_cov_18.029109_30  GT2_Glyco_transf_2(9-127)  N  N  N  1
NODE_207_length_144259_cov_18.029109_49  GT4  GT4  N  1
NODE_207_length_144259_cov_18.029109_8  GT8  GT8  N  N  1
NODE_207_length_144259_cov_18.029109_11  GT2_Glyco_transf_2(5-158)  GT2  GT2  N  3
NODE_207_length_144259_cov_18.029109_38  GT35  GT35  N  3
NODE_207_length_144259_cov_18.029109_89  GT35  GT35  N  3
NODE_207_length_144259_cov_18.029109_89  CBM48  N  3
NODE_207_length_144259_cov_18.029109_89  GH13+CBM48  CBM48  N  3
NODE_207_length_144259_cov_18.029109_90  GT8  GT8  N  1
NODE_207_length_144259_cov_18.029109_99  GT33  GT33  N  3
NODE_239_length_131518_cov_22.771415_43  CE4  GT9  N  Y(1-26)  2
NODE_25_length_383736_cov_17.808888_19  CE28  CE28  N  2
NODE_25_length_383736_cov_17.808888_130  GT30  GT30  N  3
NODE_25_length_383736_cov_17.808888_132  GT19  GT19  N  3
NODE_25_length_383736_cov_17.808888_134  CE11  CE11  N  3
NODE_25_length_383736_cov_17.808888_18  GT9  GT9  N  2
NODE_25_length_383736_cov_17.808888_19  GT6  GT6  N  1
NODE_25_length_383736_cov_17.808888_214  CE10  CE10  N  1
NODE_25_length_383736_cov_17.808888_222  GT41  GT41  N  2
NODE_25_length_383736_cov_17.808888_230  GH84  GH84  N  3
NODE_25_length_383736_cov_17.808888_27  GT2_Glyco_transf_2(5-119)  N  N  1
NODE_25_length_383736_cov_17.808888_299  GH5  GH5  N  1
NODE_25_length_383736_cov_17.808888_30  GT2_Glyco_transf_2(19-146)  GT2  N  2
NODE_25_length_383736_cov_17.808888_31  GT4  GT4  N  1
NODE_25_length_383736_cov_17.808888_38  GT9  GT9  N  1
NODE_25_length_383736_cov_17.808888_47  GT2_Glyco_transf_2_4(21-108)  N  N  1
NODE_25_length_383736_cov_17.808888_48  GT2_Glyco_transf_2(6-174)  N  N  1
