

# 1 **zol & fai: large-scale targeted detection and evolutionary investigation of gene clusters**

2  
3 Rauf Salamzade<sup>1,2</sup>, Patricia Tran<sup>3,4</sup>, Cody Martin<sup>2,3</sup>, Abigail L. Manson<sup>5</sup>, Michael S. Gilmore<sup>5,6,7</sup>,  
4 Ashlee M. Earl<sup>5</sup>, Karthik Anantharaman<sup>3</sup>, Lindsay R. Kalan<sup>1,8,9</sup>

5  
6 <sup>1</sup>Department of Medical Microbiology and Immunology, School of Medicine and Public Health, University of  
7 Wisconsin-Madison, Madison, WI, USA

8 <sup>2</sup>Microbiology Doctoral Training Program, University of Wisconsin-Madison, Madison, WI, USA

9 <sup>3</sup>Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA

10 <sup>4</sup>Freshwater and Marine Science Doctoral Program, University of Wisconsin-Madison, WI, USA

11 <sup>5</sup>Infectious Disease and Microbiome Program, Broad Institute of MIT and Harvard, Cambridge, Massachusetts, USA

12 <sup>6</sup>Department of Ophthalmology, Harvard Medical School and Mass Eye and Ear, Boston, Massachusetts, USA

13 <sup>7</sup>Department of Microbiology, Harvard Medical School and Mass Eye and Ear, Boston, Massachusetts, USA

14 <sup>8</sup>Department of Medicine, Division of Infectious Disease, School of Medicine and Public Health, University of  
15 Wisconsin-Madison, Madison, WI, USA

16 <sup>9</sup>M.G. DeGrootte Institute for Infectious Disease Research, David Braley Centre for Antibiotic Discovery, Department  
17 of Biochemistry and Biomedical Sciences, McMaster University, Hamilton, Ontario, Canada

18  
19 Address for Correspondence: Lindsay R. Kalan; kalanlr@mcmaster.ca

## 20 21 **Abstract**

22  
23 Many universally and conditionally important genes are genomically aggregated within  
24 clusters. Here, we introduce fai and zol, which together enable large-scale comparative analysis  
25 of different types of gene clusters and mobile-genetic elements (MGEs), such as biosynthetic  
26 gene clusters (BGCs) or viruses. Fundamentally, they overcome a current bottleneck to reliably  
27 perform comprehensive orthology inference at large scale across broad taxonomic contexts and  
28 thousands of genomes. First, fai allows the identification of orthologous or homologous  
29 instances of a query gene cluster of interest amongst a database of target genomes.  
30 Subsequently, zol enables reliable, context-specific inference of protein-encoding ortholog  
31 groups for individual genes across gene cluster instances. In addition, zol performs functional  
32 annotation and computes a variety of statistics for each inferred ortholog group. These  
33 programs are showcased through application to: (i) longitudinal tracking of a virus in  
34 metagenomes, (ii) discovering novel population-genetic insights of two common BGCs in a  
35 fungal species, and (iii) uncovering large-scale evolutionary trends of a virulence-associated  
36 gene cluster across thousands of genomes from a diverse bacterial genus.

## 37 38 **Introduction**

39  
40 Within bacterial genomes, genes are often co-located within smaller genetic structures  
41 such as operons<sup>1,2</sup>, phages<sup>3</sup>, metabolic gene clusters<sup>4</sup>, biosynthetic gene clusters (BGCs)<sup>5</sup>, and  
42 pathogenicity islands<sup>6,7</sup>. Although less prevalent, eukaryotic genomes also contain genes  
43 aggregated within discrete clusters<sup>5,8</sup>.

44 Sometimes gene clusters are highly conserved, encoding for products essential to the  
45 survival of the organism<sup>9</sup>. In other cases, a single gene cluster can exhibit variability in gene  
46 carriage and order across different strains or species<sup>10-12</sup>. This is often the case for BGCs

47 encoding specialized metabolites or virulence-associated gene clusters, where evolution of  
48 gene content and sequence divergence can influence fitness and contribute to adaptation within  
49 a changing ecosystem.

50 Bioinformatic toolkits to perform accurate pangenomic and comparative genomic  
51 analyses have been heavily developed over the past two decades<sup>13–18</sup>; however, tool  
52 development to aid the identification and comparative analysis of smaller homologous gene  
53 clusters has been more limited and largely designed for specific types of gene clusters<sup>19–22</sup>. In  
54 addition, while methods for comprehensive comparative genomics within species exist and are  
55 scalable<sup>17,23,24</sup>, methods for reliable, large-scale comparative genomics of thousands of  
56 genomes representing a greater breadth of taxonomic diversity are lacking and bear heavy  
57 computational costs<sup>25,26</sup>. Context-specific inference of orthologous genes within focal gene  
58 clusters offers a targeted and reliable solution to overcome challenges with scalability<sup>27,28</sup>. Such  
59 an approach was recently taken to infer orthologous genes between instances of homologous  
60 BGCs<sup>22</sup>.

61 Here, we introduce *fai* (*find-additional-instances*) and *zol* (*zoom-on-locus*), which are  
62 designed for the identification (*fai*) and in-depth evolutionary genomics investigations (*zol*) of a  
63 wide array of gene cluster types. We demonstrate the utility of these programs through  
64 application to three types of gene clusters within different genomic contexts including a novel  
65 bacteriophage within environmental metagenomes, a fungal secondary metabolite encoding  
66 biosynthetic gene clusters, and a conserved polysaccharide antigen locus within the diverse  
67 bacterial genus of *Enterococcus*.

68

## 69 Results

70

### 71 **fai and zol allow for the rapid inference of gene cluster orthologs across diverse** 72 **genomes**

73

74 The two programs, *fai* and *zol*, build upon approaches we recently reported in *IsaBGC*<sup>29</sup>  
75 that were developed to investigate evolutionary trends of BGCs in a single taxon. Within *fai* and  
76 *zol*, algorithmic adjustments have been implemented to broaden the application for searching  
77 any type of gene cluster across a diverse set of target genomes (Figure 1A). First, *fai* allows  
78 users to rapidly search for gene cluster instances in a target set of genomes. Then, *zol* can be  
79 used to compute evolutionary statistics and functional annotations of gene cluster content in  
80 table-based reports. Importantly, because *fai* has an option to filter secondary, potentially  
81 paralogous, instances of gene clusters found in target genomes, downstream *ab initio* clustering  
82 of proteins using a flexible, InParanoid-type algorithm<sup>14</sup> by *zol* can be used to reliably infer  
83 ortholog groups.

84 In addition to filtering secondary instances of query gene clusters identified in target  
85 genomes, detection criteria in *fai* can be adjusted by assessing whether gene cluster homologs  
86 lie near scaffold edges in target genomic assemblies. This feature overcomes challenges  
87 inherent to the identification of full gene-clusters in metagenomic assemblies or metagenome-  
88 assembled genomes, which can be highly fragmented (Figure S1). *fai* can further accept query  
89 gene-clusters in different formats to ease searching for gene clusters and genomic islands  
90 cataloged in databases such as ICEberg<sup>30</sup>, MIBiG<sup>31</sup>, or IslandViewer<sup>32</sup>. In addition, to promote

91 consistency in gene calling across target genomes, we have incorporated computationally light-  
92 weight dependencies for *de novo* gene prediction in prokaryotic genomes<sup>33,34</sup> and gene-  
93 mapping in eukaryotic genomes<sup>35</sup> within prepTG, to prepare and format target genomes for  
94 optimized gene-cluster searching in fai (Figure 1B). Together these unique features and options  
95 differentiate fai from other software with similar functionalities, such as cblaster<sup>21</sup> (Figure 1C,  
96 S1; Table S1; Supplementary Text).

97 zol is differentiated from *IsaBGC*<sup>29</sup>, where ortholog groups are inferred across full  
98 genomes using OrthoFinder<sup>18</sup>, by delineating ortholog groups within the context of a  
99 homologous or orthologous set of gene clusters, similar to the approach taken within  
100 CORASON<sup>22</sup> to visualize similarities between BGCs. While CORASON uses bidirectional best-  
101 hits to identify direct orthologs, zol accounts for the presence of in-paralogs and  
102 comprehensively partitions proteins into ortholog groups. Similar to *IsaBGC-PopGene*<sup>29</sup>, zol will  
103 then construct a tabular report with information on conservation, evolutionary trends, and  
104 annotation for individual ortholog groups (Figure 1D). To make annotated reports generated by  
105 zol more broadly informative for a variety of gene clusters, several databases have been  
106 included, such as VOGs<sup>36</sup>, VFDB<sup>37</sup>, ISFinder<sup>38</sup>, and CARD<sup>39</sup>. In addition, zol incorporates  
107 HyPhy<sup>40</sup> as a dependency and calculates evolutionary statistics not previously reported in  
108 *IsaBGC-PopGene*, such as sequence entropy in the 100 bp upstream of an ortholog group,  
109 where important regulatory differences could exist<sup>41</sup>. Ultimately, beyond high-throughput  
110 inference of ortholog groups across diverse genomic datasets, the rich tabular report produced  
111 by zol provides complementary information to figures generated by comparative visualization  
112 software such as clinker<sup>42</sup>, CORASON<sup>22</sup>, gggenomes<sup>43</sup>, and Easyfig<sup>44</sup>.

