



# Functional Annotation Final Results

## Team 1 Functional Annotation

Wenyi Qiu, Tianze Song, Saurabh Gulati, Ryan Place, Dongjo Ban, Qinwei Zhuang,  
Kunal Agarwal, Frank Ambrosio



# Background/Review

- Antibiotic Resistance
  - Antibiotic Resistance
  - Colistin Resistance
- Next Generation Sequencing
  - Bioinformatics
  - Scalability
- Functional Annotation
  - Multi-Tool Approach
  - Gene Clustering

# Antibiotic Resistance

## High levels of antibiotic resistance found worldwide, new data shows

News release

29 JANUARY 2018 | BANGKOK - WHO's first release of surveillance data on antibiotic resistance reveals high levels of resistance to a number of serious bacterial infections in both high- and low-income countries.

WHO's new Global Antimicrobial Surveillance System (GLASS) reveals widespread occurrence of antibiotic resistance among 500 000 people with suspected bacterial infections across 22 countries.

<http://www.who.int/mediacentre/news/releases/2018/antibiotic-resistance-found/en/>

- Antibiotic resistance is a rapidly growing problem
- Strains have been shown to possess resistance to last line antibiotics such as colistin
- The power of sequencing technologies is ever increasing
- Bioinformatic analysis techniques must scale up

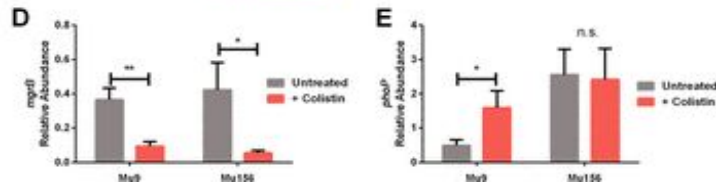
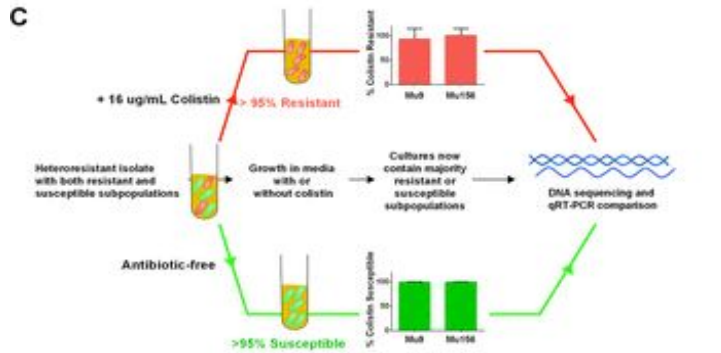
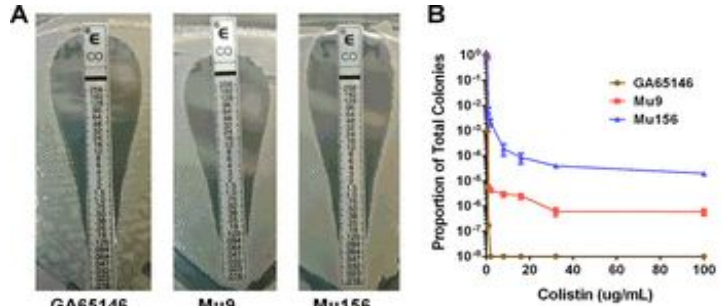
# Colistin Resistance

- Significance
  - Last-Line Drug
  - Resistance Genes are Mobile
- Mechanism
  - Efflux pump
- Related Gene Family
  - MCR-1 to MCR-5 (Lipopolysaccharide modification)
  - PhoP
  - PhoQ

# Heteroresistance

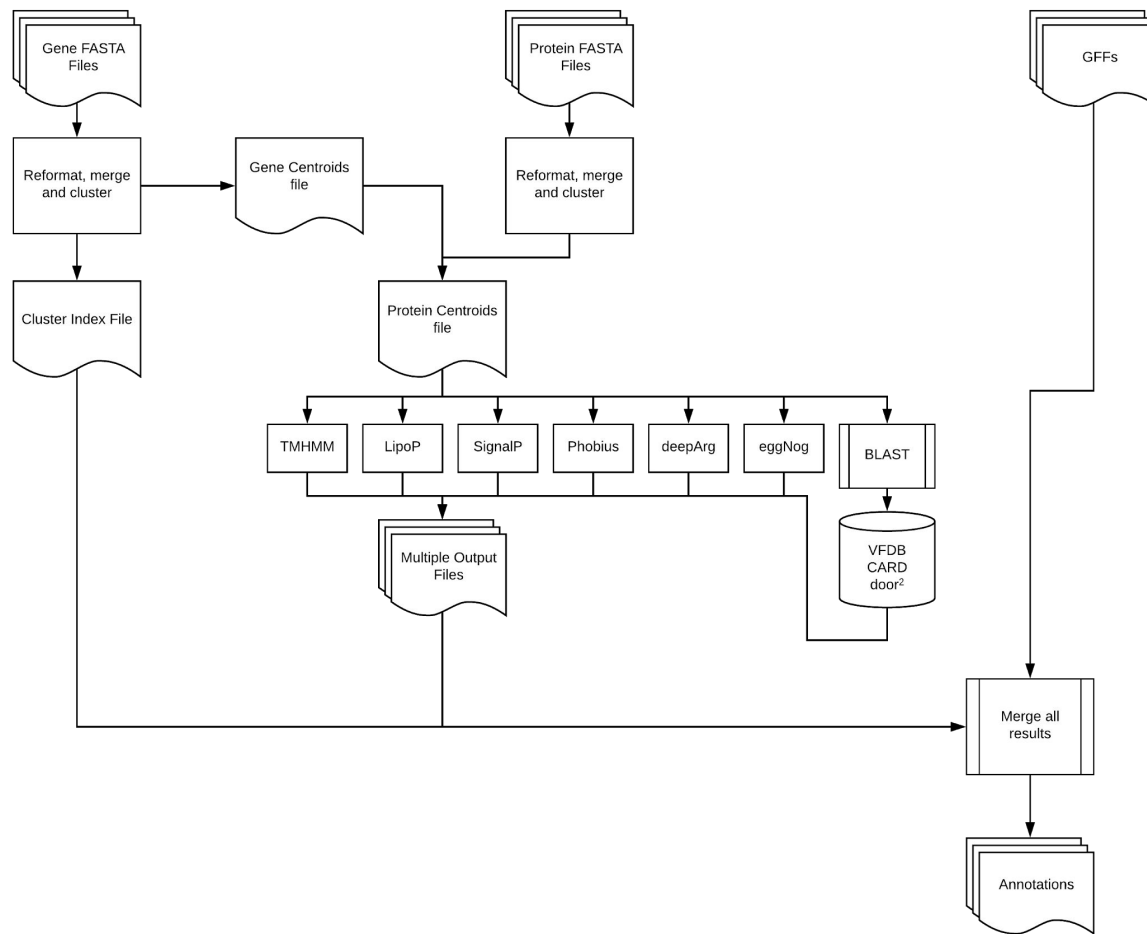
- Heteroresistance is an emerging concern for AMR research
- Isolates previously considered fully susceptible must now be reviewed
- Presence of a resistant subpopulation will affect treatment
- Ability to identify (culture independently) resistant subpopulation will save lives

# Heteroresistance and Heterosusceptibility



- Heteroresistance exists in two “states”
- Untreated Isolate
  - 95% Susceptible
  - 5% Resistant
- Antibiotic Treated Isolate
  - 5% Susceptible
  - 95% Resistant
- Relative abundance of *mgrB(D)* and *phoP(E)* determined by qPCR
- <http://mbio.asm.org/content/9/2/e02448-17.full.pdf+html>