NODE_25_length_383736_cov_17.808888_49  GT2_Glyco_transf_2_3(7-218)  N  N  1
NODE_25_length_383736_cov_17.808888_56  GT2_Glyco_transf_2_3(8-197)+GT2_Glyco_transf_2(230-363)  N  N  1
NODE_25_length_383736_cov_17.808888_57  GT2_Glyco_transf_2_3(42-220)  N  N  1
NODE_25_length_383736_cov_17.808888_60  GT41  GT41  N  3
NODE_25_length_383736_cov_17.808888_62  GT41  GT41  N  3
NODE_25_length_383736_cov_17.808888_63  GT41  GT41  N  2
NODE_25_length_383736_cov_17.808888_64  GT2_Glyco_transf_2(5-101)+GT41(931-993)+GT41(924-1454)  GT41  GT41  N  3
NODE_25_length_383736_cov_17.808888_65  GT2_Glyco_transf_2(5-101)  GT2  N  2
NODE_25_length_383736_cov_17.808888_66  GT2_Glyco_transf_2(14-107)+GT2_Glyco_transf_2(345-485)  GT2  GT2  N  3
NODE_25_length_383736_cov_17.808888_67  GT2_Glyco_transf_2(6-118)  GT2  N  2
NODE_20_length_375312_cov_16.951897_124  GH83  GH83  N  2
NODE_20_length_375312_cov_16.951897_170  GT28  GT28  N  3
NODE_20_length_375312_cov_16.951897_178  CE4  CE4  N  2
NODE_20_length_375312_cov_16.951897_219  GT9  GT9  N  1
NODE_20_length_375312_cov_16.951897_274  GT11  GT11  N  1
NODE_20_length_375312_cov_16.951897_28  GT51  GT51  N  2
NODE_20_length_375312_cov_16.951897_34  CE1  CE1  N  1
NODE_20_length_375312_cov_16.951897_34  CE1  CE1  N  1
NODE_20_length_375312_cov_16.951897_72  GH83  GH83  N  2
NODE_20_length_375312_cov_16.951897_92  CE3  CE3  N  1
NODE_299_length_111202_cov_18.027123_22  GT41  GT41  N  3
NODE_299_length_111202_cov_18.027123_90  GH13  GH13  N  1
NODE_299_length_111202_cov_18.027123_90  GT2  N  2
NODE_299_length_111202_cov_17.246075_43  GT2_Glyco_transf_2(5-163)  N  1