113 Another key feature in zol is the ability to dereplicate gene clusters directly using skani<sup>45</sup>,  
114 which was recently shown to be more reliable at estimating ANI between genomes of variable  
115 contiguity relative to comparative methods. Dereplication allows for more appropriate inference  
116 of evolutionary statistics to overcome availability or sampling biases in genomic databases<sup>46</sup>.  
117 Finally, zol allows for comparative investigations of gene-clusters based on taxonomic or  
118 ecological groupings<sup>47-49</sup>. For instance, users can designate a subset of gene clusters as  
119 belonging to a specific population to allow zol to calculate ortholog group conservation across  
120 just the focal set of gene clusters. In addition, if comparative investigations are requested, zol  
121 will also compute the fixation index<sup>50</sup>,  $F_{ST}$ , for each ortholog group to assess gene flow between  
122 the focal and complementary sets of gene clusters.

123

## 124 **Longitudinal tracking of a virus within lake metagenomic assemblies**

125

126 Viruses are important members of host and environmental microbiomes<sup>51-53</sup>, influencing  
127 the microbial composition and participating in several metabolic pathways. Targeted  
128 identification of a specific virus or bacteriophage within metagenomes can thus offer greater  
129 insight into their elusive functional roles in microbiomes.

130 Recently, changes in the composition and function of the metagenome at three different  
131 depths of a lake was reported using longitudinal shotgun metagenomics<sup>54</sup>. Using metagenome  
132 assemblies generated from this dataset, large ( $\geq 20$ kb) and predicted-circular phages were  
133 identified independently across a subset of metagenomes from the three different depths at the  
134 the earliest sampling date using VIBRANT<sup>55</sup>. Subsequent clustering based on the sequence and

135 syntenic similarity of protein domains identified a ~36kb highly conserved virus in two  
136 metagenomes sampled from lower lake depths.

137 *fai* was then used to perform a rapid, targeted search for this ~36kb *Caudovirales* virus  
138 across the full set of 16 metagenomes to identify additional instances of the virus. *fai* completed  
139 its search of the metagenomes, featuring >20 million proteins and 10.7 million contigs, in less  
140 than seven minutes using 20 threads. Of the 16 total metagenomes, spanning five distinct  
141 sampling timepoints and four distinct sampling depths, nine metagenomes containing the virus  
142 were identified (Figure 2A) exclusively from anoxic conditions ( $p=8.7e-5$ ; two-sided Fisher's  
143 exact test). This suggests the viral host likely performs anaerobic respiration. Application of *zol*  
144 further revealed that 34 (64%) of the 53 total distinct ortholog groups were core to all instances  
145 of the virus across nine metagenomes and completely conserved in sequence over the course  
146 of 2.5 months (Figure 2B; Table S2). Furthermore, seven of the 53 ortholog groups were not  
147 observed in the query viruses from the earliest sampling date, demonstrating the ability of *fai* to  
148 identify new genes within additional instances of known gene clusters.

149

### 150 **Investigating population-level and species-wide evolutionary trends of BGCs in the** 151 **eukaryotic species *Aspergillus flavus***

152

153 The fungal genus of *Aspergillus* is a source of several natural products, including  
154 aflatoxins, a common and economically impactful contaminant of food. The genus also contains  
155 species that are model organisms for studying fungal secondary metabolism<sup>56-58</sup>. Examination  
156 of the secondary metabolome of *A. flavus* has revealed that different clades or populations  
157 comprising certain species can exhibit variability in their metabolite production despite high  
158 conservation of core BGC genes encoding enzymes for synthesis of these metabolites<sup>12,59,60</sup>.  
159 For instance, population B *A. flavus* were identified as producing a greater abundance of the  
160 insecticide leporin B relative to populations A and C<sup>12,61</sup>.

161 To further understand the genomic basis for differences in metabolite content between  
162 populations, we investigated the leporin BGC using *fai* and *zol*. While the leporin cluster was  
163 previously identified as a core component of the *A. flavus* genome<sup>12</sup>, a recent study suggested  
164 that the full BGC was specific to a single clade from the species<sup>60</sup>. Low sensitivity in direct  
165 assessment of gene cluster presence in eukaryotic genome assemblies can arise from their  
166 incompleteness, leading to gene clusters being fragmented across multiple scaffolds, and  
167 challenges in *ab initio* gene prediction<sup>62,63</sup>. Further deterring the direct prediction of gene  
168 clusters in eukaryotic assemblies is the lack of gene annotations, with only 11 (5.1%) of 216 *A.*  
169 *flavus* genomes in NCBI's GenBank database having coding sequence predictions (Figure 3A).  
170 Therefore, we used *miniprot*<sup>35</sup>, which is integrated within *prepTG*, to directly map high-quality  
171 coding genes predictions based on transcriptomics data from the genome of strain *A. flavus*  
172 NRRL 3357<sup>64</sup> to the 216 genomes available for the species. Running *fai* in "draft mode" led to  
173 the identification of the leporin BGC within 212 (98.1%) assemblies, consistent with prior read  
174 mapping-based investigations<sup>12</sup>. This increase in sensitivity when *fai* is run with *miniprot*-based  
175 gene-mapping is substantial when compared to common alternate approaches for identifying  
176 homologous instances of BGCs across genomes (Figure 3B; Supplementary Text).

177 Of the 212 genomes with the leporin BGC, 202 contain instances that were not near  
178 scaffold edges. This set of 202 instances of the gene cluster were further investigated using *zol*,

179 with comparative investigation of BGC instances from *A. flavus* population B genomes to  
180 instances from other populations requested. High sequence conservation was observed for all  
181 genes in the leporin gene cluster as previously reported<sup>12</sup> (Table S3). Further, alleles for genes  
182 in the BGC from population B genomes were generally more similar to each other than to alleles  
183 from outside the population as indicated by high  $F_{ST}$  values ( $>0.85$  for 9 of 10 genes) (Figure  
184 3C; Table S3). While regulation of secondary metabolites in *Aspergillus* is complex<sup>65</sup>, zol  
185 analysis showed that the three essential genes for leporin production<sup>61</sup> also had the lowest  
186 variation in the 100 bps upstream their exonic coordinates (Figure S2). This suggests higher  
187 variability is occurring in the transcription of the accessory *lep* genes within the species. This  
188 supports experimental evidence that has shown gene knockouts depleting certain leporin  
189 species will still permit the production of others<sup>61</sup>.

190 *fai* and *zol* were also applied to the BGC encoding aflatoxin across *A. flavus*<sup>66</sup> (Table  
191 S4). Similar to the leporin BGC, the aflatoxin BGC was highly prevalent in the species and found  
192 in 71.8% of genomes. However, in contrast to the leporin BGC, the aflatoxin BGC contains  
193 several genes with positive Tajima's D values, indicating greater sequence variability for these  
194 coding regions across the species (Figure 3D). One of the genes with a positive Tajima's D  
195 value is *afIX*, which has been shown to influence conversion of the precursor vericolorin A to  
196 downstream intermediates in the aflatoxin biosynthesis pathway<sup>67</sup> (Figure 3E). An abundance of  
197 sites with mid-frequency alleles in the oxidoreductase encoding gene could represent granular  
198 control for the amount of aflatoxin relative to intermediates produced. The polyketide synthase  
199 gene *pksA* had the lowest Tajima's D value of -2.4, which suggests it is either highly conserved  
200 or under purifying selection (Figure 3F). In addition, because a recent predicted reference  
201 proteome was used to infer genomic coding regions, *fai* and *zol* detected several highly  
202 conserved genes within the aflatoxin BGC that are not represented in the original reference  
203 gene cluster input for *fai*<sup>31</sup>. This includes a gene annotated as a noranthrone monooxygenase  
204 recently characterized as contributing to aflatoxin biosynthesis<sup>68,69</sup> (Figure 3D).

205

## 206 **Large-scale identification of the Enterococcal polysaccharide antigen and assessment of** 207 **context restricted orthology inference**

208

209 The Enterococcal polysaccharide antigen (Epa) is a signature component of the cellular  
210 envelope of multiple species within *Enterococcus*<sup>70-73</sup>, which has mostly been characterized in  
211 the species *Enterococcus faecalis*<sup>70,74-77</sup>. While molecular studies have provided evidence that  
212 the locus contributes to enterococcal host colonization<sup>76</sup>, evasion of immune systems<sup>78</sup>, and  
213 sensitivity to antibiotics<sup>79</sup> and phages<sup>79,80</sup>, it was only recently that the structure of Epa was  
214 resolved and a model for its biosynthesis and localization formally proposed<sup>77</sup>. A homologous  
215 instance of the *epa* locus was identified in the other prominent pathogenic species from the  
216 genus, *Enterococcus faecium*<sup>71,73,81</sup>; however, the prevalence and conservation of *epa* across  
217 the diverse genus of *Enterococcus*<sup>82-84</sup> remains poorly studied.