# Pipeline

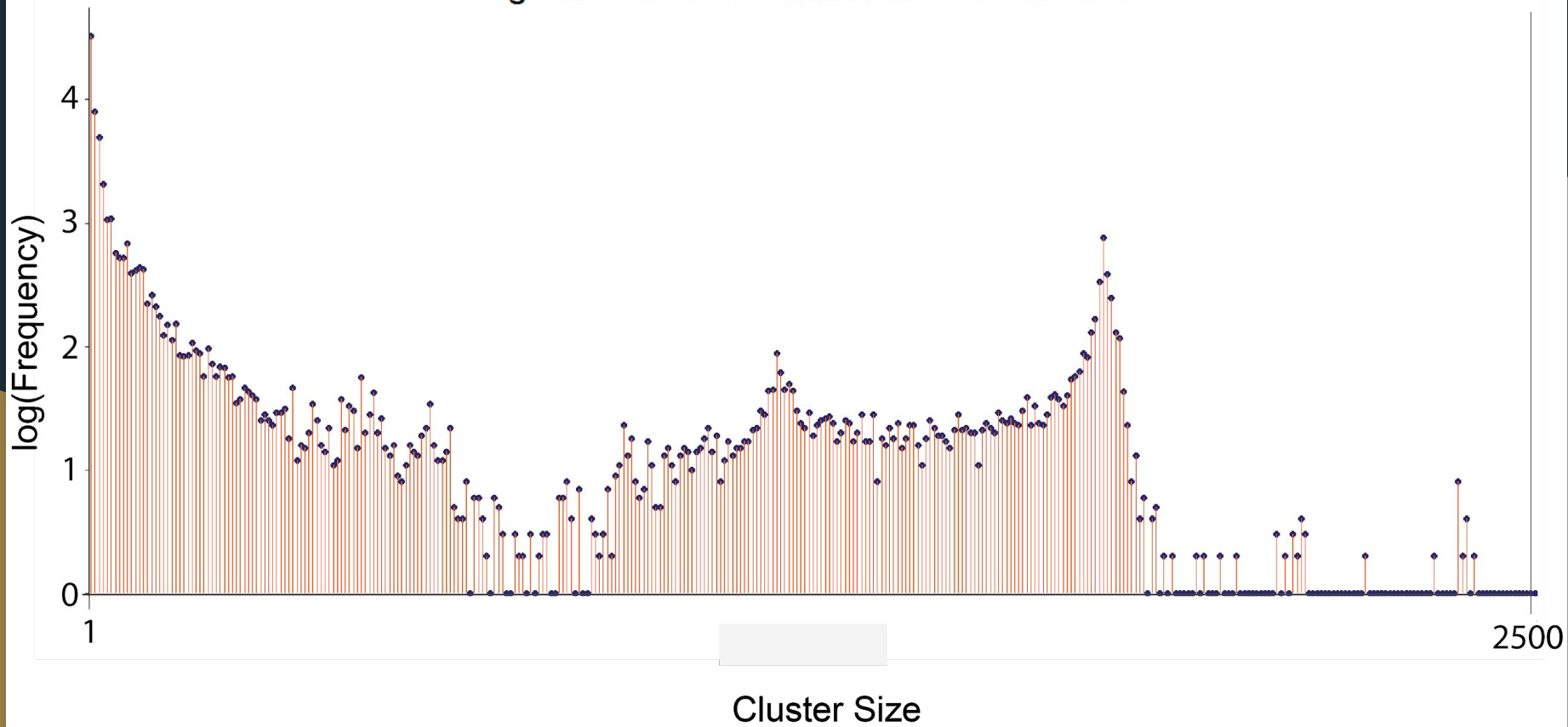


# Clustering

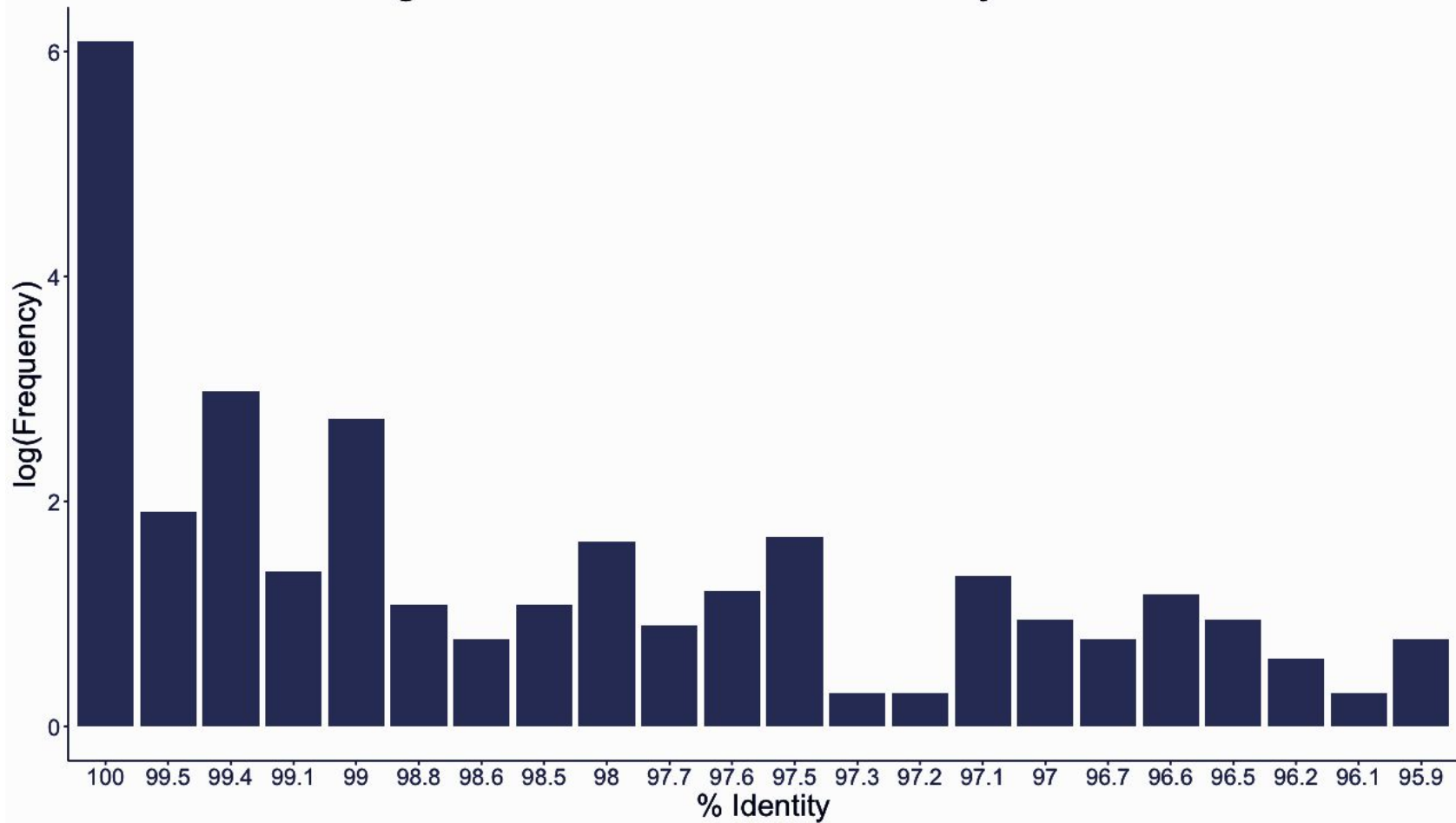
- Tool used: uCLUST (uSEARCH)
- Identity threshold: 0.99
- # of clusters: 63,127
- # of singletons: 32,792
- Max size of clusters: 2,458



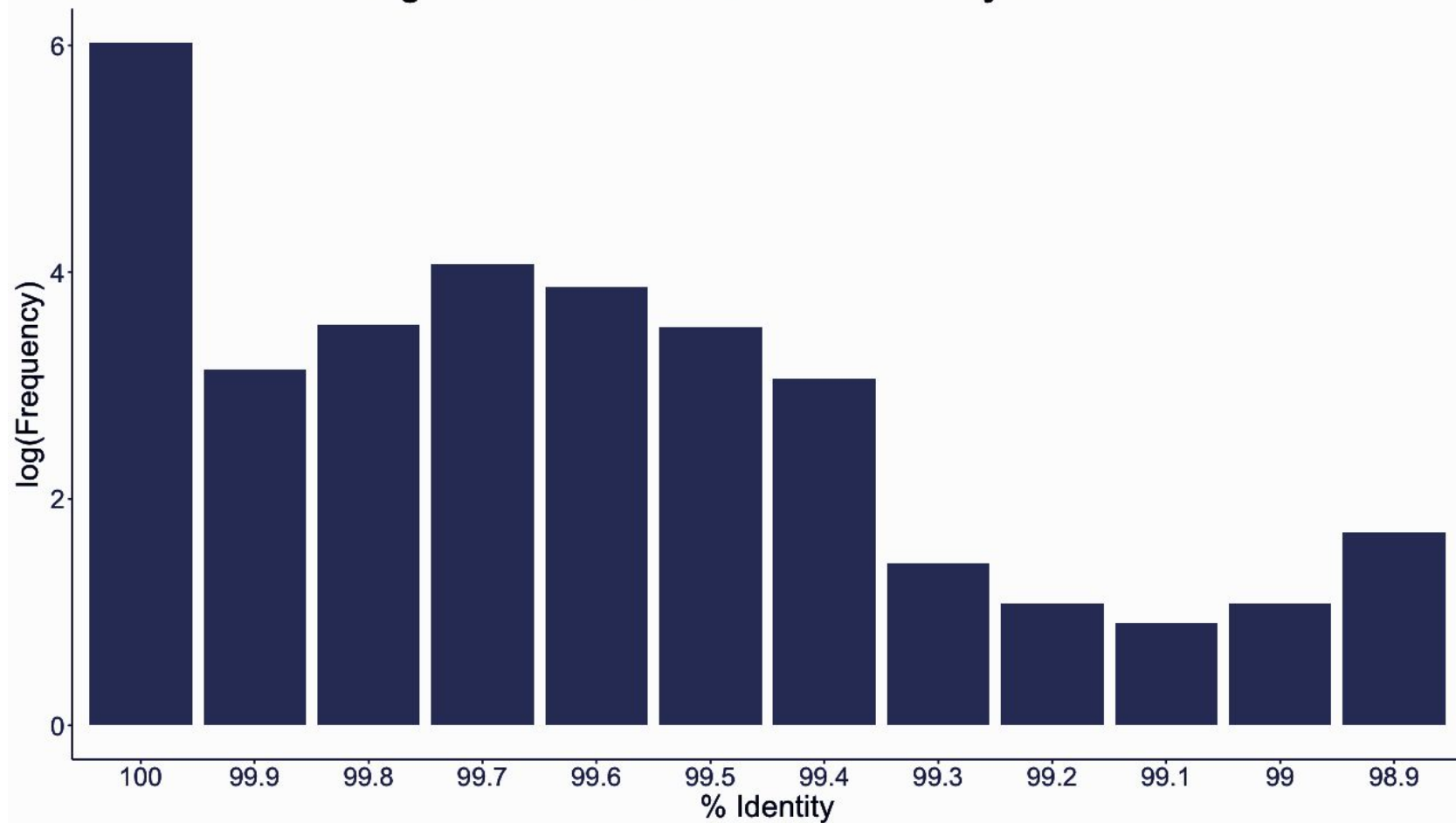
## Log Transformed Distribution of Cluster Sizes