NODE_591_length_62691_cov_17.574552_51  GT9  GT9  N  3
NODE_726_length_50179_cov_16.117540_11  CE4  CE4  N  2
NODE_799_length_18578_cov_18.293481_14  GT9  GT9  N  1
NODE_799_length_18578_cov_16.951897_107  N  GH83  N  1
NODE_207_length_144259_cov_18.029109_41  N  GT0  N  2
NODE_207_length_144259_cov_18.029109_41  N  GT0  N  1
NODE_207_length_144259_cov_18.029109_41  N  GT0  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13_30  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13_11  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1
NODE_25_length_383736_cov_17.808888_176  N  GH13  N  1

```



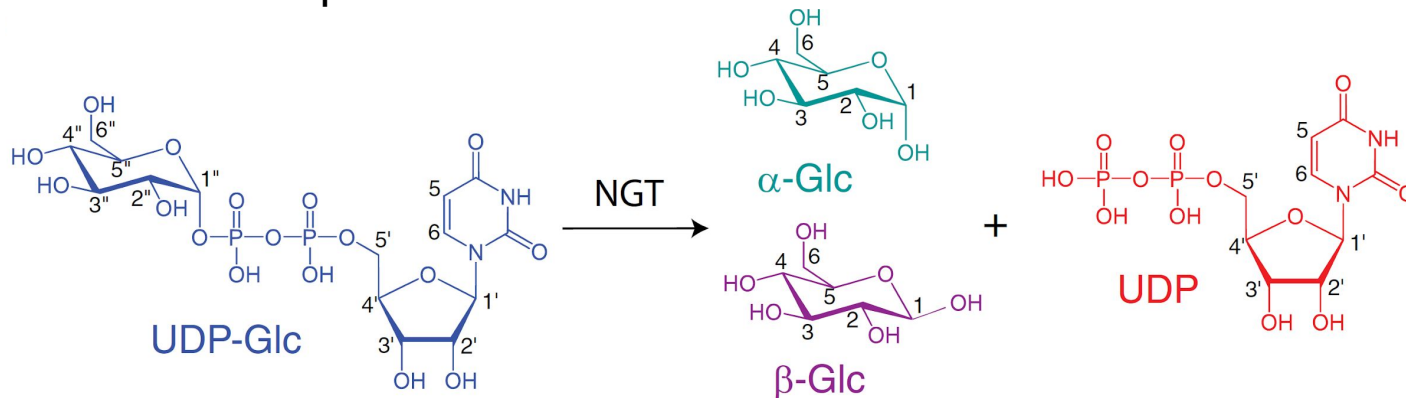
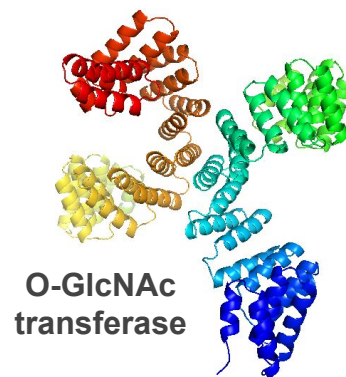
Glucosyltransferase
Glycoside Hydrolase
Carbohydrate Esterase
Carbohydrate Binding Module



GT2	6
GT51	1
GT4	1
GT35	1
CBM48	1
GT5	1
GT83	1
GT30	1
GT19	1
CE11	1
GH84	1
GT41	5
GT28	1
GT9	1
GH13	2
GH77	1

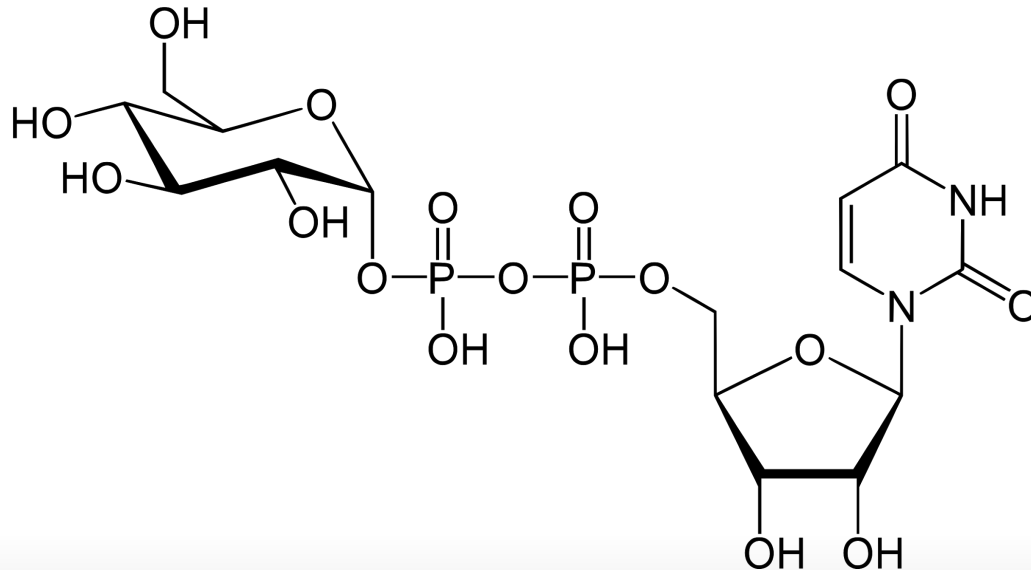
Glycosyl transferase family 41 (GT41) (N-glycosyltransferase)

- Attaches hexoses (saccharide moiety) to nitrogen atom in an amide group (asparagine) on target protein
- Found in bacteria *Actinobacillus pleuropneumoniae*
- Targets adhesion proteins for surface attachment stability
- Can hydrolyze sugar when no substrate present
- Resembles protein O-GlcNAc transferase



Glycosyl Transferase Family 2 (GT2)

- Diverse: > 3500 sequences from animal, plant, fungi, and bacteria, > 12 distinct GT functions have already been characterized
- Cellulose synthase, chitin synthase, mannosyltransferase, glucosyltransferase, galactosyltransferase, rhamnosyltransferase, etc.