218 *fai* was used to search for homologous instances of *epa* across 5,291 *Enterococcus*,  
219 genomes estimated by GTDB to represent 92 species<sup>85</sup>, using a sensitive searching criterium  
220 and coordinates of the locus along the *E. faecalis* V583 genome as a reference<sup>75,77</sup>  
221 (Supplementary Text). For detection of *epa* orthologous regions, co-location of at least seven of  
222 the 14 *epa* genes previously identified as conserved in both *E. faecalis* and *E. faecium* was

223 required. The default threshold for syntenic conservation of homologous instances to the query  
224 gene cluster was also disregarded to increase sensitivity for the detection of *epa* in more  
225 distantly related enterococcal species to *E. faecalis*. To allow for capture and downstream  
226 analysis of auxiliary genes which might be species or strain-specific but related to *Epa*  
227 production or decoration, 20 kb flanking contexts of the core *epa* genes identified in each target  
228 genome were extracted.

229 Using these criteria, 5,085 (96.1%) genomes from across the genus were found to  
230 possess an *epa* locus, confirming the locus as nearly core to the genus. Visual inspection of the  
231 *epa* genes among 463 representative *Enterococcus* genomes revealed that the core genes  
232 *epaA-epaR* are highly conserved in three of four major clades (Figure 4; Supplementary Text).  
233 Based on the detection criteria in *fai*, the *epa* locus in the fourth clade, previously referred to as  
234 the *Enterococcus columbae* group<sup>82</sup>, was either missing or encoded for highly divergent  
235 homologs of these genes. This clade includes *Enterococcus gallinarum*, one of the only other  
236 species in the genus, besides *E. faecalis* and *E. faecium*, reported to cause nosocomial  
237 outbreaks<sup>86,87</sup>.

238 Evolutionary trends and sequence diversity for individual genes with the *epa* locus, were  
239 next computed using *zol* after assessing *zol*'s reliability for gene cluster context-limited inference  
240 of orthology and the impact of dereplication on the calculation of evolutionary statistics by *zol*.

#### 241 **Gene-context specific orthology inference using *fai* and *zol* are concordant with genome-** 242 **wide ortholog group predictions**

243  
244  
245 Genome-wide orthology inference is currently difficult to scale to hundreds or thousands  
246 of genomes belonging to multiple species. However, orthology inference can be made more  
247 accessible if larger loci are first identified as orthologous between genomes, through leveraging  
248 syntenic support<sup>23,27</sup>. To assess whether ortholog group inference was reliable when *zol* is  
249 applied on orthologous gene clusters identified across multiple species, we ran *zol* on high-  
250 quality instances of the *epa* locus from 42 different species (Figure 5C). Ortholog group  
251 predictions by *zol* were then compared to genome-wide orthology predictions by OrthoFinder<sup>18</sup>,  
252 which has been shown to yield highly accurate predictions in benchmarking experiments  
253 involving genomes from multiple species<sup>88</sup>. Orthology predictions were highly concordant  
254 between *zol* and OrthoFinder for proteins from diverse instances of the *epa* locus. *zol* identified  
255 23,623 pairs of proteins within ortholog groups, of which 22,843 (96.70%) were also grouped  
256 together by OrthoFinder. Only 1,520 (6.24%) pairs of *epa*-associated proteins which were  
257 identified by OrthoFinder to belong to the same ortholog group were missed by *zol*.

258 Because the *epa* locus encodes multiple characterized and putative  
259 glycosyltransferases<sup>89</sup>, we used phylogenetics to examine the relationship between proteins  
260 belonging to ortholog groups with glycosyltransferase domains to confirm that major clades  
261 correspond to distinct ortholog group designations (Figure 5B). *zol* also has an option to “re-  
262 inflate” ortholog groups, expanding them to include proteins from gene clusters which were  
263 deemed redundant during dereplication. To demonstrate the scalability of *zol*, this “re-inflation”-  
264 based approach was next applied on the full set of high-quality and contiguous *epa* instances  
265 and a comprehensive phylogeny of ortholog groups corresponding to glycosyltransferases was  
266 constructed. In concordance with our analysis of the 42 representative genomes, distinct

267 phylogenetic clades for glycosyltransferases corresponded to different ortholog groups identified  
268 by zol (Figure 5C).

269

## 270 Dereplication can impact taxa-wide inferences of selection-informative statistics

271

272 Dereplication, or removal of redundant gene cluster instances, is important to consider  
273 when working with highly sequenced bacterial taxa, including *E. faecalis*, where certain  
274 lineages, such as those commonly isolated at clinics, can be overrepresented in genomic  
275 databases. Over-representation of select lineages will skew estimates for some evolutionary  
276 statistics, such as those informative of selective pressures, complicating evaluation of  
277 evolutionary trends across the entire taxonomic group. We thus assessed the impact of  
278 dereplication on the calculation of evolutionary statistics for instances of *epa* in *E. faecalis* using  
279 two different approaches: (i) genome-wide dereplication with dRep<sup>90,91</sup> and (ii) gene cluster  
280 specific dereplication with skani<sup>45</sup>. Dereplication at the gene cluster level with skani was  
281 performed directly in zol. The “re-inflation” option was also used to simulate comprehensive  
282 processing and calculation of evolutionary statistics while avoiding excessive computation.

283 Regardless of the approach for dereplication, genome-wide or gene cluster-specific, the  
284 estimates of evolutionary and genomic statistics for analogous ortholog groups were highly  
285 concordant (Figure 6, S3). However, gene cluster based dereplication can overestimate or  
286 underestimate selection informative statistics, such as Tajima’s D or FUBAR-based inference of  
287 the number of sites under selection, relative to genome-wide dereplication performed using  
288 similar thresholds. This is likely because the core *epa* locus is highly conserved across *E.*  
289 *faecalis* which led to fewer representative gene clusters following dereplication and a lower  
290 weight being placed on conserved alleles when estimating such statistics. In contrast, more  
291 simplistic statistics, such as average sequence entropy and the proportion of total alignment  
292 sites regarded as segregating sites, were closely estimated for genes regardless of the  
293 dereplication method used. In addition, using the “re-inflation” option in zol to infer orthology  
294 relationships across a comprehensive set of 1,232 high-quality and contiguous *epa* locus  
295 instances from the species produced concordant values for selection informative statistics to  
296 values generated using genome-wide based dereplication.

297

## 298 zol identifies genetic diversity of *epaX*-like glycosyltransferases

299

300 Because *Epa* biosynthesis and its conditional importance has mostly been investigated  
301 in *E. faecalis*<sup>70,74,75,77</sup>, we first examined evolutionary trends for proteins across instances of the  
302 *epa* locus from 75 *E. faecalis* representative genomes following genome-wide dereplication. In  
303 accordance with prior studies<sup>71,77</sup>, zol reported that one end of the locus corresponds to genes  
304 which are highly conserved and core to *E. faecalis* (*epaA-epaR*) whereas the other end contains  
305 strain-specific genes (Figure 7A; Table S5). Using zol, we further found that variably conserved  
306 genes exhibit high sequence dissimilarity, as measured using both Tajima’s D and average  
307 sequence entropy, in comparison to the core genes of the locus (Figure 7BC). Comparative and  
308 multi-species analysis of the *epa* locus between and across *E. faecalis* and *E. faecium* was next  
309 performed using gene cluster based dereplication with re-inflation using zol (Table S6). zol  
310 reported conservation statistics were consistently in agreement with previous studies<sup>71,73</sup>.

311 Twenty genes determined to be present in the majority (>95%) of *epa* clusters across both  
312 species, including *epa*ABCDEFGH, *epa*LM, and *epa*OPQR. In addition, default parameters for  
313 orthologous clustering of proteins in *zol* detected a known truncated variant of the  
314 glycosyltransferase *epaN* in *E. faecium*.

315 The gene *epaX*, encoding a glycosyltransferase, was identified as one ortholog group  
316 with the greatest sequence variation in *E. faecalis* (Figure 7BD, S4). *epaX* was previously  
317 shown to be critical for *E. faecalis* host-gut colonization and proposed to be involved in the  
318 decoration of the rhamnan backbone structure of Epa with galactose and N-acetyl  
319 glucosamine<sup>76</sup>. Comparative analysis using *E. faecium* as the focal taxa further showed that the  
320 *epaX*-containing ortholog group has a low  $F_{ST}$  value, indicating alleles from *E. faecalis* and *E.*  
321 *faecium* species are phylogenetically interspersed. This was confirmed through phylogenetic  
322 assessment of the ortholog group (Figure 7E). In addition, although some allelic clades encode  
323 sequences from both species, genes remained sub-partitioned by species. This phylogenetic  
324 structure for the ortholog group, together with our prior observation that the *epaX*-containing  
325 ortholog group in *E. faecalis* has greater sequence variability relative to other  
326 glycosyltransferases from the locus, suggests extensive and ancestral sequence evolution of  
327 *epaX*-like glycosyltransferases. Further, while only 70% of *E. faecium* found to carry *epa*  
328 possess an *epaX*-like ortholog group, approximately 7% of them encode the ortholog in multi-  
329 copy (Figure 7F), suggesting the occurrence of intra-locus gene duplication.