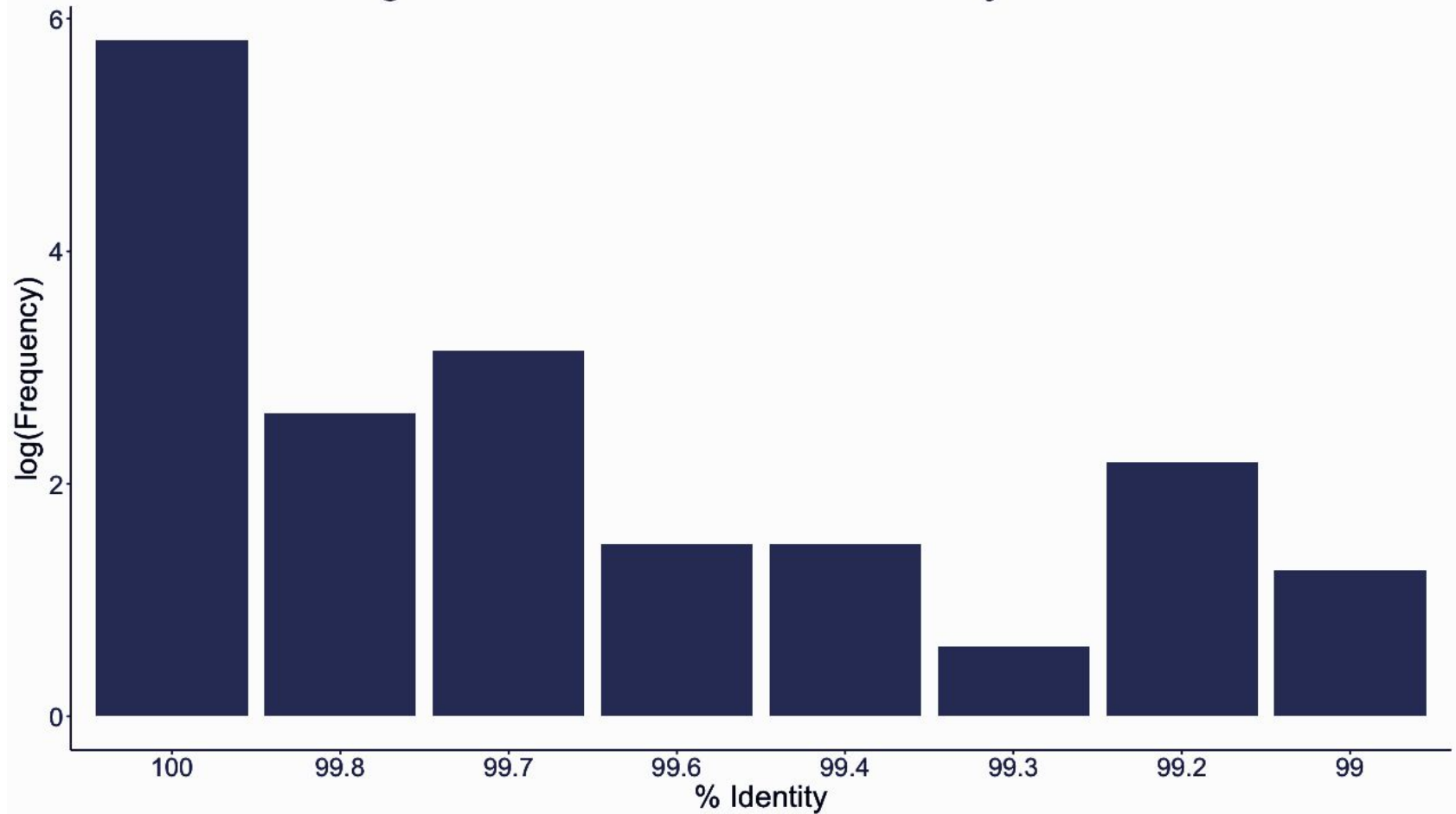
## Log Transformed Distribution % Identity - Cluster 1



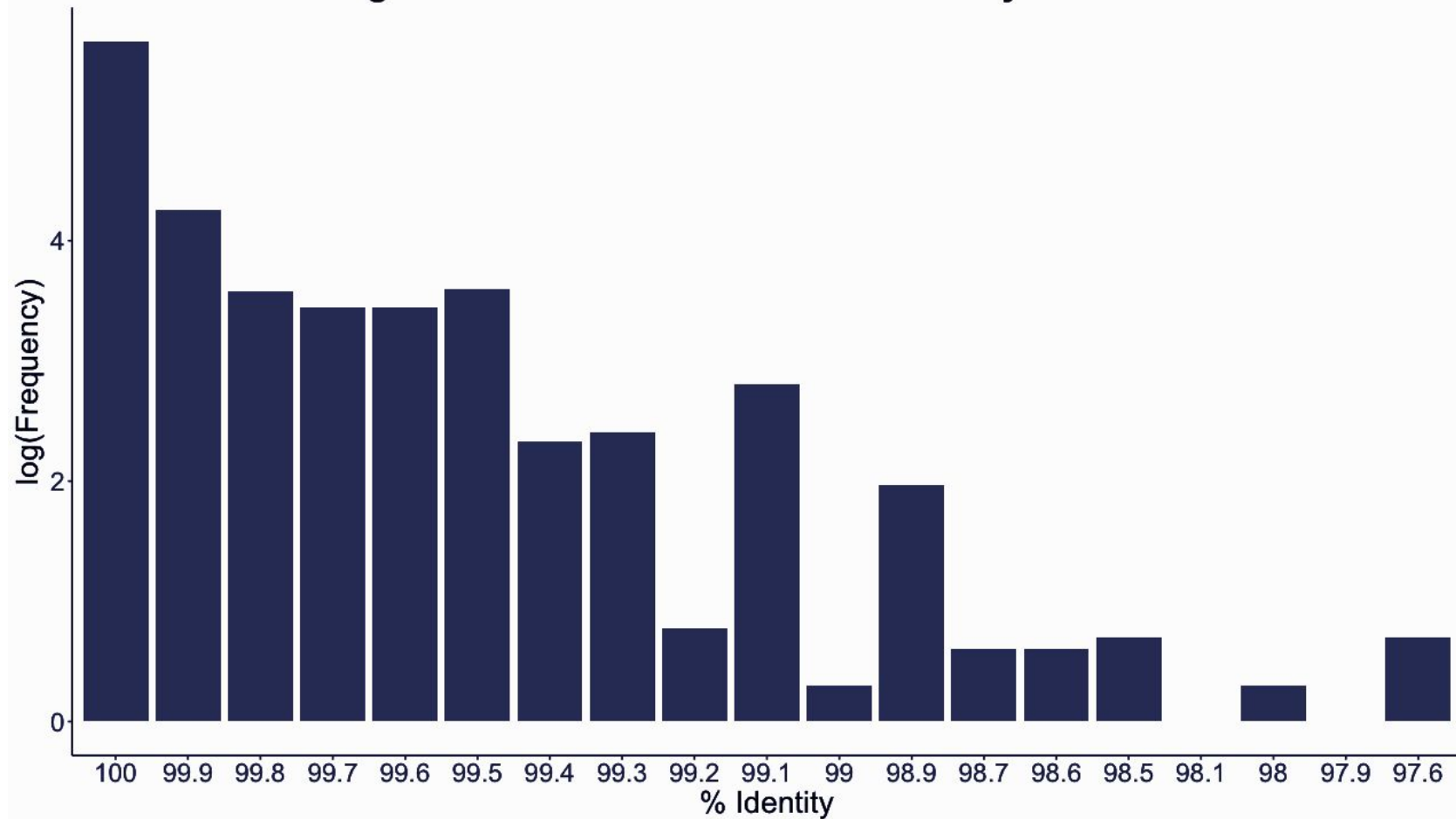
**Log Transformed Distribution % Identity - Cluster 2**



