

# 1 **Going through phages: A Computational approach to Revealing the role of prophage**

## 2 **in *Staphylococcus aureus***

3 Tyrome Sweet<sup>1a</sup>, Suzanne Sindi<sup>2b</sup>, Mark Sstrom<sup>3a</sup>

4 <sup>a</sup> Department of Life and Environmental Sciences, University of California, Merced,

5 California, USA

6 <sup>b</sup> Department of Applied Mathematics, University of California, Merced, California,

7 USA

## 8 **Abstract**

9 Prophages have important roles in virulence, antibiotic resistance and genome  
10 evolution in *Staphylococcus aureus*. Rapid growth in the number of sequenced *S. aureus*  
11 genomes allows for an investigation of prophage sequences in *S. aureus* at an  
12 unprecedented scale. We developed a computational pipeline to detect and analyze  
13 prophage sequences in nearly 10,011 *S. aureus* genomes, discovering thousands of  
14 putative prophage sequences with genes encoding virulence factors and antibiotic  
15 resistance.

## 16 **Importance**

17 Bacteriophages (phages) play key roles in bacterial evolution, governing abundance,  
18 adaptation and diversity of bacterial communities. Temperate phage can facilitate

19 bacterial adaptation via transduction of novel genes. This study takes advantage of the  
20 unprecedented quantity of genomic sequencing in public repositories to analyze viral  
21 genes in 10,000 *Staphylococcus aureus* genomes. We found 196,727 prophage sequences,  
22 with an estimated total of 129,935 genes. We determined the function of these genes,  
23 identifying a large quantity of novel genes that benefit the host such as beta-lactamase,  
24 enterotoxins and cytotoxins. These results will inform studies of bacterial evolution and  
25 adaptation, by identifying the mechanism of horizontal transfer of genes that confer  
26 adaptive traits to bacteria, especially in the context of antibiotic resistance.

## 27 **Acknowledgments**

- 28 ● Multi-Environment Computer for Exploration and Discovery (MERCED) cluster  
29 at UC Merced, funded by National Science Foundation Grant No. ACI-1429783
- 30 ● National Science Foundation –National Research Traineeship in Intelligent  
31 Adaptive Systems (NRT-IAS) (Award No. 1633722)
- 32 ● National GEM Consortium/Georgia Tech Research Institute - GEM PhD  
33 Engineering and Science Fellowship
- 34 ● Ali Heydari

## 35 **Introduction**

36 Bacteriophages are the most abundant self replicating organisms on earth, with an  
37 estimated global population of  $10^{31}$ , phages outnumber bacteria by 10 to 1<sup>1</sup>(**Liu, 2014**).  
38 Lysogenic phage are transduced into the host bacterial genome as prophage sequences,  
39 and can have a range of selectional impacts on the host, spanning the breadth of the  
40 mutualism-parasitism continuum (**Blair, Webber, Baylay, Ogbolu and Piddock 2015**). It  
41 is hypothesized that prophage sequences that confer a selective advantage to their host  
42 are more likely to be conserved in the bacterial genomes than those that are neutral or  
43 deleterious to their hosts (**Gandon 2016**). The resultant expectation is that prophage  
44 sequences will contain an elevated quantity of genes conferring adaptive functions to  
45 host bacteria.

46 A well-studied example of an adaptive trait conferred by transduction by lysogenic  
47 phage is the *mecA* gene encoded in the phage *Staphylococcus SCIURI7* (**Zeman,**  
48 **Mašlaňová, Indráková, Šiborová, Mikulášek, Bendíčková and Pantůček 2017**).

49 Transduction of this temperate phage into the *Staphylococcus aureus* genome confers  
50 resistance to broad spectrum beta-lactam antibiotics (**Scharn, Tenover and Goering**  
51 **2013**). Methicillin Resistant *Staphylococcus aureus* (MRSA) is one of the major causes of  
52 antibiotic resistant clinical infections. Between 1999 and 2005, hospitalizations for *S.*  
53 *aureus* increased from 294,570 patients to 477, 927. Moreover, MRSA was responsible for

54 127,036 patients in 1999 increasing to 278,203 by 2005 (**Klein, Smith and Laxminarayan**  
55 **2007**).

## 56 **Bacteriophages impact *S. aureus* evolution**

57 Temperate bacteriophages, bacteriophages whose genome is incorporated into the host  
58 bacterium, can switch between the lytic and lysogenic life cycle (**Liu, 2014**). The lytic  
59 cycle destroys the host, but as long as the phage stays lysogenic it provides several  
60 benefits. One benefit is protection from secondary phage attacks from other prophage.  
61 Temperate phages can lose their switching ability if there are mutations in the  
62 attachment sites. Changes to the gene that encode the recombinase responsible for the  
63 excision of phage can result in ‘grounding’ of the phage (**Ramisetty, and Sudhakari**  
64 **2019**). Grounded phages offer the host benefits, without the risk of entering the lytic  
65 cycle.

66 *S. aureus* has a mesh-like cell wall composed of cross-linked polymer peptidoglycans  
67 (PG). Penicillin-binding proteins (PBPs), mediate the final stages of PG synthesis  
68 (**Scheffers and Pinho 2005**). Methicillin is a  $\beta$ -lactam antibiotic that inhibits the  
69 transpeptidation domain of PBPs, which weakens the cell wall (**Fishovitz, Hermoso,**  
70 **Chang and Mobashery 2014**). MRSA produces PBP2A due to the *mecA* gene that  
71 encodes it. Furthermore, this *mecA* gene is transducible by prophage (**Scharn, Tenover**  
72 **and Goering 2013**).

### 73 **Computational advances for Whole Genome Sequence (WGS) analysis**

74 The number of sequenced and annotated phage genomes is relatively small with 40,981  
75 phage sequences, and 266,129 prokaryotic genome sequences (*Staphylococcus Aureus*  
76 **(ID 154) - Genome - NCBI, n.d.**) on August 18th, 2018. A significant proportion of the  
77 genes encoded by both free living and prophage sequences are of unknown function  
78 (Touchon et al., 2016)**Moura, Criscuolo, Pouseele, Maury, Leclercq, Tarr, Brisse 2016**).  
79 There is a large possibility for novel functions to be conferred to bacterial hosts by  
80 transduction by lysogenic phage (Scharn, Tenover and Goering 2013). Given the  
81 exponential increase in the number of genome sequences deposited in public  
82 repositories, it is timely to take advantage of these sequences to analyze them for novel  
83 functions. In this study we analyze 10,011 *S. aureus* genomes downloaded from NCBI in  
84 2018 for prophage sequences, and determine their functions. The total number of  
85 genome sequences for all organisms numbered 528,859 for 1 online repository(**Genome**  
86 **List - Genome - NCBI, n.d.**). Advances in computational techniques for the analysis of  
87 large data sets have advanced the omics field by enabling researchers to analyze larger  
88 datasets at lower costs (**Krassowski, Das, Sahu, and Misra 2020**).  
89 In this study, we developed a computational pipeline to detect and analyze prophage  
90 sequences in nearly 10,000 *S. aureus* genomes. We discovered thousands of putative  
91 prophage sequences with genes encoding virulence factors and antibiotic resistance. In  
92 particular, we found genes encoding *mecA*, genes encoding toxins/antitoxins and

93 clusters of prophage sequences that had genes in common . Our results, and methods  
94 developed, will facilitate similar studies for other bacterial species and promise to be a  
95 useful tool in the study of prophage host evolution.