330

## 331 Discussion

332

333 Here *fai* and *zol* are introduced to enable large scale evolutionary investigations of gene  
334 clusters in diverse taxa. Together these tools overcome current bottlenecks in computational  
335 biology to infer orthologous sets of genes at scale across thousands of diverse genomes.

336 Both *fai* and *cblaster*<sup>21</sup> can be used to identify additional gene clusters within target  
337 genomes and extract them as GenBanks for downstream investigations using *zol*. For those  
338 lacking computational resources needed for *fai* analysis, *cblaster* offers remote searching of  
339 BGCs using NCBI's BLAST infrastructure and non-redundant databases. More recently,  
340 CAGECAT<sup>92</sup>, a highly accessible web-application for running *cblaster*, was also developed and  
341 can similarly be used to identify and extract gene cluster instances from genomes represented  
342 in NCBI databases. In contrast to these tools, *fai* contains algorithms and options for users  
343 interested in: (i) identifying gene clusters across a comprehensive or redundant set of genomic  
344 assemblies, (ii) improved sensitivity for gene cluster detection in draft-quality assemblies, and  
345 (iii) automated filtering of secondary, or paralogous, matches to query gene clusters. In addition,  
346 users can apply *zol* to further investigate homologous sets of gene clusters identified from  
347 IslandCompare<sup>93</sup>, BiG-SCAPE<sup>22</sup>, or vConTACT2<sup>94</sup> analyses, which perform comprehensive  
348 clustering of predicted genomic islands, BGCs, or viruses.

349 The utility of *fai* is demonstrated here through rapid, targeted detection of a virus directly  
350 from lake metagenomic assemblies. Targeted detection of specific viruses longitudinally  
351 presents an efficient and tractable approach to understand how viral pangenomes evolve over  
352 time. In addition, by permitting fragmented detection of gene clusters and detection of proximity  
353 to scaffold edges, users can assess whether phages or other gene clusters corresponding to  
354 MGEs are present in their metagenomes. *fai* and *zol* will continue to compliment metagenomic



355 applications as long-read sequencing becomes more economical and commonly used to profile  
356 microbial communities. For example, their application could be useful for assessing the  
357 presence of concerning MGEs conferring antimicrobial resistance traits<sup>95-97</sup> and identifying novel  
358 auxiliary genes within known BGCs which may tailor the resulting specialized metabolites and  
359 expand chemical diversity<sup>98,99</sup>.

360 Reidentifying gene-clusters in eukaryotic genomes remains difficult due to technical  
361 challenges in gene prediction owing to the presence of alternative splicing. The ability of fai and  
362 zol to perform population-level genetics on common BGCs from the eukaryotic species *A. flavus*  
363 was demonstrated. While there are over 200 genomes of *A. flavus* on NCBI, only 5.1% have  
364 coding-sequence information readily available. We used miniprot<sup>35</sup> to map high quality gene  
365 coordinate predictions from a representative genome in the species<sup>64</sup> to the remainder of  
366 genomic assemblies within prepTG which enabled high sensitivity detection of BGCs with fai.  
367 Our analysis provides additional support that the leporin BGC is conserved in full across the  
368 species<sup>12</sup> using an assembly-based approach.

369 Application of fai and zol to exopolysaccharide encoding gene clusters from pathogens  
370 of interest allows a better understanding of their conservation and evolutionary trends. This  
371 information can then aid the identification of potential genes to target for antivirulence  
372 efforts<sup>103,104</sup> or genes underlying host-pathogen interactions<sup>76,105</sup>. fai was used to identify  
373 orthologous instances of the *epa* locus, encoding for an extracellular polysaccharide antigen,  
374 across thousands of diverse genomes from the genus of *Enterococcus*. Subsequently,  
375 application of zol reliably produced comparable orthology predictions to OrthoFinder, a highly  
376 dependable genome-wide orthology inference software<sup>18,88</sup>. While zol missed a small  
377 percentage of orthologous instances identified by OrthoFinder in our testing, this could be due  
378 to threshold settings for percent identity and coverage between pairs of proteins set in zol. Such  
379 thresholds are not enforced in OrthoFinder. However, parameters controlling these thresholds  
380 are adjustable in zol and allow users to increase or decrease orthology sensitivity at the  
381 expense of incurring false positives as they deem appropriate for their research objective.

382 Using zol, it was determined that an ortholog group containing *epaX*-like  
383 glycosyltransferases possess high sequence divergence relative to other glycosyltransferases  
384 within the *epa* locus in *E. faecalis*. In addition to influencing the ability of *E. faecalis* to colonize  
385 hosts<sup>76</sup>, mutations in *epaX* and other genes from the ortholog group have also been shown to  
386 impact susceptibility to phage predation<sup>100-102</sup>. Thus, because similar *epaX*-like  
387 glycosyltransferases are found in both *E. faecalis* and *E. faecium*, we hypothesize that  
388 extensive ancestral evolution of the *epaX*-containing ortholog group may have occurred to  
389 support evasion from phages and confer colonization of new hosts. In this study, we further  
390 found that the *E. columbae* group might lack or possess highly divergent versions of core *epa*  
391 genes found in *E. faecalis* and *E. faecium*, suggesting that development of anti-virulence  
392 approaches to broadly target Epa in all pathogenic enterococci might be difficult to achieve.  
393 Similar investigations with fai and zol can readily be performed for other exopolysaccharide  
394 encoding gene clusters of pathogens to better understand their conservation, evolutionary  
395 trends, identify appropriate genes to target for antivirulence efforts<sup>103,104</sup>, and infer whether  
396 certain genes underlie host-pathogen interactions<sup>76,105</sup>.

397 Options for dereplication and re-inflation provided within zol enable scalability to  
398 thousands of gene cluster instances. The usage of these options can further aid in performing

399 more accurate evolutionary investigations for genes broadly across focal taxa or between  
400 clades, by overcoming biases due to overrepresentation of certain lineages in genomic  
401 databases<sup>12,47</sup>. Depending on the underlying origin of input gene clusters, zol can also be used  
402 to assess temporal<sup>48,106</sup> or spatial<sup>49</sup> evolutionary trends.

403 Practically, zol presents a comprehensive analysis tool for comparative genetics of  
404 related gene clusters to facilitate detection of evolutionary patterns that might be less apparent  
405 from visual analysis. Fundamentally, the algorithms presented within fai and zol enable the  
406 reliable detection of orthologous gene clusters, and subsequently orthologous proteins, across  
407 multi-species datasets spanning thousands of genomes and help overcome a key barrier in  
408 scalability for comparative genomics.

409

## 410 **Acknowledgments**

411

412 This work was supported by grants from the National Institutes of Health awarded to L.R.K  
413 (NIAID U19AI142720 and NIGMS R35GM137828) and the Broad Institute (U19AI110818). The  
414 content is solely the responsibility of the authors and does not necessarily represent the official  
415 views of the National Institutes of Health. The authors are grateful to James Kosmopoulos, Dr.  
416 Caitlin Sande, and Mary Hannah Swaney for feedback and assistance with data acquisition as  
417 well as Dr. Devon Ryan and Dr. Robert A. Petit III for assistance with incorporation of the suite  
418 into Bioconda.

419

## 420 **Methods**

421

### 422 **Software availability**

423

424 zol is provided as an open-source software suite, developed primarily in Python3 on GitHub at:  
425 <https://github.com/Kalan-Lab/zol>. Docker and Bioconda<sup>107</sup> based installations of the suite are  
426 supported. For the analyses presented in this paper, we used v1.2.0 of the zol software  
427 package. Minor patches, since incorporated into the software since v1.25, were added  
428 retrospectively to this version pertaining to safer acquisition of stored statistics when generating  
429 the final report. Version information for major dependencies of the zol suite<sup>33,35,40,45,108–115</sup> or  
430 software generally used<sup>22,55,116</sup> for analyses in this study is provided in [Supplementary Table S7](#).

431

### 432 **Data availability**

433

434 Genomes and metagenomes used to showcase the application of fai and zol are listed with  
435 GenBank accession identifiers in [Supplementary Table S8](#). Total metagenomes and their  
436 associated information from Lake Mendota microbiome samplings were originally described in  
437 Tran *et al.* 2023<sup>54</sup> and deposited in NCBI under BioProject PRJNA758276. Genomic assemblies  
438 available for *A. flavus* in NCBI's GenBank database on Jan 31st, 2023 were downloaded in  
439 GenBank format using ncbi-genome-download ([https://github.com/kblin/ncbi-genome-](https://github.com/kblin/ncbi-genome-download)  
440 [download](#)). Genomic assemblies for *Enterococcus* that met quality and taxonomic criteria for  
441 belonging to the genus or related genera (e.g. Enterococcus\_A, Enterococcus\_B, etc.) in  
442 GTDB<sup>85</sup> release R207 were similarly downloaded from NCBI's GenBank database using ncbi-  
443 genome-download in FASTA format.