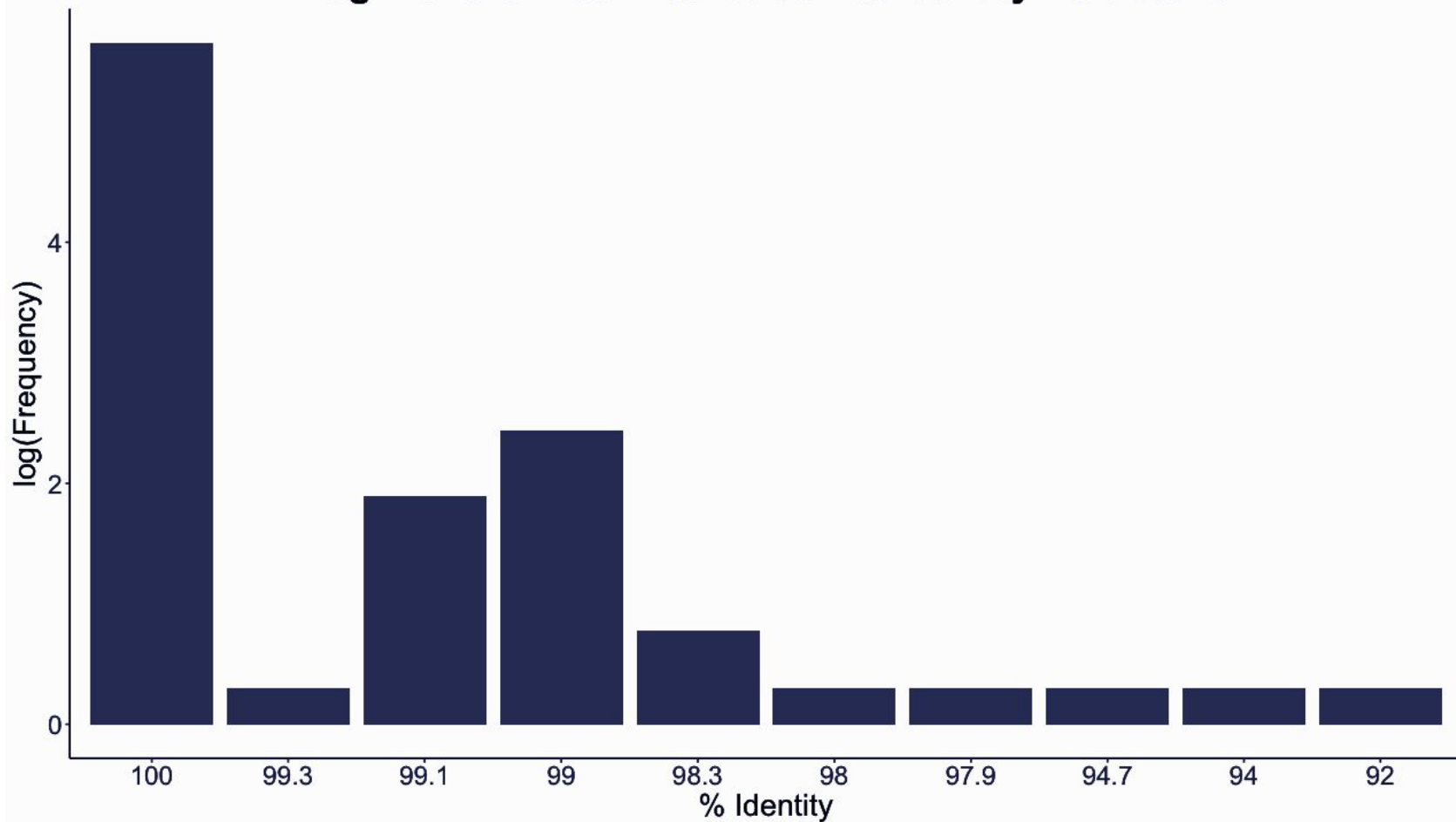
### Log Transformed Distribution % Identity - Cluster 3



## Log Transformed Distribution % Identity - Cluster 4



## Log Transformed Distribution % Identity - Cluster 5



# Annotation

- Egnog
- SignalP
- LipoP
- VFDB
- CARD
- DeepARG
- DOOR<sup>2</sup>

# Homology Based Functional Annotation

Uses databases of genomic features with known function

- Accuracy is dependent on database quality
  - Garbage in garbage out
- Databases for AMR genomic features are added to on a regular basis
- CARD and VFDB are examples of databases of homologous genomic features



# Interproscan

Command:

```
time /projects/data/team1_functionalAnnotation/interproscan-5.28-67.0/interproscan.sh  
-appIPfamA, CDD, HAMAP, PROSITEPATTERNS, PROSITEPROFILES, SFLD, SMART,  
SUPERFAMILY, TIGRFAM, Phobius, TMHMM, SignalP_GRAM_NEGATIVE -i  
<input.faa> -f gff3
```

Number of annotations for clustered centroids: 49831 (total sequences for analysis: 63127)

Running time: ~48h (could have been way faster if using cluster mode instead of standalone mode)

```
real    2895m30.857s  
user    6041m13.085s  
sys     377m0.426s
```

# eggNog



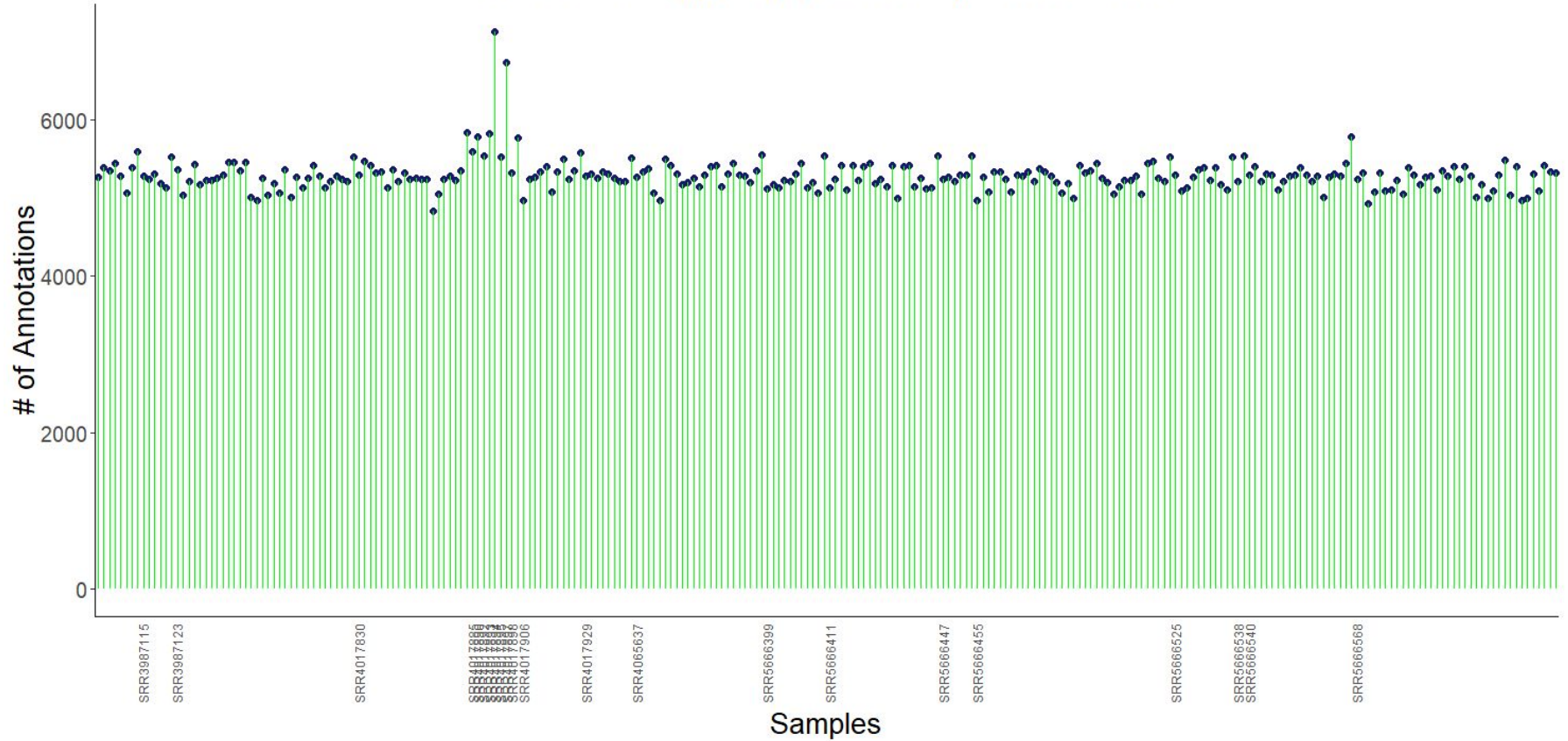
Command:

```
time python emapper.py -i <input.faa> --output test1_faa_maNOG  
-m diamond --translate --usemem --cpu 4
```

Number of annotations for clustered centroids: 49353 (total sequences for analysis: 63127)

Running time: ~2h

# Number of eggNog (Diamond) Annotations



# eggNog vs. InterProScan

- eggNog able to perform just as well in terms of proteome coverage and precision
- 2.5 times faster than InterProScan
- Annotation sources kept up-to-date

**Fast Genome-Wide Functional Annotation through Orthology Assignment by eggNOG-Mapper**

Jaime Huerta-Cepas,<sup>†,1</sup> Kristoffer Forslund,<sup>†,1</sup> Luis Pedro Coelho,<sup>1</sup> Damian Szklarczyk,<sup>2,3</sup> Lars Juhl Jensen,<sup>4</sup> Christian von Mering,<sup>2,3</sup> and Peer Bork<sup>\*,1,5,6,7</sup>

# CARD

Input: Protein sequence collection of all samples (mergedProtein.fasta)