## 96 **Materials and Methods**

### 97 ***S. aureus* Genomes**

98 *S. aureus* genomes were obtained from the National Center for Biotechnology  
99 Information NCBI's Genbank repository on August 18, 2018 (***Staphylococcus Aureus (ID***  
100 ***154) - Genome - NCBI, n.d.***). All available genome sequences (n=10,011 including  
101 complete and partial assemblies) were downloaded for this study. (Accession numbers  
102 are provided in Supplementary Tables 1&2).

### 103 **Viral Detection**

104 Putative prophage sequences were detected using PhiSpy, Version 3.2 (**Akhter, Aziz**  
105 **and Edwards 2012**). PhiSpy uses a random forest algorithm that has been trained on  
106 seven distinct features of prophage: protein length, transcription strand directionality,  
107 AT and GC skew, the abundance of unique phage words, phage insertion points and  
108 the similarity of phage proteins. PhiSpy has 49 available training sets to increase  
109 accuracy for specific genomes. We used the *S. aureus* training dataset (option 24) and  
110 identified 196,727 phage regions in our 10,000 *S. aureus* genomes.

## 111 **Prophage Clustering**

112 Prophage sequences identified by PhiSpy were unique within a genome, but highly  
113 redundant between genomes. We identified highly-similar prophages between genomes  
114 through a reciprocal blast (**Johnson, Zaretskaya, Raytselis, Merezhuk, McGinnis and**  
115 **Madden 2008**) search. We increased the max\_target\_seqs to 12,000 (higher than our total  
116 number of *S. aureus* genomes) to ensure we captured all possible matches. We also used  
117 a custom output format which provided additional information on the alignment.  
118 We then grouped prophages by using an undirected graph approach with nodes of the  
119 form: Genome *i*, Prophage *j*. Edges were added between nodes if they had a blast  
120 alignment which exceeded 90% similarity and 90% coverage of both source and target.  
121 We then identified genomes sharing the same prophage by determining the connected  
122 components, resulting in 191 unique phage clusters.

## 123 **Cluster Validation**

124 Each of the 191 phage clusters were aligned with Muscle v3.8.1551 (**Edgar, 2004**) and  
125 ClustalW v2.1 (**Thompson et al., 1994**) to ensure each phage was similar. A score of  
126 0.0000 indicates that the undirected graph script formed accurate phage clusters.

## 127 **Genome Annotation**

128 One representative was selected from each of the 191 phage clusters and analyzed with  
129 2 different tools for gene annotation: VGAS (**Zhang et al., 2019**), and Prokka (**Seemann  
130 2014**). VGAS and PROKKA identified ORFs in each of the phage genome sequences.  
131 VGAS identifies ORFs through an enhanced version of the ZCurve algorithm (**Guo et  
132 al., 2003**) that was customized by adding 13 additional identifying variables (45 total)  
133 for the classification model, and BLASTP (**Mahram & Herbordt, 2015**) searches for gene  
134 prediction. The all ORFs were annotated by all both tools with default settings.

### 135 **Heat Maps**

136 We identified shared genes between phage through a reciprocal blast search using the  
137 annotated phage sequences. We constructed a new undirected graph with the nodes  
138 being the phage genome and the edges representing genes shared between the phage.  
139 The output was a .csv file that listed each of the 191 phage with the genes shared with  
140 other phages. A pairwise count matrix was formed with the x and y axis' are the 191  
141 phage, and the cells show the number of genes shared. Using this matrix, a heatmap is  
142 created using the seaborn python package (**Bisong 2019**). (See **Figures 4 and 7 in the  
143 results section**)



## 144 **Network Visualization**

145 We used the `layout_with_mds` option for the layout function of the R package Igraph  
146 (Csardi and Nepusz 2006) to visualize the phages with shared genes using the pairwise  
147 count matrix for both PROKKA and VGAS.

## 148 **Jaccard Index**

149 The Jaccard Index (Hamers 1989) was calculated using a modified version of the  
150 Jaccard index function in R (Woudenberg 2021) to compare the Prokka and VGAS  
151 networks.

## 152 **Results**

153 Of the 10,011 genomes initially analyzed, there were 12 that were not annotated  
154 completely and did not pass the conversion to SEED (Aziz, Devoid, Disz, Edwards,  
155 Henry, Olsen and Xia 2012) due to missing locus tags (*Prokaryotic Genome*  
156 *Annotation Guide, n.d.*). A further four genome sequences were too short, and did not  
157 have enough base pairs for PhiSpy to detect phage regions, resulting in a total of 9,995  
158 genomes which passed initial QC and were used for subsequent analysis. Within these,  
159 we detected a total of 196,727 prophage sequences across the 10,011 genomes, with an  
160 average of 39 (standard deviation = 1967.468) prophage sequences per genome (**Figure**  
161 **2**).

## 162 **Analysis Uncovers 191 Unique Prophage Sequences**

163 Reciprocal BLAST analysis coupled with undirected graph analysis (see Methods)  
164 found that the 196,727 prophage sequences corresponded to 191 unique prophage  
165 sequences. Each unique prophage sequence appeared in an average of 1024 host  
166 genomes (Standard deviation = 2581.33)(**Figure 3**). Each prophage contained an average  
167 of 16.83 putative coding regions, resulting in a total of 3,208 (VGAS) and 3,203 (Prokka)  
168 unique open reading frames (ORFs). One phage appeared in all 9,995 genome  
169 sequences, while 42 different phages were found in only a single genome sequence.

## 170 **Analysis Detects Thousands of ORFs with Potential Gene Function**

171 One representative was selected from each of the 191 phage clusters and analyzed with  
172 two different tools for gene annotation: VGAS (**Zhang et al., 2019**), and Prokka  
173 (**Seemann, 2014**). VGAS identified 3,208 genes in the 191 unique phage, and PROKKA  
174 detected 3,203 genes. For the PROKKA results, 1,155 ORFs did not have an identified  
175 function. 807 predicted ORFs corresponded to known ORFs with accession numbers  
176 matching known databases ISfinder (**Siguier, Pérochon, Lestrade, Mahillon and**  
177 **Chandler 2006**), NCBI (*National Database of Antibiotic Resistant Organisms*  
178 *(NDARO) - Pathogen Detection - NCBI, n.d.*), UniProtKB (**Boutet, Lieberherr,**  
179 **Tognolli, Schneider and Bairoch 2007**). 2041 genes had a predicted gene function.  
180 VGAS predicted 2935 ORFs, 362 of which corresponded to known accession numbers

181 matching databases Swissprot and refseq (**Pruitt and Maglott 2001**) and 309 other  
182 predicted ORFs had predicted gene functions.

### 183 **Analysis Shows Shared ORFs between Unique Prophage Sequences**

184 2 undirected graphs based on the genes identified by PROKKA and VGAS were  
185 created. The approach outlined in the “**prophage clustering**” section with nodes of the  
186 form: Genome *i*, identified gene *j*. Edges were added between nodes if they had a  
187 matching identified gene. We then Compared the edges produced by both tools  
188 PROKKA and VGAS with each other.