UDP-glucose

Other CAzyme genes found in MAG 5

- Glycoside Hydrolase (GH13, GH77, GH84): hydrolysis and rearrangement of glycosidic bonds
- Carbohydrate Esterase(CE11): UDP-3-O-acyl N-acetylglucosamine deacetylase
- Carbohydrate Binding Module (CBM48): glycogen-binding function, appended to GH13 modules

CAZyme Gene Clusters (CGCs)

Show 15 ▾ entries							Search: <input type="text"/>	
Cluster	Contig ID	# of Genes	# of CAZymes	# of TC	# of TF	Cluster Start	Cluster End	
CGC1		12	1	6	0	184995	200444	
CGC10		2	1	1	0	151765	154876	
CGC11		3	1	1	0	12762	15952	
CGC12		6	1	2	0	125774	131698	
CGC13		3	1	1	0	170482	174002	
CGC2		4	1	1	0	295175	298678	
CGC3		2	1	1	0	454119	456221	
CGC4		3	1	2	0	9284	12401	
CGC5		8	3	4	0	65624	82334	
CGC6		3	2	1	0	142971	148385	
CGC7		12	2	6	0	240595	257338	
CGC8		2	1	1	0	1667	3655	
CGC9		3	1	1	0	32391	35584	
Showing 1 to 13 of 13 entries							First	Previous
							1	Next
							Last	



VirSorter 1.0.3

Simon Roux

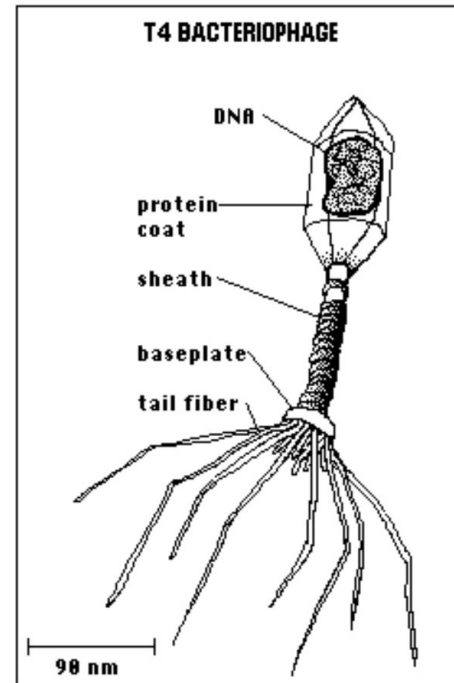
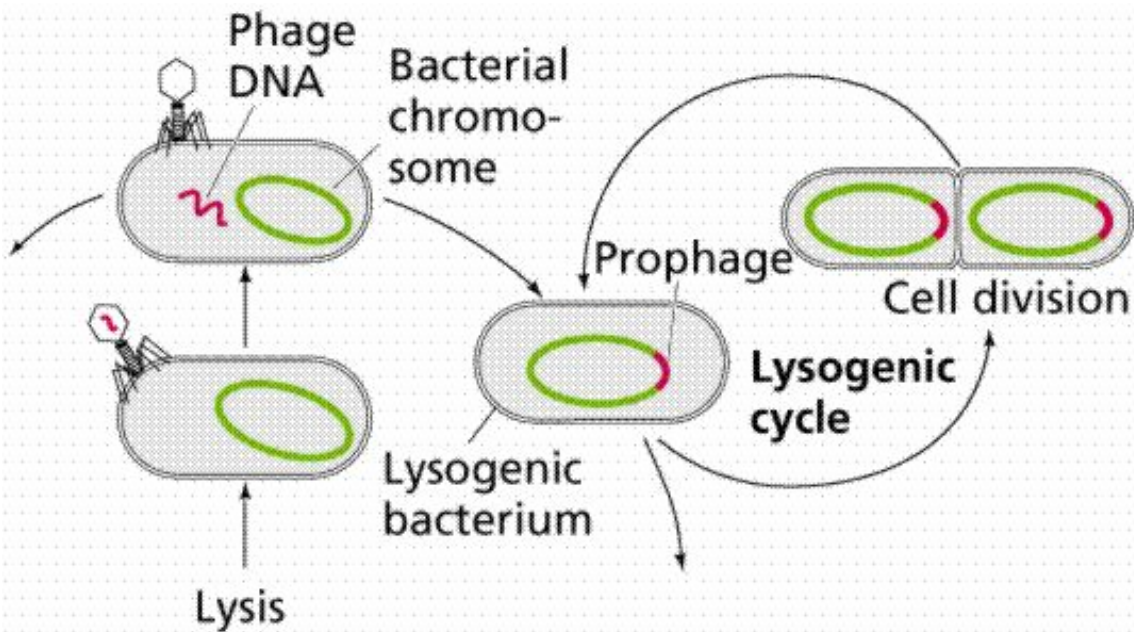
de

★★★★☆ (55)

Viral database used: **Viromes**

All bacterial and archaeal virus genomes in Refseq, plus non-redundant predicted genes from viral metagenomes

- Sequences **entirely viral**, from the more to less confident predictions (category 1, 2, and 3)
- **Prophages** (viral regions detected in a cellular contig), again from the more to less confident predictions (category 4, 5, and 6)



VIRSorter_global-phage-signal.csv - Excel												
File Home Insert Page Layout Formulas Data Review View Help Acrobat Search Share Comments												
E2 Category												
	A	B	C	D	E	F	G	H	I	J	K	L
1	## 1 - Complete phage contigs - category 1 (pretty sure)											
2	## Contig_id	Nb genes contigs	Fragment	Nb genes	Category	Nb phage	Phage ger	Non-Caud Pfam depletion sig	Uncharacterized enrichment sig	Strand swi	Short genes e	
3	## 2 - Complete phage contigs - category 2 (quite sure)											
4	## Contig_id	Nb genes contigs	Fragment	Nb genes	Category	Nb phage	Phage ger	Non-Caud Pfam depletion sig	Uncharacterized enrichment sig	Strand swi	Short genes e	
5	## 3 - Complete phage contigs - category 3 (not so sure)											