Command used:

```
time blastp -query ../final_fna_faa/mergedProtein.fasta -db
protein_fasta_protein_homolog_model.fasta -outfmt "6 qseqid qstart qend qlen length qcovs pident evalue
stitle" -max_hsps 1 -max_target_seqs 1 -num_threads 5 > CARD_larger_top1_nucl.out

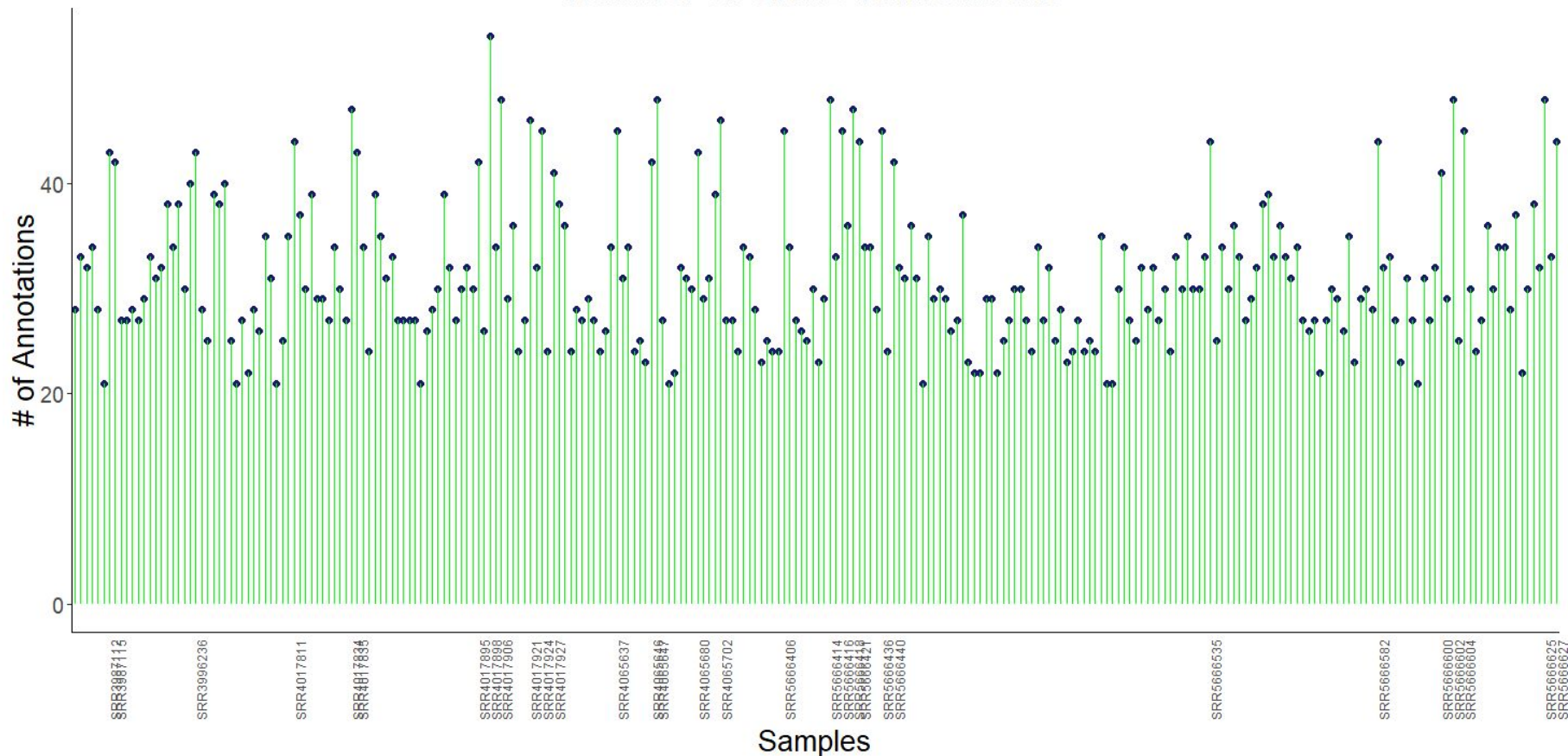
awk '{if($6>=90 && $7>=90) print $0}' CARD_larger_top1_nucl.out >card_90_90.out
```

Size of database: 2239

Number of queries return >90% coverage and >90% identity: 8051 (# of queries: 1,501,569)

Running time: 248 mins

# Number of CARD Annotations



# VFDB

Input: Nucleotide sequence collection of all samples (mergedNucleotide.fasta)

Command used:

```
time blastn -query <input.fna> -db ./VFDB_setB_nt.fas -outfmt "6 qseqid qstart qend qlen length  
qcovs pident evalue stitle" -dust no -max_hsps 1 -max_target_seqs 1 > VFDB_setB_top1_nucl.out
```

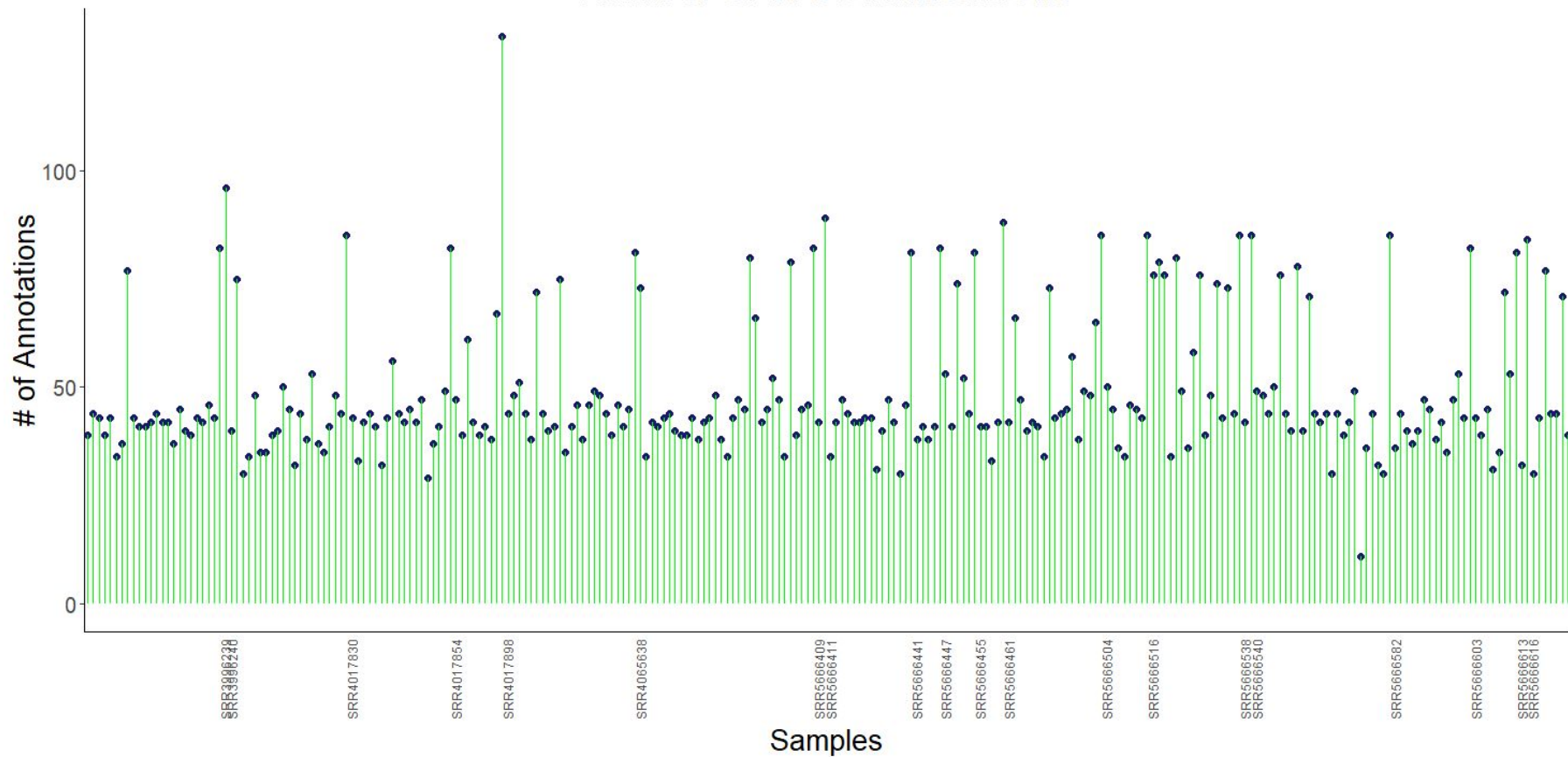
```
awk '{if($6>=90 && $7>=90) print $0}' CARD_larger_top1_nucl.out >card_90_90.out
```

Size of database: 30318

Number of queries return >90% coverage and >90% identity: 12,517(# of queries: 1,501,569)

Running time: ~7 mins

# Number of VFDB Annotations





# DOOR<sup>2</sup>

Input: Protein sequence collection of clustered centroids (centroidsProtein.fasta)

Command used:

```
blastp -db kop_final.table -query ../centroidsProtein.fasta -out $out -max_target_seqs 1 -max_hsp 1  
-num_threads 6 -outfmt "6 qseqid sseqid qstart qend evalue pident qcovs" > operon_intermediate.txt
```