189 We found a total of 1,335 shared edges defined by PROKKA and VGAS. The lowest  
190 number of shared edges between phage sequences was 1, and the highest was 73  
191 (**Figure 4A**). There were 1,306 shared edges between PROKKA and VGAS, and 28  
192 shared edges unique to PROKKA (**Figure 6**) out of the total 1,335 (**Figure 5**). In the 28  
193 unique PROKKA the numbers of shared edges between each node ranged from 1 to 22.  
194 VGAS defined a total of 1,334 connected components. The lowest amount of genes  
195 shared between phage sequences was 1, and the highest was 75 (**Figure 4**). There were  
196 27 shared edges unique to VGAS (**Figure 9**) out of the total 1334 (**Figure 8**). The 27  
197 unique VGAS shared edges ranged from 1 to 22 as well.

### 198 **Discussion**

199 There were several virulence factors and toxins identified in the 191 unique prophage  
200 representatives, 1% of the total 196,727 phage detected. *Staphylococcus aureus subsp.*  
201 *aureus strain NCTC 8325* is referenced several times throughout the dataset. It was used  
202 as a propagating strain for bacteriophage 47 of the international typing set of  
203 bacteriophages and is considered the prototypical strain for most genetic research on *S.*  
204 *aureus* (*Staphylococcus Aureus Subsp. Aureus NCTC 8325 (ID 57795) - BioProject -*  
205 *NCBI, n.d.*).

#### 206 **Genes Encoding *mecA* Found in 2 of the 191 Unique Prophage**

207 There were several traces of antimicrobial resistance found in the 191 phage clusters.  
208 The *mecA* ancestral gene specifically was identified in 2 sequences. The first sequence,  
209 accession number ASM900v1(*Staphylococcus Aureus (ID 154) - Genome - NCBI, n.d.*),  
210 cluster group has 1023 phage, 10% of the total *S. aureus* genomes. ASM900v1, or RF122  
211 (ET3-1) provides a framework for the identification of specific factors associated with  
212 host specificity in this major human and animal pathogen (**Herron et al., 2002**). RF122  
213 (ET3-1) has several genes involved in host colonization, toxin production, iron  
214 metabolism, antibiotic resistance, and gene regulation (**Herron-Olson et al, 2007**).  
215 ASM323779v1 (*Staphylococcus Aureus (ID 400143) - BioProject - NCBI, n.d.*) is the only  
216 phage in the cluster, making it individually unique compared to the 196,727 total

217 detected. It is a part of 184 *S. aureus* isolates collected from 135 patients over a  
218 timespan of 3 years at an Italian paediatric hospital (**Manara et al., 2018**).

#### 219 **48 Unique Gene Functions Supporting Antimicrobial Resistance**

220 48 unique encoding traces of Antimicrobial Resistance (**Table 1**). 4 genes stuck out the  
221 most due to the number of clusters they appeared in. GDAEFEPF\_00005 Staphylococcal  
222 complement inhibitor, a gene found in ASM2514v1 appeared in 10 clusters (**Nübel et**  
223 **al., 2010**). GHDFECEE\_00007 Superantigen-like protein 13 was found in ASM17451v1  
224 and appeared in 8 clusters (**Cameron et al., 2012**). ASM17451 also contained  
225 GHDFECEE\_00008 Superantigen-like protein 13 which appeared in 7 clusters.  
226 GAIDFPLK\_00004 Superantigen-like protein 13 was found in ASM1150v1 and was  
227 identified in 7 clusters (**Holden et al., 2004**).

#### 228 **4 Genes Showing Traces of Toxin/Antitoxin (TA) System**

229 Toxin/Antitoxin (TA) systems encode toxin proteins that interfere with vital cellular  
230 functions and are counteracted by antitoxins. There are 6 different types of TA systems,  
231 *S. aureus* has genes identified showing types I, II and III (**Schuster & Bertram, 2016**).  
232 Type I toxin-antitoxin systems have the base-pairing of antitoxin RNA with the toxin  
233 mRNA (**Fozo et al., 2008**). Type III systems toxic proteins and an RNA antitoxin have a

234 direct iteration where the toxic proteins are neutralized by the RNA gene (**Labrie et al.,**  
235 **2010**).

236 Type II, the most studied TA system, has proteic antitoxin that tightly binds and inhibits  
237 the activity of a stable toxin (**Hayes, 2003**). The TA system yoeB-yefM has been detected  
238 as genes MBJHDCJA\_00021 Toxin YoeB and MBJHDCJA\_00022 Antitoxin YefM in  
239 ASM900v1 (**Herron et al., 2002; Herron-Olson et al., 2007**). yoeB inhibits bacterial  
240 growth and translation by cleavage of mRNA molecules and is repressed by antitoxin  
241 yefM (**Schuster & Bertram, 2016**). Enterotoxin Type A causes food poisoning and was  
242 identified in 3 genome sequences (**Ono et al., 2012**). M1022 (NCTC 8325) was identified  
243 in 2 genome sequences (Staphylococcus Aureus Subsp. Aureus NCTC 8325 (ID 57795) -  
244 BioProject - NCBI, n.d.). CAFLMJIC\_00063 Enterotoxin type A was identified in 1  
245 genome sequence (**Herron et al., 2002; Herron-Olson et al., 2007**). (**See Tables 1 and 2**  
246 **the supplemental materials section**)

### 247 **13 Most Shared Genes in the 191 Unique Phage**

248 4 genes that stand out the most due to the amount of phage they were found in (**table**  
249 **2**). KHDAMHGJ\_00009 Chorismate synthase, found in M0471 (*Staphylococcus Aureus*  
250 *Subsp. Aureus NCTC 8325 (ID 57795) - BioProject - NCBI, n.d.*), was identified in 17 phage  
251 clusters. Its gene function is shikimate pathway, which shows signs of AMR in plants  
252 (**Mander and Liu 2010**). EOLKNJBM\_00007 Nucleoside diphosphate kinase in

253 ASM1150v1\_genomic.gbff\_pp18.ffn (**Holden et al., 2004**) was found in 16 phage  
254 clusters. MIIMDJNA\_00002 Heptaprenyl diphosphate synthase component 2 in  
255 ASM24879 (*Staphylococcus Aureus Subsp. Aureus CIG1612 (ID 60683) - BioProject - NCBI,*  
256 n.d.) was identified in 15 clusters. HGDEFLKI\_00006 3-dehydroquinase synthase in  
257 M0877\_V1\_genomic.gbff\_pp18.ffn (*Staphylococcus Aureus Subsp. Aureus NCTC 8325 (ID*  
258 *57795) - BioProject - NCBI, n.d.*) was identified in 14 phage clusters.

## 259 **References**

- 260 1. Liu, L. (2014). Fields Virology, 6th Edition. Clinical Infectious Diseases, 59(4),  
261 613–613. <https://doi.org/10.1093/cid/ciu346>
- 262 2. Blair, J. M., Webber, M. A., Baylay, A. J., Ogbolu, D. O., & Piddock, L. J. (2015).  
263 Molecular mechanisms of antibiotic resistance. Nature reviews microbiology,  
264 13(1), 42-51.
- 265 3. Zeman, M., Mašlaňová, I., Indráková, A., Šiborová, M., Mikulášek, K.,  
266 Bendíčková, K., ... & Pantůček, R. (2017). Staphylococcus sciuri bacteriophages  
267 double-convert for staphylokinase and phospholipase, mediate interspecies  
268 plasmid transduction, and package mecA gene. Scientific reports, 7(1), 1-11.
- 269 4. Scharn, C. R., Tenover, F. C., & Goering, R. V. (2013). Transduction of  
270 Staphylococcal Cassette Chromosome mec Elements between Strains of  
271 Staphylococcus aureus. *Antimicrobial Agents and Chemotherapy*, 57(11), 5233–5238.