444

## 445 **Application of fai and zol to identify phages within metagenomes**

446

447 VIBRANT was used to identify viral contigs or sub-contigs in the three total metagenomes from  
448 Tran *et al.* 2023<sup>54</sup> sampled on the earliest date of 07/24. Afterwards, predicted circular contigs  
449 were clustered using BiG-SCAPE<sup>22</sup> which revealed a ~36 kb virus was found in two of the three  
450 metagenomes.

451

452 prepTG was run on all 16 total metagenomic assemblies from the Tran *et al.* 2023 study,  
453 performing gene calling with pyrodigal in metagenomics mode<sup>33</sup> to prepare for comprehensive  
454 targeted searching of the virus. Afterwards, fai was run with default settings, with filtering of  
455 paralogous (or secondary) instances of the phage requested to retain only the best matching  
456 scaffold or scaffold segment resembling the queries.

457

## 458 **Microevolutionary investigations of leporin and aflatoxin BGCs in *Aspergillus flavus***

459

460 Genomic assemblies downloaded from NCBI GenBank were processed using prepTG. Of the  
461 217 genomic assemblies downloaded, one, GCA\_000006275.3, was dropped from the analysis  
462 because the original GenBank had multiple CDS features with the same name, leading to  
463 difficulties in performing BGC prediction with antiSMASH<sup>116</sup>, and because alternate assemblies  
464 were available for the isolate. prepTG was run on all assemblies with minimap2<sup>35</sup> based gene-  
465 mapping of the high-quality gene coordinate predictions available for *A. flavus* NRRL 3357  
466 (GCA\_009017415.1)<sup>64</sup> requested. Target genomes were then searched for the leporin  
467 (BGC0001445) and aflatoxin (BGC0000008) BGCs using GenBanks provided on MIBiGv3<sup>31</sup>.  
468 For leporin, AFLA\_066840, as represented in the MIBiG database, was treated as a key protein  
469 required for detection of the BGC. Similarly, for aflatoxin, PksA (AAS90022.1), as represented in  
470 the MIBiG database, was treated as a key protein required for detection of the BGC. Draft-mode  
471 and filtering of paralogous segments was requested but turned off by default.

472

473 We reidentified population B as previously delineated<sup>12</sup> using k-mer based ANI estimation<sup>117</sup> and  
474 neighbor-joining tree construction<sup>118</sup>. A discrete clade (n=81) in the tree was validated to feature  
475 all isolates previously determined as part of population B<sup>12</sup> and thus regarded as such.

476

477 For comprehensive and *de novo* BGC prediction, antiSMASH was run on the 216 genomic  
478 assemblies with 'glimmerhmm' requested for the option '--genefinding-tool'. BGCs were  
479 clustered using default settings in BiG-SCAPE with MIBiG reference BGC integration requested  
480 and a PKS-NRPS hybrid GCF was found to feature the leporin B BGC representative  
481 (BGC0001445). Only 65 (30.1%) of the 216 genomic assemblies featured this GCF, likely  
482 resulting from the use of distant gene models based on *Cryptococcus* genomes with  
483 glimmerhmm<sup>119</sup>. For remote clinker analysis, CAGECAT<sup>92</sup> was used to search NCBI's nr  
484 database with proteins from the leporin BGC representative (BGC0001445) provided as a  
485 query. Only 13 scaffolds, belonging to 12 assemblies (including GCA\_000006275.3), were  
486 identified.

487

## 488 **Evolutionary investigations of the *epa* locus across *Enterococcus***

489

490 All *Enterococcus* genomes represented in GTDB R207<sup>85</sup> (n=5,291) were downloaded using  
491 ncbi-genome-download and processed in prepTG with gene-calling performed using pyrodigal<sup>33</sup>.  
492 Coordinates extending from 2,071,671 to 2,115,174 along the *E. faecalis* V583 chromosome,  
493 corresponding to genes EF2164 to EF2200. When using direct coordinates along a reference,  
494 fai re-performs gene-calling along the reference and extracts a local GenBank corresponding to  
495 the region between the coordinates. Gene calling is performed using pyrodigal. Because prior  
496 comparative analyses had shown that gene-conservation and gene-order can be slightly  
497 variable between *epa* loci from *E. faecalis* and *E. faecium*<sup>71</sup>, we relaxed the syntenic similar to  
498 query in fai from 0.6 to 0.0 and minimum percentage of query proteins needed to report a  
499 homologous instance of the *epa* locus to 10%. Instead, we required the presence of 50% of key  
500 *epa* proteins found in both *E. faecalis* and *E. faecium*, *epa*ABCDEFGHIJLMOPQR, for the  
501 identification of valid homologous instances of the *epa* locus. To gather auxiliary genes flanking  
502 the core *epa* regions detected, we further requested the inclusion of CDS features found within  
503 20 kb of the boundary core *epa* genes.

504

### 505 Genome selection for comparing ortholog grouping of proteins by *zol* with OrthoFinder:

506 Genome-wide dereplication of all *Enterococcus* genomes using dRep<sup>90</sup> with fastANI<sup>91</sup> and a  
507 secondary ANI clustering threshold of 99.0% led to the identification of 463 distinct genomes,  
508 including 101 *E. faecalis* genomes. Of these 101 genomes, 75 had high-quality *epa* instances  
509 which were not located near scaffold edges. *zol* was run on the 75 high-quality *epa* instances  
510 using default ortholog grouping parameters and similarly OrthoFinder v2.5.4 was run using  
511 default settings on the full, genome-wide set of 75 proteomes. To assess the concordance  
512 between OrthoFinder and *zol* for more diverse gene-clusters, gathered from multiple species,  
513 dRep was applied a second time on the set of 463 *Enterococcus* genomes using an ANI  
514 threshold of 95.0% to approximate selection of one representative genome per species<sup>120</sup>. This  
515 secondary dereplication identified 89 genomes, of which 42 featured highly-quality instances of  
516 the *epa* locus.

517

518 Phylogenetic analysis of glycosyltransferases found in or near the *epa* locus: Ortholog groups  
519 from the *zol* analysis on the 42 representative and 2,442 comprehensive multi-species *epa*  
520 instances (Figure 5BC), as well as the 75 representative *E. faecalis* *epa* instances (Figure S4),  
521 were identified as glycosyltransferases if they featured the key words: “glycosyl” and  
522 “transferase” in Pfam protein domain annotations<sup>121</sup>. For each gene cluster set, protein  
523 sequences belonging to the ortholog groups were extracted, retaining association information  
524 with particular ortholog groups, and subsequently aligned using MUSCLE<sup>115</sup>. Alignment filtering  
525 was next performed using trimal with options “-keepseqs -gt 0.9”, sequences with greater than  
526 25% of sites being gaps were filtered, and an approximate maximum-likelihood phylogeny was  
527 finally constructed using FastTree2<sup>110</sup>, midpoint rooted, and visualized using iTol<sup>122</sup>. Ortholog  
528 groups were assigned to specific *epa* gene designations based on sequence alignment of *E.*  
529 *faecalis* V583 proteins.

530

531 Assessing the impact of dereplication on the calculation of evolutionary statistics computed by  
532 zol: To assess the impact of dereplication on the estimation of evolutionary statistics using zol,  
533 we focused on high-quality instances (<10% of bases ambiguous) of the *epa* locus that were not  
534 near scaffold edges from *E. faecalis* genomes. We ran dereplication at the genome scale using  
535 dRep<sup>90</sup> with fastANI<sup>91</sup> and a secondary ANI clustering threshold of 99.0% and dereplication at  
536 the gene-cluster scale using skani<sup>45</sup> at 99.0% identity and 99.0% coverage with single-linkage  
537 clustering. We additionally simulated comprehensive processing of all high-quality gene-clusters  
538 distant from scaffold edges using the re-inflation option in zol, which allows expansion of  
539 ortholog groups determined in the dereplicated gene cluster set to the full listing of gene-  
540 clusters. Comparisons of estimates for various evolutionary statistics by zol between the  
541 different dereplication approaches were performed by first identifying the best matching ortholog  
542 groups from the three distinct analyses to each *epa*-associated gene from EF2164 to EF2200 in  
543 the *E. faecalis* V583 reference genome based on E-value. Only ortholog groups which were  
544 found in single-copy within the *epa* context were considered.

545

## 546 **Figure Legends**

547

548 **Figure 1: Overviews of fai and zol.** **A)** A schematic of how prepTG, fai, and zol are integrated  
549 to perform evolutionary investigations by searching for gene clusters. An overview of the  
550 prepTG **(B)**, fai **(C)** and zol **(D)** algorithms and workflows.