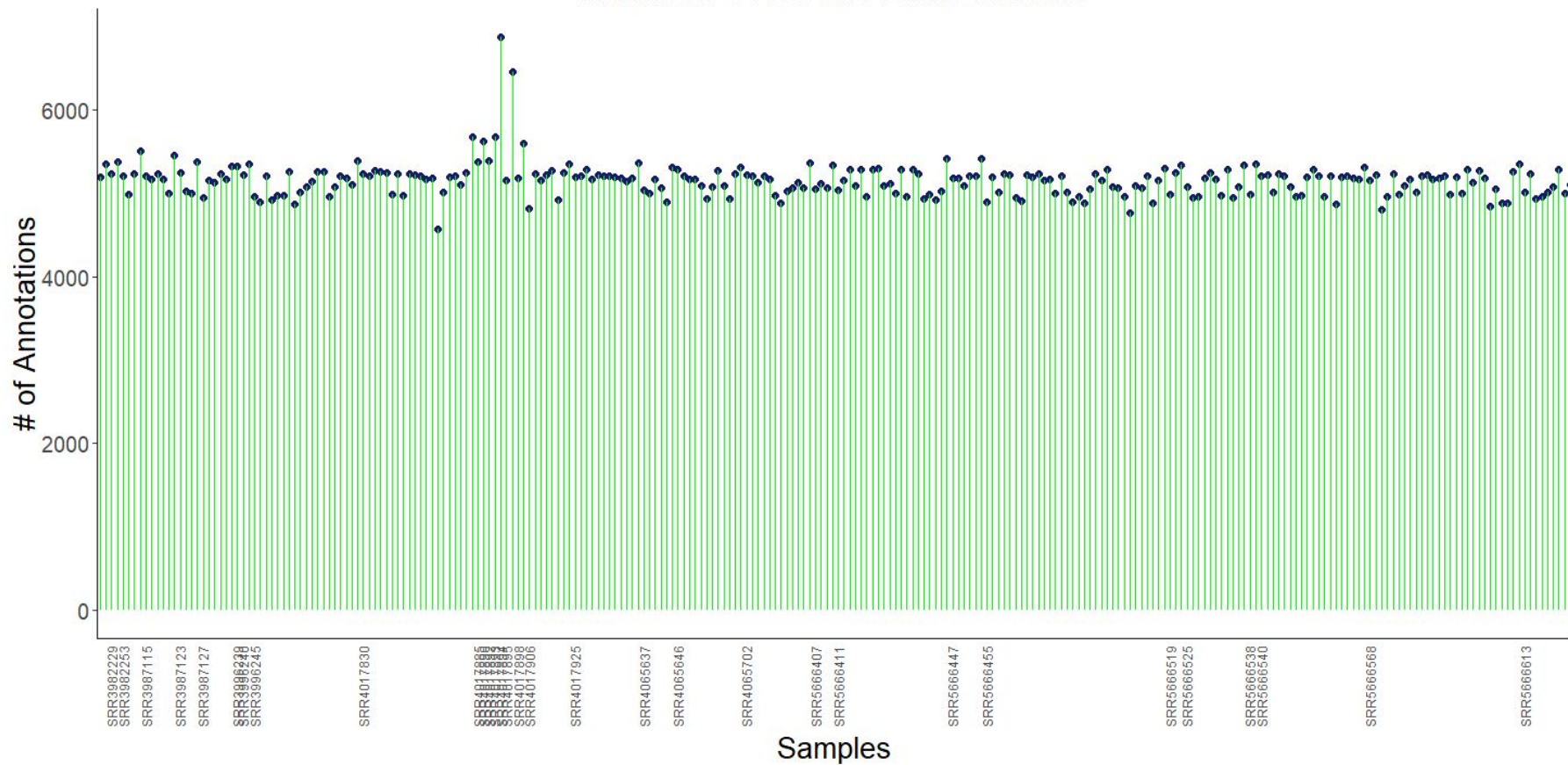
Created a perl script that uses other files to link the blast output sseqid to its GID located in the the DOOR<sup>2</sup> operon tables and adds the operon id, cog, and product to the blast results

Size of database: 97073

Number of queries return  $\geq 90\%$  coverage and  $\geq 90\%$  identity: 38469(# of queries: 63127)

Running time: ~270 minutes

# Number of Door2 Annotations



# DeepARG

Output files: (% identity range: 30% to 100%; e-value < 1e-10)

- 1) ARG: probability  $\geq 0.8$ ; best-hit: found in DeepARG-DB (from CARD, ARDB, and UNIPROT)
- 2) Potential ARG: a) probability  $\geq 0.8$ ; best-hit: "undefined" (one of the 30 categories defined by the deepARG model).  
b) probability < 0.8

Command used:

```
time python ../deeparg-ss/deepARG.py --align --type nucl --genes --input  
mergedNucleotide.fasta --out ./raw_fasta_gene/larged_merged_nucl.out
```

Running time: 44 mins

# DeepARG

Output files: (% identity range: 30% to 100%; e-value < 1e-10)

1) ARG: probability  $\geq 0.8$ ; best-hit: found in DeepARG-DB (from CARD, ARDB, and UNIPROT)

2) Potential ARG: a) probability  $\geq 0.8$ ; best-hit: undefined (one of the 30 categories defined by the deepARG model).  
b) probability < 0.8

% identity distribution	ARG	Potential ARG
[30, 50)	30070 (77.5%)	43960 (94.7%)
[50, 70)	4827 (12.4%)	861 (1.9%)
[70, 90)	3201 (8.3%)	1077 (2.3%)
[90, 100]	701 (1.8%)	538 (1.1%)
Total	38799 (100%)	46436 (100%)

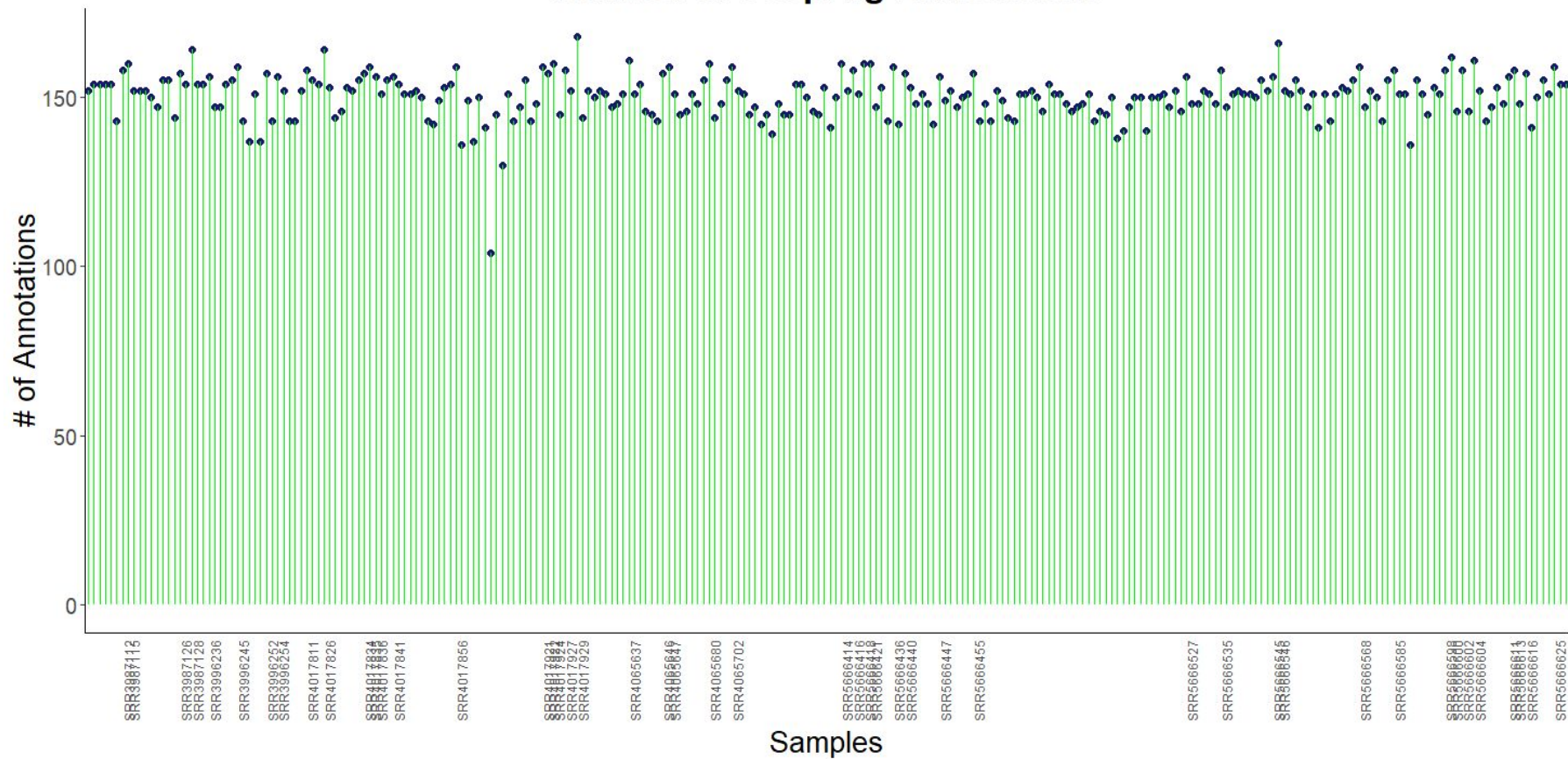
	ARG	Potential ARG
Total ARGs	38799 (100%)	46436 (100%)
Unique ARGs	261	67
Average number	~149	~693