- 272 <https://doi.org/10.1128/AAC.01058-13>
- 273 5. Klein, E., Smith, D. L., & Laxminarayan, R. (2007). Hospitalizations and deaths  
274 caused by methicillin-resistant *Staphylococcus aureus*, United States, 1999–2005.  
275 *Emerging infectious diseases*, 13(12), 1840.
- 276 6. Ramisetty, B. C. M., & Sudhakari, P. A. (2019). Bacterial ‘grounded’ prophages:  
277 hotspots for genetic renovation and innovation. *Frontiers in genetics*, 10, 65.
- 278 7. Scheffers, D. J., & Pinho, M. G. (2005). Bacterial cell wall synthesis: new insights  
279 from localization studies. *Microbiology and molecular biology reviews*, 69(4),  
280 585-607.
- 281 8. Fishovitz, J., Hermoso, J. A., Chang, M., & Mobashery, S. (2014).  
282 Penicillin-binding protein 2a of methicillin-resistant *Staphylococcus aureus*.  
283 *IUBMB life*, 66(8), 572-577.
- 284 9. (*Staphylococcus Aureus* (ID 154) - Genome - NCBI, n.d.)
- 285 10. Moura, A., Criscuolo, A., Pouseele, H., Maury, M. M., Leclercq, A., Tarr, C., ... &  
286 Brisse, S. (2016). Whole genome-based population biology and epidemiological  
287 surveillance of *Listeria monocytogenes*. *Nature microbiology*, 2(2), 1-10.
- 288 11. U.S. National Library of Medicine. (n.d.). Genome list - genome - NCBI. National  
289 Center for Biotechnology Information. Retrieved August 18, 2018, from  
290 <https://www.ncbi.nlm.nih.gov/genome/browse/#!/overview/>.
- 291 12. Krassowski, M., Das, V., Sahu, S. K., & Misra, B. B. (2020). State of the field in



- 292 multi-omics research: From computational needs to data mining and sharing.  
293 *Frontiers in Genetics*, 11.
- 294 13. Akhter, S., Aziz, R. K., & Edwards, R. A. (2012). PhiSpy: a novel algorithm for  
295 finding prophages in bacterial genomes that combines similarity- and  
296 composition-based strategies. *Nucleic acids research*, 40(16), e126.  
297 <https://doi.org/10.1093/nar/gks406>
- 298 14. Johnson, M., Zaretskaya, I., Raytselis, Y., Merezhuk, Y., McGinnis, S., & Madden,  
299 T. L. (2008). NCBI BLAST: a better web interface. *Nucleic acids research*, 36(Web  
300 Server issue), W5–W9. <https://doi.org/10.1093/nar/gkn201>
- 301 15. Edgar R. C. (2004). MUSCLE: multiple sequence alignment with high accuracy  
302 and high throughput. *Nucleic acids research*, 32(5), 1792–1797.  
303 <https://doi.org/10.1093/nar/gkh340>
- 304 16. Thompson, J. D., Higgins, D. G., & Gibson, T. J. (1994). CLUSTAL W: improving  
305 the sensitivity of progressive multiple sequence alignment through sequence  
306 weighting, position-specific gap penalties and weight matrix choice. *Nucleic  
307 acids research*, 22(22), 4673–4680. <https://doi.org/10.1093/nar/22.22.4673>
- 308 17. Zhang, K. Y., Gao, Y. Z., Du, M. Z., Liu, S., Dong, C., & Guo, F. B. (2019). Vgas: a  
309 viral genome annotation system. *Frontiers in microbiology*, 10, 184.
- 310 18. Seemann T. (2014). Prokka: rapid prokaryotic genome annotation. *Bioinformatics  
311 (Oxford, England)*, 30(14), 2068–2069.

- 312 <https://doi.org/10.1093/bioinformatics/btu153>
- 313 19. Guo, F. B., Ou, H. Y., & Zhang, C. T. (2003). ZCURVE: a new system for  
314 recognizing protein-coding genes in bacterial and archaeal genomes. *Nucleic  
315 acids research*, 31(6), 1780-1789.
- 316 20. Mahram, A., & Herbordt, M.C. (2015). NCBI BLASTP on High-Performance  
317 Reconfigurable Computing Systems. *ACM Trans. Reconfigurable Technol. Syst.*,  
318 7, 33:1-33:20.
- 319 21. Bisong, E. (2019). Matplotlib and Seaborn. In *Building machine learning and  
320 deep learning models on google cloud platform* (pp. 151-165). Apress, Berkeley,  
321 CA.
- 322 22. Csardi, G., & Nepusz, T. (2006). The igraph software package for complex  
323 network research. *InterJournal, complex systems*, 1695(5), 1-9.
- 324 23. Hamers, L. (1989). Similarity measures in scientometric research: The Jaccard  
325 index versus Salton's cosine formula. *Information Processing and Management*,  
326 25(3), 315-18.
- 327 24. Woudenberg, T. van. (2021, January 8). Custom function for jaccard index in  
328 Igraph. Thabo J. van Woudenberg. Retrieved November 9, 2021, from  
329 [https://www.tvanwoudenberg.com/post/custom-function-for-jaccard-index-in-ig  
330 raph/](https://www.tvanwoudenberg.com/post/custom-function-for-jaccard-index-in-igraph/).
- 331 25. Aziz, R. K., Devoid, S., Disz, T., Edwards, R. A., Henry, C. S., Olsen, G. J., ... &

- 332 Xia, F. (2012). SEED servers: high-performance access to the SEED genomes,  
333 annotations, and metabolic models. *PloS one*, 7(10), e48053.
- 334 26. U.S. National Library of Medicine. (n.d.). Prokaryotic Genome Annotation  
335 Guide. National Center for Biotechnology Information. Retrieved November 9,  
336 2021, from [https://www.ncbi.nlm.nih.gov/genbank/genomesubmit\\_annotation/](https://www.ncbi.nlm.nih.gov/genbank/genomesubmit_annotation/).
- 337 27. Siguier, P., Pérochon, J., Lestrade, L., Mahillon, J., & Chandler, M. (2006).  
338 ISfinder: the reference centre for bacterial insertion sequences. *Nucleic acids*  
339 *research*, 34(suppl\_1), D32-D36.
- 340 28. U.S. National Library of Medicine. (n.d.). National database of antibiotic resistant  
341 organisms (NDARO) - pathogen detection - NCBI. National Center for  
342 Biotechnology Information. Retrieved November 9, 2021, from  
343 <https://www.ncbi.nlm.nih.gov/pathogens/antimicrobial-resistance/>.
- 344 29. Boutet, E., Lieberherr, D., Tognolli, M., Schneider, M., & Bairoch, A. (2007).  
345 UniProtKB/Swiss-Prot. *Methods in molecular biology* (Clifton, N.J.), 406, 89–112.  
346 [https://doi.org/10.1007/978-1-59745-535-0\\_4](https://doi.org/10.1007/978-1-59745-535-0_4)
- 347 30. Pruitt, K. D., & Maglott, D. R. (2001). RefSeq and LocusLink: NCBI gene-centered  
348 resources. *Nucleic acids research*, 29(1), 137-140.
- 349 31. U.S. National Library of Medicine. (n.d.). Taxonomy browser (*Staphylococcus*  
350 *aureus* subsp. *aureus* NCTC 8325). National Center for Biotechnology  
351 Information. Retrieved November 9, 2021, from