6	## Contig_id	Nb genes contigs	Fragment	Nb genes	Category	Nb phage	Phage ger	Non-Caud Pfam depletion sig	Uncharacterized enrichment sig	Strand swi	Short genes e	
7	## 4 - Prophages - category 1 (pretty sure)											
8	## Contig_id	Nb genes contigs	Fragment	Nb genes	Category	Nb phage	Phage ger	Non-Caud Pfam depletion sig	Uncharacterized enrichment sig	Strand swi	Short genes e	
9	## 5 - Prophages - category 2 (quite sure)											
10	## Contig_id	Nb genes contigs	Fragment	Nb genes	Category	Nb phage	Phage ger	Non-Caud Pfam depletion sig	Uncharacterized enrichment sig	Strand swi	Short genes e	
11	## 6 - Prophages - category 3 (not so sure)											
12	## Contig_id	Nb genes contigs	Fragment	Nb genes	Category	Nb phage	Phage ger	Non-Caud Pfam depletion sig	Uncharacterized enrichment sig	Strand swi	Short genes e	
13	VIRSorter_NODE_112_length_196477_cov_17_935031	178	VIRSorter_	100	3				gene_74-gene_173:7	gene_146-gene_156:2.21203832126184		
14												
15												

Prophage ranking of 3, not so sure

Conclusions

- *MAG 5* is a high-quality assembly: **99.4%** completeness and **0%** contamination
- *MAG 5* represents a bacterial genome in the phylum *Firmicutes*
- Closest relative on the phylogenetic tree: ***Selenomonas sputigena***
- No potential for antibiotics production
- **26 CAZyme genes** predicted by *dbCAN2* - sugar metabolism and glycosylation
- **1 not-so-sure prophage** predicted by *Virsorter*

References

1. Grass, S., Buscher, A. Z., Swords, W. E., Apicella, M. A., Barenkamp, S. J., Ozchlewski, N., and Geme, J. W. S. (2003). "The Haemophilus influenzae HMW1 adhesin is glycosylated in a process that requires HMW1C and phosphoglucomutase, an enzyme involved in lipooligosaccharide biosynthesis." *Molecular Microbiology*, 48(3), 737–751.
2. Naegeli, A., Neupert, C., Fan, Y.-Y., Lin, C.-W., Poljak, K., Papini, A. M., Schwarz, F., and Aebi, M. (2014). "Molecular Analysis of an Alternative N-Glycosylation Machinery by Functional Transfer from Actinobacillus pleuropneumoniae to Escherichia coli." *Journal of Biological Chemistry*, 289(4), 2170–2179.
3. "Pfam: Family: Glyco_transf_41 (PF13844)." (n.d.). <http://pfam.xfam.org/family/Glyco_transf_41> (May 29, 2019).