# DeepARG

Interesting results: “fosfomycin” is one of the thirty categories generated by deepARG (already trained).

```
qzhuang8@ubuntu:~/7210/func_annotation/final_fna_faa$ grep 'colistin' ./raw_fasta_gene/larged_merged_nucl.out.mapping.potential.ARG | wc -l
0
qzhuang8@ubuntu:~/7210/func_annotation/final_fna_faa$ grep 'colistin' ./raw_fasta_gene/larged_merged_nucl.out.mapping.potential.ARG | wc -l
0
qzhuang8@ubuntu:~/7210/func_annotation/final_fna_faa$ grep 'fosfomycin' ./raw_fasta_gene/larged_merged_nucl.out.mapping.ARG | wc -l
514
qzhuang8@ubuntu:~/7210/func_annotation/final_fna_faa$ grep 'fosfomycin' ./raw_fasta_gene/larged_merged_nucl.out.mapping.potential.ARG | wc -l
0
```

# Number of DeepArg Annotations



# Ab Initio Functional Annotation

- Looks for intrinsic characteristics of particular gene feature types
- Signal Peptide and Transmembrane Proteins can be identified in this way
  - These regions are of particular importance to this project because of their significance to AMR

# SignalP

Input: Protein sequence collection of clustered centroids  
(centroidsProtein.fasta)

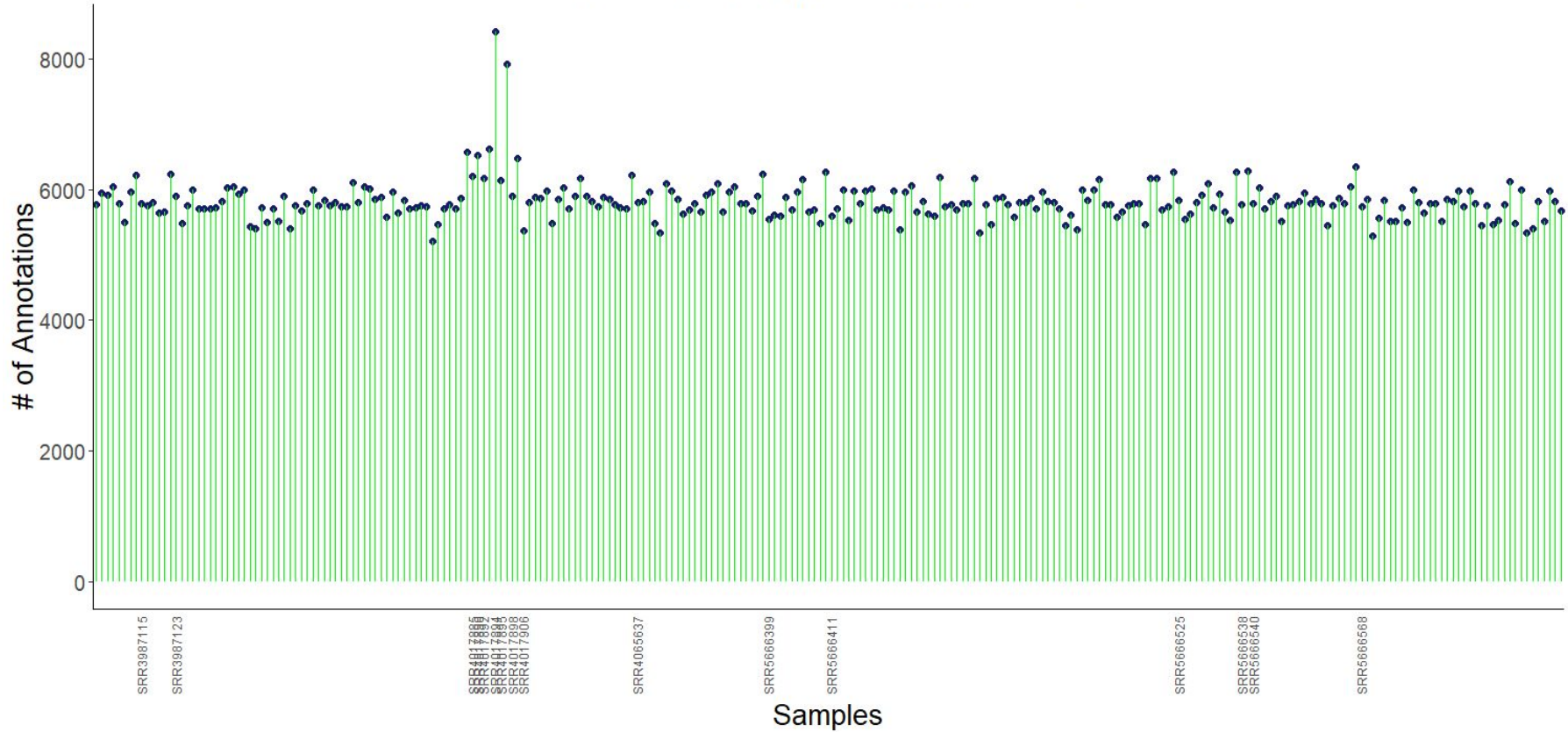
Algorithm: combination of several trained neural networks

Command used:

```
signalp -t <organism_type> -f <output_format> <input_file> >  
<output_file>
```



# Number of SignalP Annotations



# Phobius

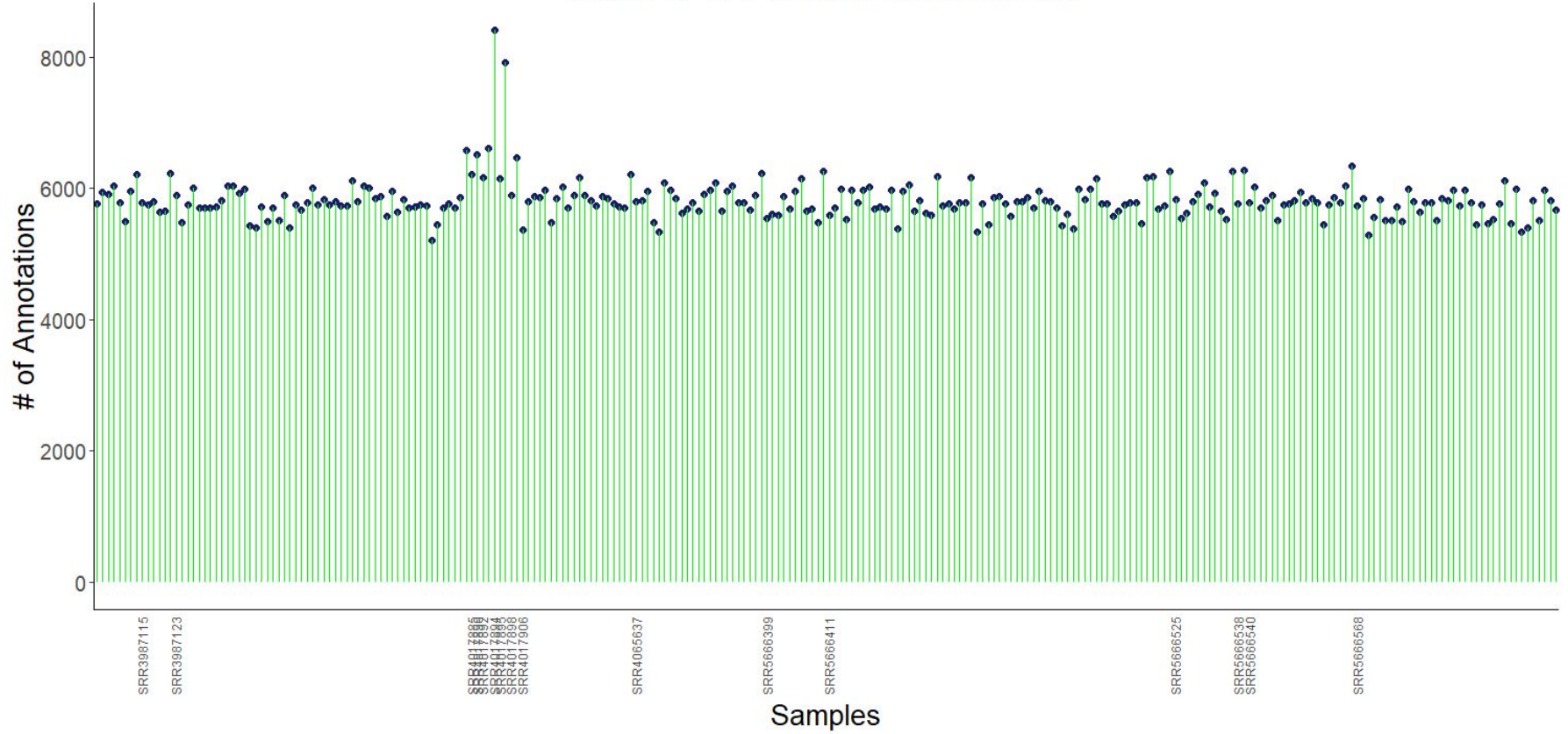
Input: Protein sequence collection of clustered centroids  
(centroidsProtein.fasta)