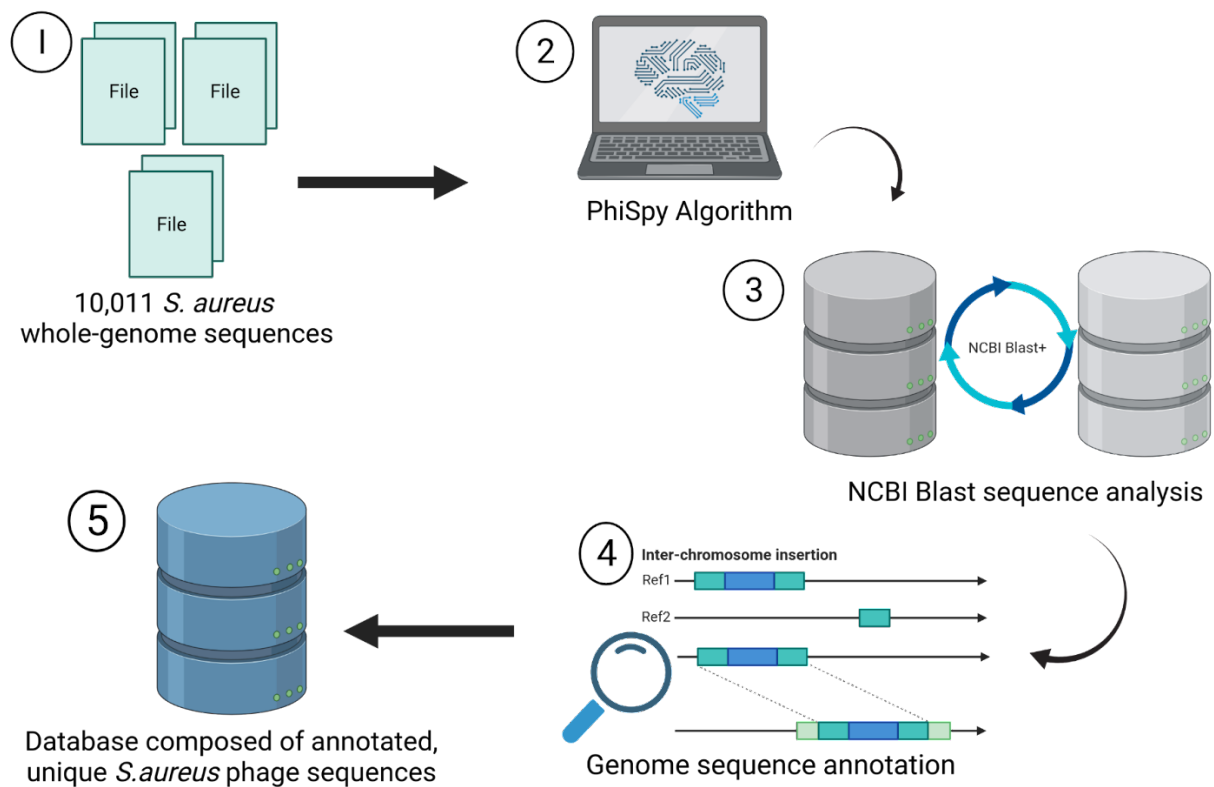
- 352 <https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?lvl=0&id=9>  
353 [3061](#).
- 354 32. Herron LL et al., "Genome sequence survey identifies unique sequences and key  
355 virulence genes with unusual rates of amino Acid substitution in bovine  
356 *Staphylococcus aureus*.", *Infect Immun*, 2002 Jul;70(7):3978-81
- 357 33. Herron-Olson L et al., "Molecular correlates of host specialization in  
358 *Staphylococcus aureus*.", *PLoS One*, 2007 Oct 31;2(10):e1120
- 359 34. U.S. National Library of Medicine. (n.d.). *Staphylococcus aureus* (ID 400143).  
360 National Center for Biotechnology Information. Retrieved November 9, 2021,  
361 from <https://www.ncbi.nlm.nih.gov/bioproject/400143>.
- 362 35. Manara, S., Pasolli, E., Dolce, D., Ravenni, N., Campana, S., Armanini, F., ... &  
363 Segata, N. (2018). Whole-genome epidemiology, characterisation, and  
364 phylogenetic reconstruction of *Staphylococcus aureus* strains in a paediatric  
365 hospital. *Genome medicine*, 10(1), 1-19.
- 366 36. Nübel, U., Dordel, J., Kurt, K., Strommenger, B., Westh, H., Shukla, S. K., ... &  
367 Witte, W. (2010). A timescale for evolution, population expansion, and spatial  
368 spread of an emerging clone of methicillin-resistant *Staphylococcus aureus*. *PLoS*  
369 *pathogens*, 6(4), e1000855.
- 370 37. Cameron, D. R., Ward, D. V., Kostoulias, X., Howden, B. P., Moellering Jr, R. C.,  
371 Eliopoulos, G. M., & Peleg, A. Y. (2012). Serine/threonine phosphatase *Stp1*

- 372 contributes to reduced susceptibility to vancomycin and virulence in  
373 *Staphylococcus aureus*. *The Journal of infectious diseases*, 205(11), 1677-1687.
- 374 38. Holden, M. T., Feil, E. J., Lindsay, J. A., Peacock, S. J., Day, N. P., Enright, M. C.,  
375 ... & Parkhill, J. (2004). Complete genomes of two clinical *Staphylococcus aureus*  
376 strains: evidence for the rapid evolution of virulence and drug resistance.  
377 *Proceedings of the National Academy of Sciences*, 101(26), 9786-9791.
- 378 39. Schuster, C. F., & Bertram, R. (2016). Toxin-Antitoxin Systems of *Staphylococcus*  
379 *aureus*. *Toxins*, 8(5), 140. <https://doi.org/10.3390/toxins8050140>
- 380 40. Fozo, E. M., Hemm, M. R., & Storz, G. (2008). Small toxic proteins and the  
381 antisense RNAs that repress them. *Microbiology and Molecular Biology*  
382 *Reviews*, 72(4), 579-589.
- 383 41. Labrie, S. J., Samson, J. E., & Moineau, S. (2010). Bacteriophage resistance  
384 mechanisms. *Nature reviews. Microbiology*, 8(5), 317-327.  
385 <https://doi.org/10.1038/nrmicro2315>
- 386 42. Hayes, F. (2003). Toxins-antitoxins: plasmid maintenance, programmed cell  
387 death, and cell cycle arrest. *Science*, 301(5639), 1496-1499.
- 388 43. Ono, H. K., Nishizawa, M., Yamamoto, Y., Hu, D. L., Nakane, A., Shinagawa, K.,  
389 & Omoe, K. (2012). Submucosal mast cells in the gastrointestinal tract are a target  
390 of staphylococcal enterotoxin type A. *FEMS Immunology & Medical*  
391 *Microbiology*, 64(3), 392-402.