551 **Figure 2: Targeted viral detection in metagenomes using fai.** **A)** Total metagenomes from a  
552 single site in Lake Mendota across multiple depths and timepoints from Tran et al. 2023 were  
553 investigated using fai for the presence of a virus found in two of the three earliest microbiome  
554 samplings (red box). The presence of the virus is indicated by a phage icon. Metagenome  
555 samples are colored according to whether they corresponded to oxic, oxycline, or anoxic. The  
556 most shallow sampling depths varied for different dates and consolidated as a single row  
557 corresponding to a sampling depth of either 5 or 10 meters. **B)** The pangenome of the virus is  
558 shown based on the consensus order and directionality of coding sequences inferred by zol. Bar  
559 heights correspond to the median length of coding sequences and are colored based on the  
560 percentages of the nine metagenomes the virus was detected in. Figure 2a was created with  
561 BioRender.com.

562 **Figure 3: Evolutionary trends of common BGCs in *A. flavus*.** **A)** The proportion of 216 *A.*  
563 *flavus* genomes from NCBI's GenBank database with coding-sequence predictions available. **B)**  
564 Comparison of the sensitivity of fai and alternate approaches based on assemblies for detecting  
565 the leporin BGC. The red-line indicates the total number of genomes (n=216) assessed. A  
566 schematic of the **(C)** leporin and **(D)** aflatoxin BGCs is shown with genes present in  $\geq 10\%$  of  
567 samples shown in consensus order and relative directionality. Coloring of genes in **(C)**  
568 corresponds to  $F_{ST}$  values and in **(D)** to Tajima's D values, as calculated by zol. Grey bars in  
569 the legends, at **(C)** 0.92 and **(D)** -0.98, indicate the mean values for the statistics across genes  
570 in the BGC. \*For the leporin BGC, *lepB* corresponds to an updated open-reading frame (ORF)  
571 prediction by Skerker et al. 2021 which was the combination of AFLA\_066860 and  
572 AFLA\_066870 ORFs in the MIBiG entry BGC0001445 used as the query for fai. For the  
573 aflatoxin BGC, ORFs which were not represented in the MIBiG entry BGC0000008 but  
574 predicted to be within the aflatoxin BGC by mapping of gene-calls from *A. flavus* NRRL 3357 by  
575 Skerker et al. 2021 are shown in gold. The major allele frequency distributions are shown for **(E)**  
576 *afIX* and **(F)** *pksA*, which depict opposite trends in sequence conservation according to their

577 respective Tajima's D calculations.

578 **Figure 4: The *epa* locus is conserved across most enterococcal species.** The distribution  
579 of the *epa* locus and associated genes, based on criteria used for running fai, is shown across  
580 463 representative genomes across *Enterococcus*. Coloring of the heatmap corresponds to the  
581 normalized bitscore of the best alignment to coding sequences from *E. faecalis* V583.

582  
583 **Figure 5: Assessment of gene-cluster restricted ortholog grouping by fai and zol.** A) zol  
584 gene-cluster constricted ortholog group predictions for *epa* locus proteins from 42 distinct  
585 representative enterococcal species were compared to genome-wide predictions of ortholog  
586 groups by OrthoFinder. A phylogeny based on gap-filtered protein alignments of ortholog groups  
587 with domains featuring "glycosyl" and "transferase" as key words is shown from (B) *epa* loci in  
588 the 42 representative genomes and (C) a more comprehensive set of 2,442 *epa* loci. Each node  
589 represents a specific protein and coloring of the track corresponds to their ortholog group  
590 designations by zol. Note, (B) 2 (0.07%) and (C) 79 proteins (0.4%) were removed prior to  
591 phylogeny construction due to an abundance of gaps in the trimmed alignment.

592  
593 **Figure 6: Effects of dereplication on the calculation of evolutionary statistics by zol.** The  
594 heatmap shows the correlation of values for analogous ortholog groups for various evolutionary  
595 statistics computed by zol when different approaches to dereplication are used. See Methods  
596 for further details. \*To simulate no dereplication, gene cluster dereplication with re-inflation  
597 parameters were used in zol.

598  
599 **Figure 7: Distribution of the *epa* locus and associated genes across the genus of**  
600 ***Enterococcus*.** A) A schematic is shown for the *epa* locus in *E. faecalis* for genes which were  
601 found in  $\geq 25\%$  of 83 representative genomes for the species presented in consensus order with  
602 consensus directionality as inferred by zol. The coloring corresponds to the conservation of  
603 individual genes. Genes upstream and/or including *epaR* were recently proposed to be involved  
604 in decoration of Epa by Guerardel *et al.* 2020. "/" indicates that the ortholog group was not  
605 single-copy in the context of the gene cluster. The tracks below the gene showcase their  
606 sequence similarity across the *E. faecalis* genomes measured using (B) Tajima's D and (C) the  
607 average sequence alignment entropy. D) The major allele frequency is depicted across the  
608 alignment for the ortholog group featuring *epaX*. Sites predicted to be under negative selection  
609 by FUBAR,  $\text{Prob}(\square > \square) \geq 0.9$ , are marked in red. E) An approximate maximum-likelihood  
610 phylogeny based on gap-filtered codon alignments for the ortholog group corresponding to *epaX*  
611 and *epaX*-like proteins in the joint *E. faecalis* and *E. faecium* investigation of the *epa* locus using  
612 zol. F) Conservation of *epaX* is shown amongst *E. faecalis* and *E. faecium* genomes with a  
613 high-quality representation of the *epa* locus available. Coloring of the bars corresponds to the  
614 proportion of genomes with a certain copy-count of the *epaX*-like ortholog group. G) The  
615 distribution of the *epa* locus and associated genes, based on high-sensitivity criteria used for  
616 running fai, is shown across 463 representative genomes across *Enterococcus*. Coloring of the  
617 heatmap corresponds to the normalized bitscore of the best alignment to coding sequences  
618 from *E. faecalis* V583.

619  
620 **Figure S1: Example illustrations for assessing quality of homologous gene clusters**  
621 **produced by fai.** A) Gene calling or frame-shift differences between the query gene cluster and  
622 coding-sequence predictions in the target genome have resulted in a discrepancy for OG\_1  
623 (highlighted) from the query being regarded as two separate coding-sequences in the target  
624 genome. B) Three candidate gene cluster segments located near scaffold edges which match  
625 the query gene cluster and meet the thresholds needed for detection as requested in fai in  
626 aggregate.

627

628 **Figure S2: Conservation in the upstream regions of coding sequences of genes in the**  
629 **leporin BGC.** The average entropy of the 100 bp upstream regions is shown for each of the  
630 genes from the leporin BGC. Coloring of the bars corresponds to effects on BGC expression (for  
631 *lepE* and *lepB*) or metabolite production (using a mutant with overexpression of *lepE*) when  
632 genes were knocked out as determined by Cary *et al.* 2015.

633

634 **Figure S3: Influence of dereplication on evolutionary statistics computed by zol.** The  
635 relationship in values for analogous ortholog groups which map to query proteins from *E.*  
636 *faecalis* V583 for different evolutionary statistics (**A-I**) when different sets of gene clusters  
637 corresponding to different approaches in dereplication are shown. Only ortholog groups which  
638 lacked any paralogous proteins are shown and accounted for. A line is shown in each plot  
639 corresponding to a 1:1 ratio.

640

641 **Figure S4: The ortholog group with *epaX* features greater diversity relative to other**  
642 **glycosyl transferase related ortholog groups from the *epa* locus in *E. faecalis*.** An  
643 approximate maximum-likelihood phylogeny based on gap-filtered protein alignments of  
644 ortholog groups with domains featuring “glycosyl” and “transferase” as key words. Ortholog  
645 groupings (coloring of phylogeny branches) by zol were largely consistent with phylogenetic  
646 clades. Association of clades to genes from *E. faecalis* V583 based on sequence alignment are  
647 noted.