Algorithm: Hidden Markov Model

Command used:

```
phobius.pl -<output_format> <input_file> > <output_file>
```

# Number of Phobius Annotations



# LipoP

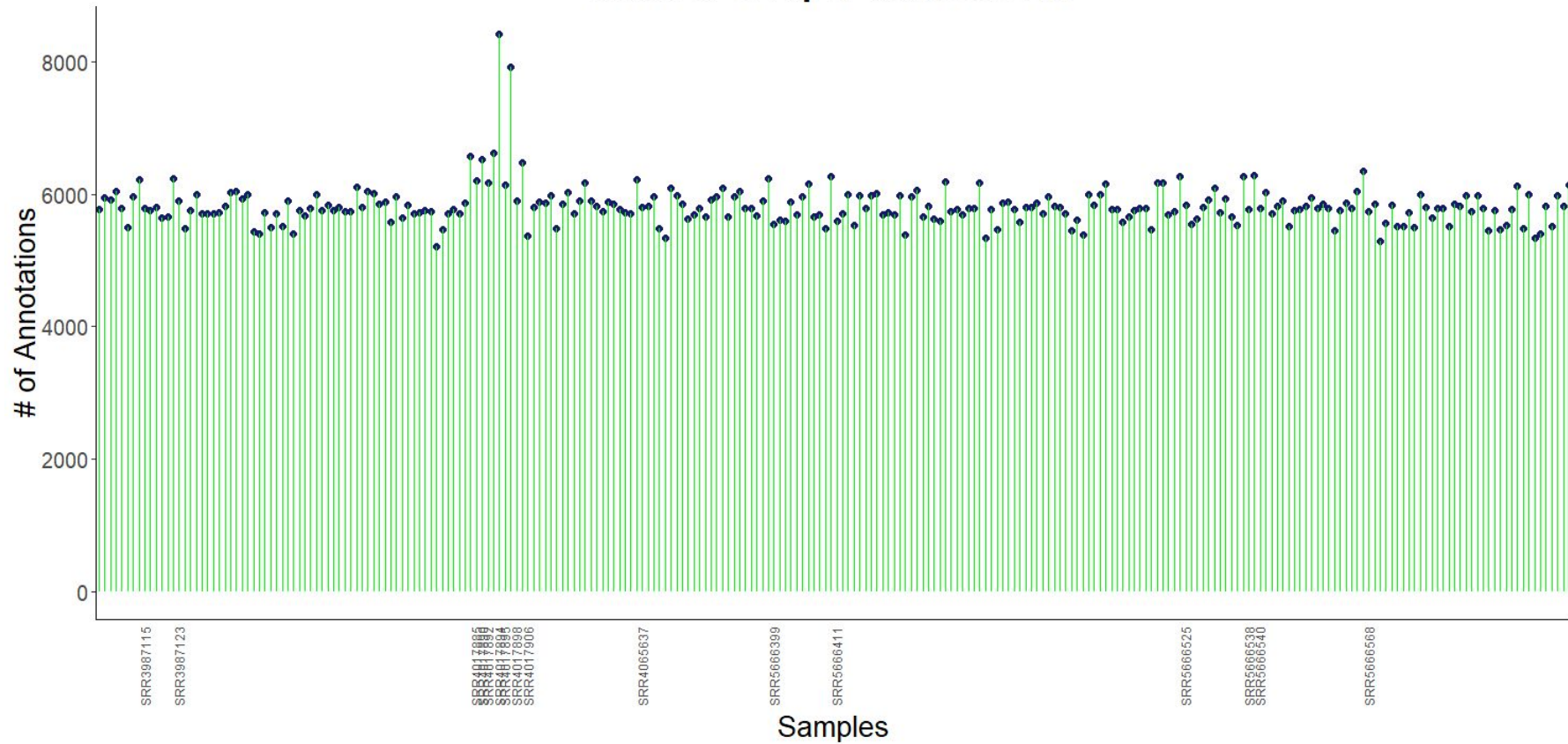
Input: Protein sequence collection of all samples  
(mergedProtein.fasta)

Algorithm: Hidden Markov Model

Command used:

```
LipoP -<output_format> -<input_file> > <output_file>
```

# Number of LipOP Annotations



# TMHMM

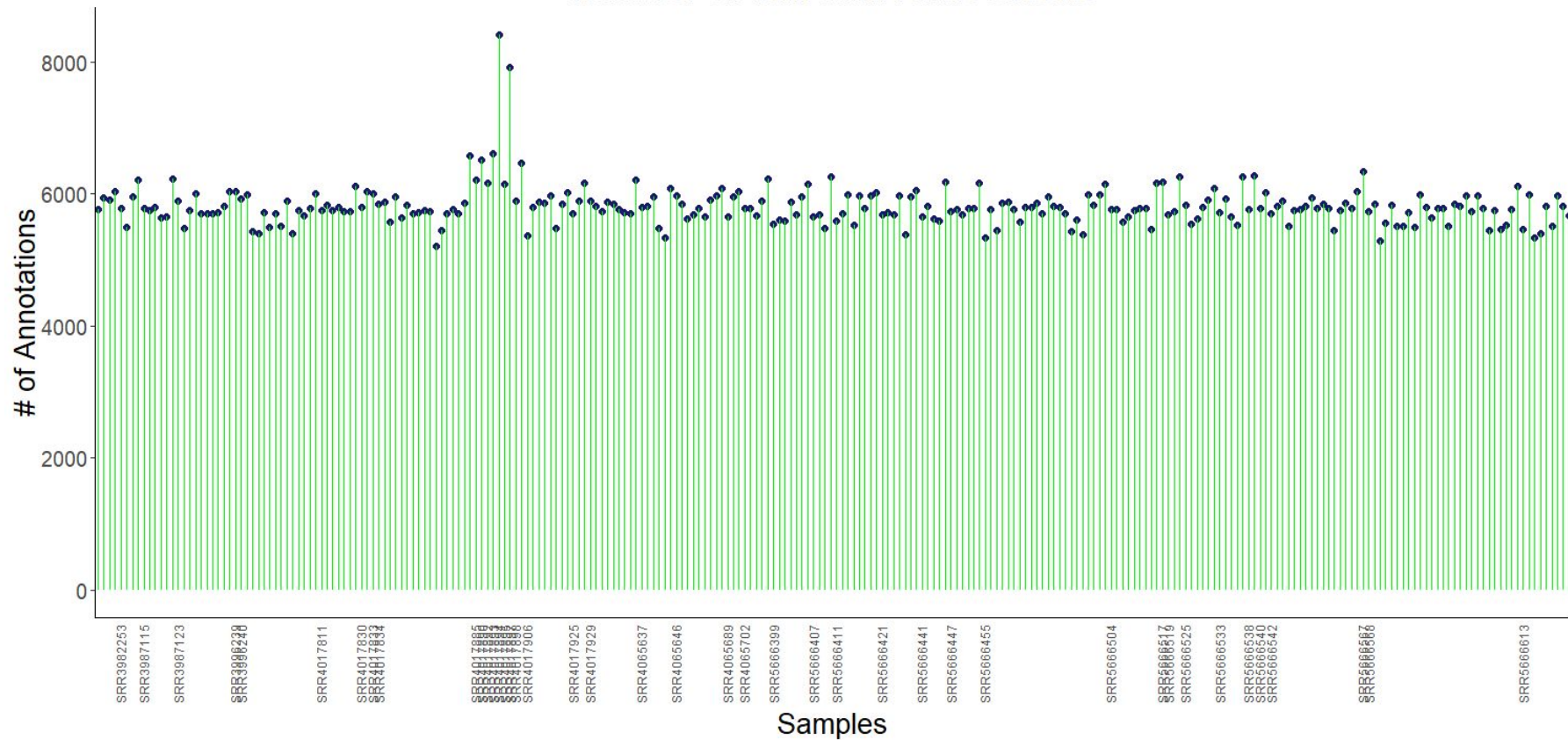
Input: Protein sequence collection of clustered centroids  
(centroidsProtein.fasta)

Algorithm: Hidden Markov Model

Command used:

```
tmhmm -<output_format> <input_file> > <output_file>
```

# Number of TMHMM Annotations



# Overall statistics

<b>Tools / Database</b>	<b>Total # of annotations</b>	<b>Average # of annotations</b>
<b>CARD</b>	<b>8,051</b>	<b>31</b>
<b>DeepARG</b>	<b>38,799</b>	<b>150</b>
<b>Door2</b>	<b>1,330,879</b>	<b>5158</b>
<b>Eggnog (diamond)</b>	<b>1,364,546</b>	<b>5289</b>
<b>GeneMark.hmm</b>	<b>110,235</b>	<b>427</b>
<b>LipoP</b>	<b>1,502,024</b>	<b>5822</b>
<b>Phobius</b>	<b>1,501,560</b>	<b>5820</b>
<b>Prodigal</b>	<b>1,391,789</b>	<b>5395</b>
<b>SignalP</b>	<b>1,501,569</b>	<b>5820</b>
<b>TMHMM</b>	<b>1,501,577</b>	<b>5820</b>
<b>VFDB</b>	<b>12,517</b>	<b>49</b>