392 44. Mander, L., & Liu, H. W. (2010). Comprehensive natural products II: chemistry  
393 and biology (Vol. 1). Elsevier.

394 **Supplemental Material**

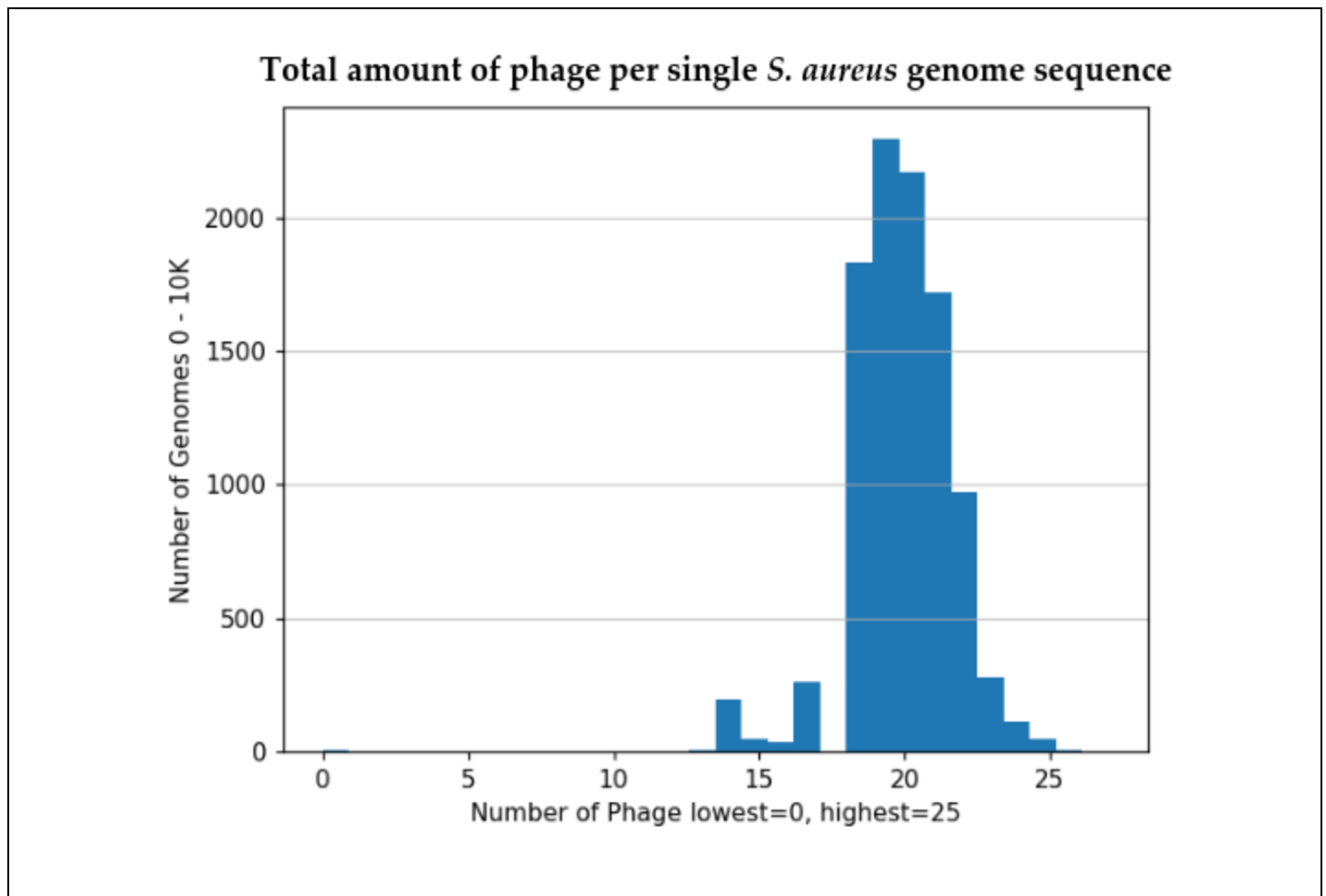
Pipeline Identifying and Characterizing Unique Prophage in *S. aureus* sequence data



**Figure 1: Pipeline Identifying and Characterizing Unique Prophage in *S. aureus* sequence data.**

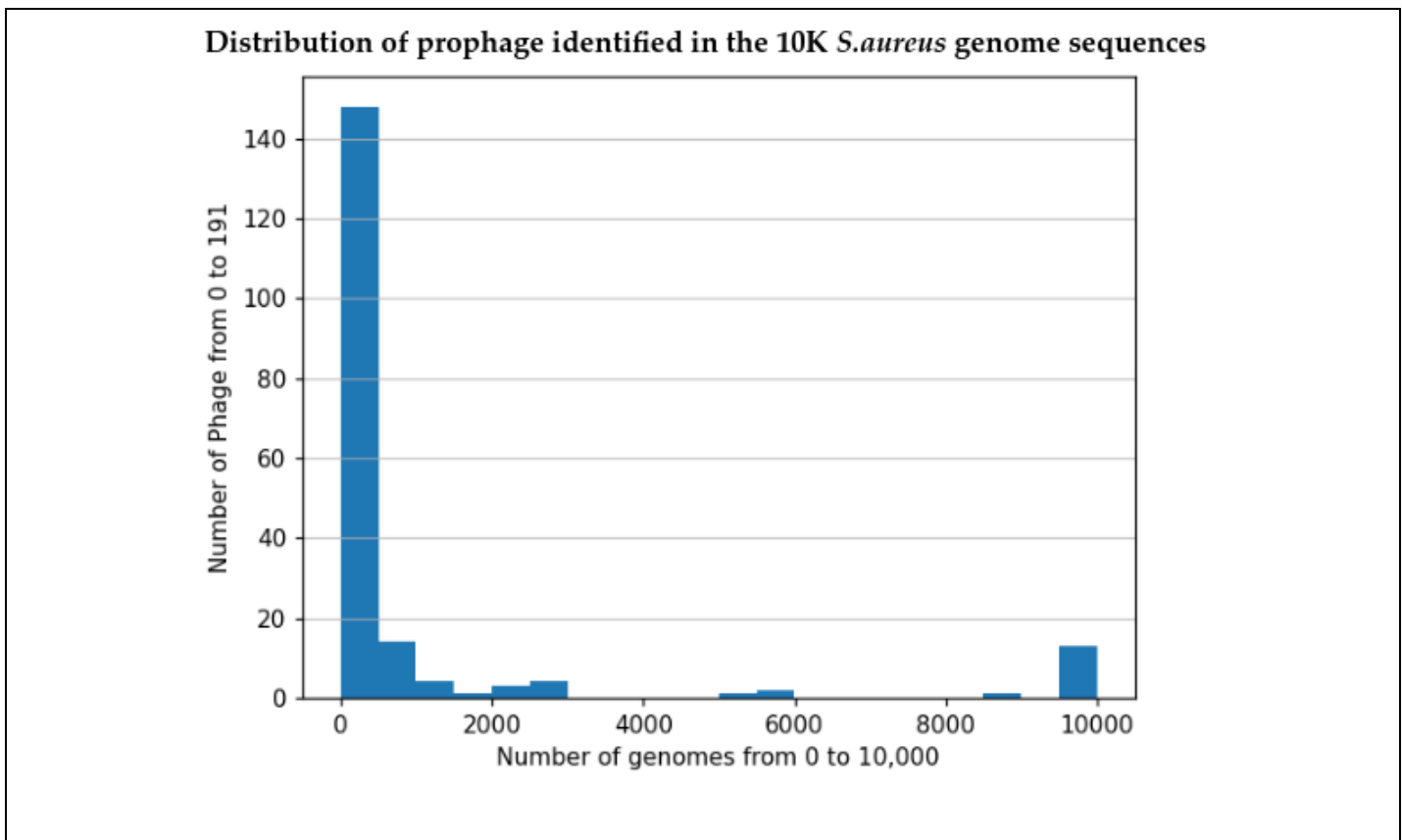
A visualization of the workflow used to identify unique prophage sequences. **1)** 10,011 *S. aureus* genome sequences were downloaded from the National Center for Biotechnology information (NCBI). **2)** The sequences were analyzed by PhiSpy. **3)** The fasta files for each predicted prophage were compared against each other using NCBI Blast nucleotide alignment tool. Prophage sequences that had 90% similarity along their full length were considered to be the same. **4)** Phage sequences were annotated using two independent methods (VGAS, Prokka) e. **5)** The resulting database of annotated, unique phage sequence allows for the identification of gene function encoded within prophage in *S. aureus*. (See **materials and methods** section above for more information)