648

649

## 650 **References**

651

- 652 1. Snyder, L., Henkin, T. M., Peters, J. E. & Champness, W. Molecular Genetics of Bacteria,  
653 4th Edition. Preprint at <https://doi.org/10.1128/9781555817169> (2013).
- 654 2. Price, M. N., Arkin, A. P. & Alm, E. J. The life-cycle of operons. *PLoS Genet.* **2**, e96 (2006).
- 655 3. Ptashne, M. *A genetic switch: Gene control and phage. lambda.* (Palo Alto, CA (US);  
656 Blackwell Scientific Publications, 1986).
- 657 4. Andreu, V. P. *et al.* gutSMASH predicts specialized primary metabolic pathways from the  
658 human gut microbiota. *Nature Biotechnology* Preprint at [https://doi.org/10.1038/s41587-](https://doi.org/10.1038/s41587-023-01675-1)  
659 [023-01675-1](https://doi.org/10.1038/s41587-023-01675-1) (2023).
- 660 5. Fischbach, M. A., Walsh, C. T. & Clardy, J. The evolution of gene collectives: How natural  
661 selection drives chemical innovation. *Proc. Natl. Acad. Sci. U. S. A.* **105**, 4601–4608  
662 (2008).
- 663 6. Gal-Mor, O. & Finlay, B. B. Pathogenicity islands: a molecular toolbox for bacterial

- 664 virulence. *Cell. Microbiol.* **8**, 1707–1719 (2006).
- 665 7. Kaper, J. B., Nataro, J. P. & Mobley, H. L. Pathogenic *Escherichia coli*. *Nat. Rev. Microbiol.*  
666 **2**, 123–140 (2004).
- 667 8. Bolwell, G. P. & Paul Bolwell, G. *Biochemistry & Molecular Biology of Plants.*  
668 *Phytochemistry* vol. 58 185 Preprint at [https://doi.org/10.1016/s0031-9422\(01\)00095-4](https://doi.org/10.1016/s0031-9422(01)00095-4)  
669 (2001).
- 670 9. Lindahl, L. & Zengel, J. M. Operon-specific regulation of ribosomal protein synthesis in  
671 *Escherichia coli*. *Proc. Natl. Acad. Sci. U. S. A.* **76**, 6542–6546 (1979).
- 672 10. Cordero, O. X. & Polz, M. F. Explaining microbial genomic diversity in light of evolutionary  
673 ecology. *Nat. Rev. Microbiol.* **12**, 263–273 (2014).
- 674 11. Salamzade, R. *et al.* IsaBGC provides a comprehensive framework for evolutionary  
675 analysis of biosynthetic gene clusters within focal taxa. *bioRxiv* 2022.04.20.488953 (2022)  
676 doi:10.1101/2022.04.20.488953.
- 677 12. Drott, M. T. *et al.* Microevolution in the pansecondary metabolome of *Aspergillus flavus* and  
678 its potential macroevolutionary implications for filamentous fungi. *Proc. Natl. Acad. Sci. U.*  
679 *S. A.* **118**, (2021).
- 680 13. Tatusov, R. L., Galperin, M. Y., Natale, D. A. & Koonin, E. V. The COG database: a tool for  
681 genome-scale analysis of protein functions and evolution. *Nucleic Acids Res.* **28**, 33–36  
682 (2000).
- 683 14. Remm, M., Storm, C. E. & Sonnhammer, E. L. Automatic clustering of orthologs and in-  
684 paralogs from pairwise species comparisons. *J. Mol. Biol.* **314**, 1041–1052 (2001).
- 685 15. Li, L., Stoeckert, C. J., Jr & Roos, D. S. OrthoMCL: identification of ortholog groups for  
686 eukaryotic genomes. *Genome Res.* **13**, 2178–2189 (2003).
- 687 16. Edwards, D. J. & Holt, K. E. Beginner's guide to comparative bacterial genome analysis  
688 using next-generation sequence data. *Microb. Inform. Exp.* **3**, 2 (2013).
- 689 17. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics*



- 690           **31**, 3691–3693 (2015).
- 691   18. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative  
692           genomics. *Genome Biol.* **20**, 238 (2019).
- 693   19. Medema, M. H., Takano, E. & Breitling, R. Detecting sequence homology at the gene  
694           cluster level with MultiGeneBlast. *Mol. Biol. Evol.* **30**, 1218–1223 (2013).
- 695   20. Abby, S. S., Néron, B., Ménager, H., Touchon, M. & Rocha, E. P. C. MacSyFinder: a  
696           program to mine genomes for molecular systems with an application to CRISPR-Cas  
697           systems. *PLoS One* **9**, e110726 (2014).
- 698   21. Gilchrist, C. L. M. *et al.* Cblaster: A remote search tool for rapid identification and  
699           visualization of homologous gene clusters. *Bioinformatics Advances* **1**, (2021).
- 700   22. Navarro-Muñoz, J. C. *et al.* A computational framework to explore large-scale biosynthetic  
701           diversity. *Nat. Chem. Biol.* **16**, 60–68 (2020).
- 702   23. Georgescu, C. H. *et al.* SynerClust: a highly scalable, synteny-aware orthologue clustering  
703           tool. *Microb Genom* **4**, (2018).
- 704   24. Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo  
705           pipeline. *Genome Biol.* **21**, 180 (2020).
- 706   25. Cosentino, S. & Iwasaki, W. SonicParanoid: fast, accurate and easy orthology inference.  
707           *Bioinformatics* **35**, 149–151 (2019).
- 708   26. Hu, X. & Friedberg, I. SwiftOrtho: A fast, memory-efficient, multiple genome orthology  
709           classifier. *Gigascience* **8**, (2019).
- 710   27. Vallenet, D. *et al.* MaGe: a microbial genome annotation system supported by synteny  
711           results. *Nucleic Acids Res.* **34**, 53–65 (2006).
- 712   28. Stam, M. *et al.* NetSyn: genomic context exploration of protein families. *bioRxiv* (2023)  
713           doi:10.1101/2023.02.15.528638.
- 714   29. Salamzade, R. *et al.* Evolutionary investigations of the biosynthetic diversity in the skin  
715           microbiome using IsaBGC. *Microb Genom* **9**, (2023).

- 716 30. Liu, M. *et al.* ICEberg 2.0: an updated database of bacterial integrative and conjugative  
717 elements. *Nucleic Acids Res.* **47**, D660–D665 (2019).
- 718 31. Terlouw, B. R. *et al.* MIBiG 3.0: a community-driven effort to annotate experimentally  
719 validated biosynthetic gene clusters. *Nucleic Acids Res.* **51**, D603–D610 (2023).
- 720 32. Bertelli, C. *et al.* IslandViewer 4: expanded prediction of genomic islands for larger-scale  
721 datasets. *Nucleic Acids Res.* **45**, W30–W35 (2017).
- 722 33. Larralde, M. Pyrodigal: Python bindings and interface to Prodigal, an efficient method for  
723 gene prediction in prokaryotes. *J. Open Source Softw.* **7**, 4296 (2022).
- 724 34. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site  
725 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 726 35. Li, H. Protein-to-genome alignment with miniprot. *Bioinformatics* **39**, (2023).
- 727 36. Graziotin, A. L., Koonin, E. V. & Kristensen, D. M. Prokaryotic Virus Orthologous Groups  
728 (pVOGs): a resource for comparative genomics and protein family annotation. *Nucleic*  
729 *Acids Res.* **45**, D491–D498 (2017).
- 730 37. Liu, B., Zheng, D., Jin, Q., Chen, L. & Yang, J. VFDB 2019: a comparative pathogenomic  
731 platform with an interactive web interface. *Nucleic Acids Res.* **47**, D687–D692 (2019).
- 732 38. Siguier, P., Perochon, J., Lestrade, L., Mahillon, J. & Chandler, M. ISfinder: the reference  
733 centre for bacterial insertion sequences. *Nucleic Acids Res.* **34**, D32–6 (2006).
- 734 39. Alcock, B. P. *et al.* CARD 2023: expanded curation, support for machine learning, and  
735 resistome prediction at the Comprehensive Antibiotic Resistance Database. *Nucleic Acids*  
736 *Res.* **51**, D690–D699 (2023).
- 737 40. Kosakovsky Pond, S. L. *et al.* HyPhy 2.5—A Customizable Platform for Evolutionary  
738 Hypothesis Testing Using Phylogenies. *Mol. Biol. Evol.* **37**, 295–299 (2019).
- 739 41. Thorpe, H. A., Bayliss, S. C., Sheppard, S. K. & Feil, E. J. Piggy: a rapid, large-scale pan-  
740 genome analysis tool for intergenic regions in bacteria. *Gigascience* **7**, 1–11 (2018).
- 741 42. Gilchrist, C. L. M. & Chooi, Y.-H. clinker & clustermap.js: automatic generation of gene

- 742 cluster comparison figures. *Bioinformatics* **37**, 2473–2475 (2021).
- 743 43. Hackl, T., Duponchel, S., Barenhoff, K., Weinmann, A. & Fischer, M. G. Virophages and  
744 retrotransposons colonize the genomes of a heterotrophic flagellate. *Elife* **10**, (2021).
- 745 44. Sullivan, M. J., Petty, N. K. & Beatson, S. A. Easyfig: a genome comparison visualizer.  
746 *Bioinformatics* **27**, 1009–1010 (2011).
- 747 45. Shaw, J. & Yu, Y. W. Fast and robust metagenomic sequence comparison through sparse  
748 chaining with skani. *bioRxiv* 2023.01.18.524587 (2023) doi:10.1101/2023.01.18.524587.
- 749 46. Blackwell, G. *et al.* Exploring bacterial diversity via a curated and searchable snapshot of  
750 archived DNA. *Access Microbiol.* **4**, (2022).
- 751 47. Lebreton, F. *et al.* Emergence of epidemic multidrug-resistant *Enterococcus faecium* from  
752 animal and commensal strains. *MBio* **4**, (2013).
- 753 48. Lieberman, T. D. *et al.* Genetic variation of a bacterial pathogen within individuals with  
754 cystic fibrosis provides a record of selective pressures. *Nat. Genet.* **46**, 82–87 (2014).
- 755 49. Crits-Christoph, A., Olm, M. R., Diamond, S., Bouma-Gregson, K. & Banfield, J. F. Soil  
756 bacterial populations are shaped by recombination and gene-specific selection across a  
757 grassland meadow. *ISME J.* **14**, 1834–1846 (2020).
- 758 50. Hudson, R. R., Slatkin, M. & Maddison, W. P. Estimation of levels of gene flow from DNA  
759 sequence data. *Genetics* **132**, 583–589 (1992).
- 760 51. Tran, P. Q. & Anantharaman, K. Biogeochemistry Goes Viral: towards a Multifaceted  
761 Approach To Study Viruses and Biogeochemical Cycling. *mSystems* **6**, e0113821 (2021).
- 762 52. Barr, J. J. *et al.* Bacteriophage adhering to mucus provide a non-host-derived immunity.  
763 *Proc. Natl. Acad. Sci. U. S. A.* **110**, 10771–10776 (2013).
- 764 53. Lefeuvre, P. *et al.* Evolution and ecology of plant viruses. *Nat. Rev. Microbiol.* **17**, 632–644  
765 (2019).
- 766 54. Tran, P. Q. *et al.* Viral impacts on microbial activity and biogeochemical cycling in a  
767 seasonally anoxic freshwater lake. *bioRxiv* 2023.04.19.537559 (2023)