395



**Figure 2: Total amount of phage per single *S. aureus* genome sequence.** Detected a total of 196,727 prophage sequences across the 10,000 genomes. There is an average of 39 (standard deviation = 1967.468) prophage sequences per genome. 1125 *S. aureus* sequences had 25 phage regions present, and 4 genomes had 0 phage detected. The x-axis reflects the exact totals each genome contained per genome (y-axis). (see **Analysis Uncovers 191 Unique Prophage Sequences** section above for more information).

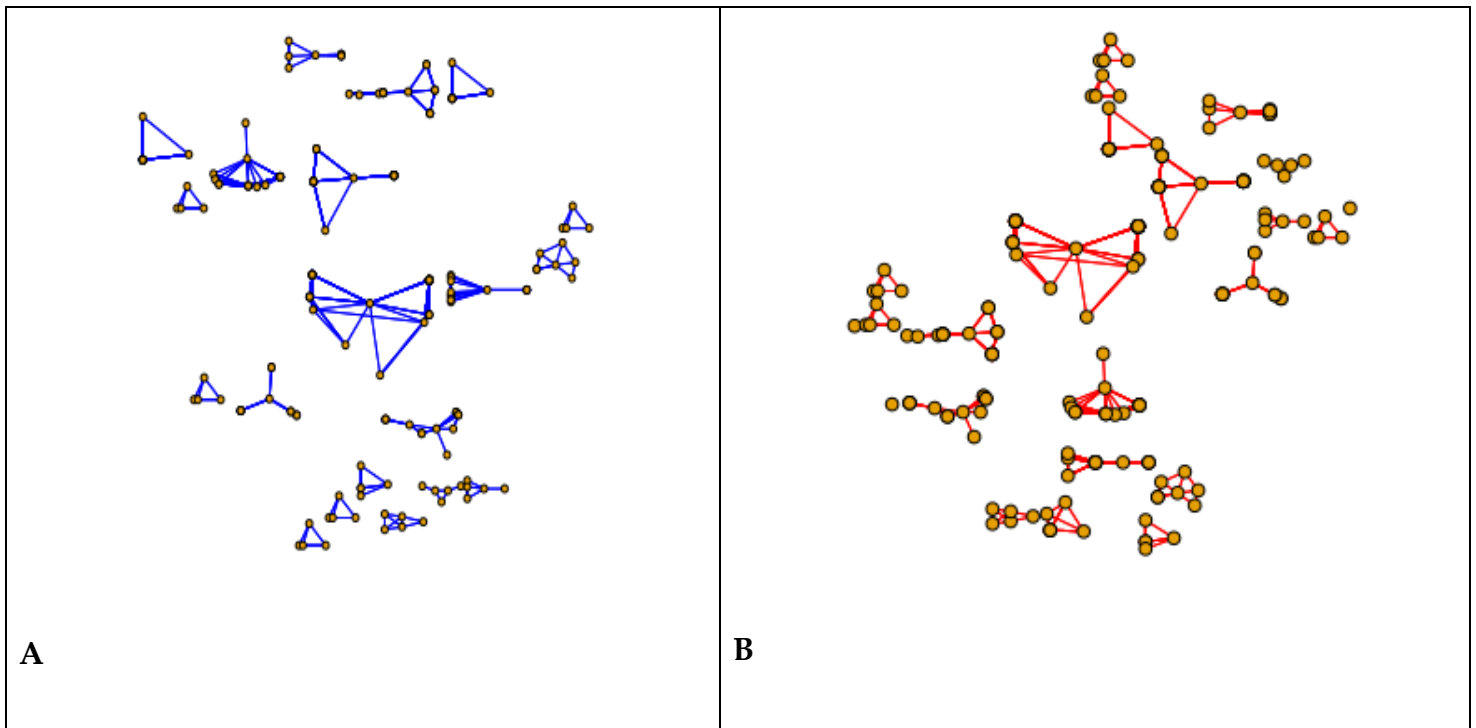
396



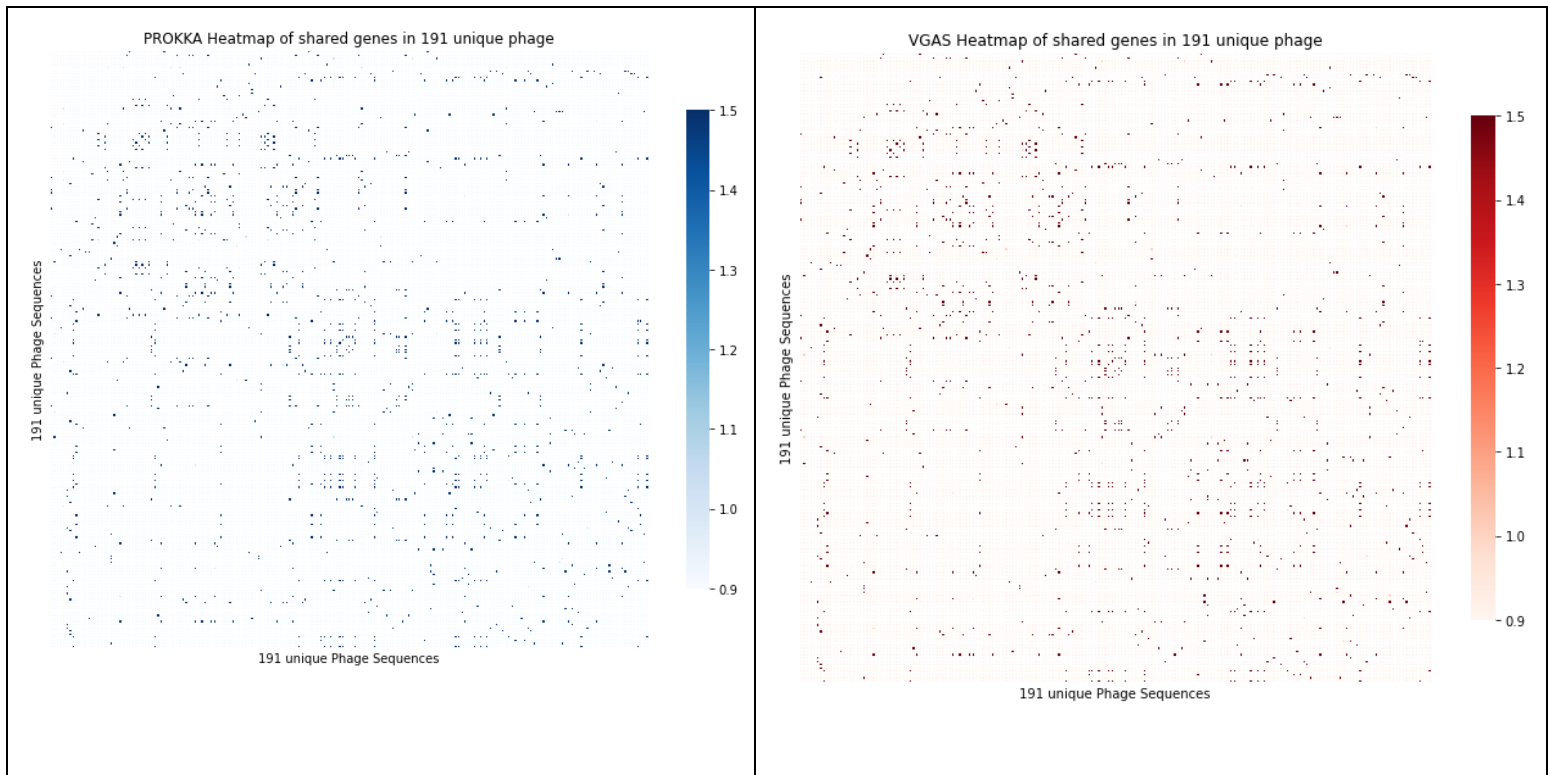
**Figure 3: Distribution of prophage identified in the 10K *S.aureus* genome sequences.** There were a few phage sequences detected in nearly every genome sequence. 191 prophages with an average of 1024 (Standard deviation = 2581.33). The highest number of genomes a single phage was detected is 9995, and the lowest was 1. (see **Analysis Uncovers 191 Unique Prophage Sequences** section above for more information).



397

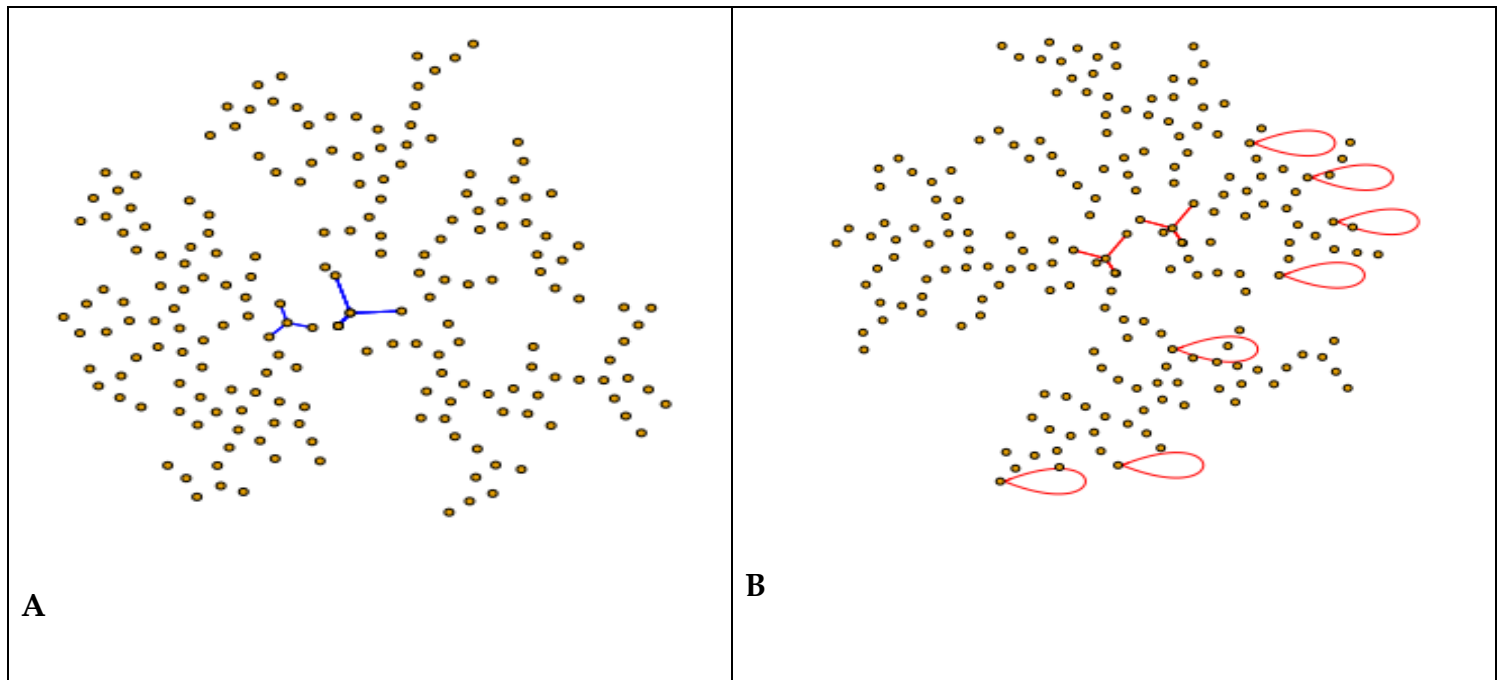


**Figure 4: PROKKA and VGAS Undirected Graphs shows shared ORFs between unique prophage sequences. A)** This graph shows the relationship between phage genomes by their gene content. Specifically, the nodes represent the 191 phage genomes and edges between nodes indicate the two phage share a gene (as annotated by Prokka). We determined that there were 1335 connected components with shared genes ranging from size 1 to 73. **B)** This graph shows the relationship between phage genomes by their gene content. Specifically, the nodes represent the 191 phage genomes and edges between nodes indicate the two phage share a gene (as annotated by VGAS). We determined that there were 1334 connected components with shared genes ranging from size 1 to 75. (See **Analysis Shows Shared ORFs between Unique Prophage Sequences** section above for more information).



**Figure 5: Heatmaps showing distribution of shared orfs between 191 unique phage sequences. Left)** This heatmap shows the numbers of genes shared between phage genomes as annotated by Prokka. These numbers ranged from 1 to 73. The X and Y axis are the 191 unique phage sequences. **Right)** This heatmap shows the numbers of genes shared between phage genomes as annotated by VGAS. These numbers ranged from 1 to 75. The X and Y axis are the 191 unique phage sequences. (See **Heat Maps** section above for more information).

399



**Figure 6: PROKKA and VGAS Undirected graph shows differences in shared ORFs between unique prophage sequences. A) 28 unique edges, 1306 shared with VGAS B) 27 unique edges, 1306 shared with Prokka.** Both the PROKKA and VGAS graphs shared the same range of connections with 1 as the lowest, and the highest at 22. (See **Analysis Shows Shared ORFs between Unique Prophage Sequences** section above for more information).

400