- 768 doi:10.1101/2023.04.19.537559.
- 769 55. Kieft, K., Zhou, Z. & Anantharaman, K. VIBRANT: automated recovery, annotation and  
770 curation of microbial viruses, and evaluation of viral community function from genomic  
771 sequences. *Microbiome* **8**, 90 (2020).
- 772 56. Bok, J. W. *et al.* Genomic mining for *Aspergillus* natural products. *Chem. Biol.* **13**, 31–37  
773 (2006).
- 774 57. Vadlapudi, V. *et al.* *Aspergillus* Secondary Metabolite Database, a resource to understand  
775 the Secondary metabolome of *Aspergillus* genus. *Sci. Rep.* **7**, 7325 (2017).
- 776 58. Robey, M. T., Caesar, L. K., Drott, M. T., Keller, N. P. & Kelleher, N. L. An interpreted atlas  
777 of biosynthetic gene clusters from 1,000 fungal genomes. *Proc. Natl. Acad. Sci. U. S. A.*  
778 **118**, (2021).
- 779 59. Hatmaker, E. A. *et al.* Genomic and Phenotypic Trait Variation of the Opportunistic Human  
780 Pathogen *Aspergillus flavus* and Its Close Relatives. *Microbiol Spectr* **10**, e0306922 (2022).
- 781 60. Xie, H. *et al.* Global multi-omics profiling reveals evolutionary drivers of phylogeographic  
782 diversity of fungal specialized metabolism. (2023) doi:10.21203/rs.3.rs-2471999/v1.
- 783 61. Cary, J. W. *et al.* An *Aspergillus flavus* secondary metabolic gene cluster containing a  
784 hybrid PKS-NRPS is necessary for synthesis of the 2-pyridones, leporins. *Fungal Genet.*  
785 *Biol.* **81**, 88–97 (2015).
- 786 62. Drăgan, M.-A., Moghul, I., Priyam, A., Bustos, C. & Wurm, Y. GeneValidator: identify  
787 problems with protein-coding gene predictions. *Bioinformatics* **32**, 1559–1561 (2016).
- 788 63. Scalzitti, N., Jeannin-Girardon, A., Collet, P., Poch, O. & Thompson, J. D. A benchmark  
789 study of ab initio gene prediction methods in diverse eukaryotic organisms. *BMC Genomics*  
790 **21**, 293 (2020).
- 791 64. Skerker, J. M. *et al.* Chromosome assembled and annotated genome sequence of  
792 *Aspergillus flavus* NRRL 3357. *G3* **11**, jkab213 (2021).
- 793 65. Yang, K., Tian, J. & Keller, N. P. Post-translational modifications drive secondary

- 794 metabolite biosynthesis in *Aspergillus*: a review. *Environ. Microbiol.* **24**, 2857–2881 (2022).
- 795 66. Klich, M. A. *Aspergillus flavus*: the major producer of aflatoxin. *Mol. Plant Pathol.* **8**, 713–  
796 722 (2007).
- 797 67. Cary, J. W., Ehrlich, K. C., Bland, J. M. & Montalbano, B. G. The Aflatoxin Biosynthesis  
798 Cluster Gene, *afIX*, Encodes an Oxidoreductase Involved in Conversion of Versicolorin A  
799 to Demethylsterigmatocystin. *Applied and Environmental Microbiology* vol. 72 1096–1101  
800 Preprint at <https://doi.org/10.1128/aem.72.2.1096-1101.2006> (2006).
- 801 68. Cleveland, T. E. *et al.* Potential of *Aspergillus flavus* genomics for applications in  
802 biotechnology. *Trends Biotechnol.* **27**, 151–157 (2009).
- 803 69. Ehrlich, K. C., Li, P., Scharfenstein, L. & Chang, P.-K. HypC, the anthrone oxidase involved  
804 in aflatoxin biosynthesis. *Appl. Environ. Microbiol.* **76**, 3374–3377 (2010).
- 805 70. Xu, Y., Murray, B. E. & Weinstock, G. M. A cluster of genes involved in polysaccharide  
806 biosynthesis from *Enterococcus faecalis* OG1RF. *Infect. Immun.* **66**, 4313–4323 (1998).
- 807 71. Palmer, K. L. *et al.* Comparative Genomics of Enterococci: Variation in *Enterococcus*  
808 *faecalis*, Clade Structure in *E. faecium*, and Defining Characteristics of *E. gallinarum* and *E.*  
809 *casseliflavus*. *mBio* vol. 3 Preprint at <https://doi.org/10.1128/mbio.00318-11> (2012).
- 810 72. Hancock, L. E., Murray, B. E. & Sillanpää, J. Enterococcal Cell Wall Components and  
811 Structures. in *Enterococci: From Commensals to Leading Causes of Drug Resistant*  
812 *Infection* (eds. Gilmore, M. S., Clewell, D. B., Ike, Y. & Shankar, N.) (Massachusetts Eye  
813 and Ear Infirmary, 2014).
- 814 73. Qin, X. *et al.* Complete genome sequence of *Enterococcus faecium* strain TX16 and  
815 comparative genomic analysis of *Enterococcus faecium* genomes. *BMC Microbiol.* **12**, 135  
816 (2012).
- 817 74. Teng, F., Jacques-Palaz, K. D., Weinstock, G. M. & Murray, B. E. Evidence that the  
818 Enterococcal Polysaccharide Antigen Gene ( *epa* ) Cluster Is Widespread in *Enterococcus*  
819 *faecalis* and Influences Resistance to Phagocytic Killing of *E. faecalis*. *Infection and*

- 820 *Immunity* vol. 70 2010–2015 Preprint at <https://doi.org/10.1128/iai.70.4.2010-2015.2002>  
821 (2002).
- 822 75. Teng, F., Singh, K. V., Bourgoigne, A., Zeng, J. & Murray, B. E. Further Characterization of  
823 the *epa* Gene Cluster and Epa Polysaccharides of *Enterococcus faecalis*. *Infection and*  
824 *Immunity* vol. 77 3759–3767 Preprint at <https://doi.org/10.1128/iai.00149-09> (2009).
- 825 76. Rigottier-Gois, L. *et al.* The surface rhamnopolysaccharide *epa* of *Enterococcus faecalis* is  
826 a key determinant of intestinal colonization. *J. Infect. Dis.* **211**, 62–71 (2015).
- 827 77. Guerardel, Y. *et al.* Complete structure of the enterococcal polysaccharide antigen (EPA) of  
828 vancomycin-resistant *Enterococcus faecalis* V583 reveals that EPA decorations are  
829 teichoic acids covalently linked to a rhamnopolysaccharide backbone. *MBio* **11**, (2020).
- 830 78. Smith, R. E. *et al.* Decoration of the enterococcal polysaccharide antigen EPA is essential  
831 for virulence, cell surface charge and interaction with effectors of the innate immune  
832 system. *PLoS Pathog.* **15**, e1007730 (2019).
- 833 79. Singh, K. V. & Murray, B. E. Loss of a Major Enterococcal Polysaccharide Antigen (Epa) by  
834 *Enterococcus faecalis* Is Associated with Increased Resistance to Ceftriaxone and  
835 Carbapenems. *Antimicrob. Agents Chemother.* **63**, (2019).
- 836 80. Ho, K., Huo, W., Pas, S., Dao, R. & Palmer, K. L. Loss-of-Function Mutations in *epaR*  
837 Confer Resistance to  $\square$ NPV1 Infection in *Enterococcus faecalis* OG1RF. *Antimicrobial*  
838 *Agents and Chemotherapy* vol. 62 Preprint at <https://doi.org/10.1128/aac.00758-18> (2018).
- 839 81. Fiore, E., Van Tyne, D. & Gilmore, M. S. Pathogenicity of Enterococci. *Microbiol Spectr* **7**,  
840 (2019).
- 841 82. Lebreton, F., Willems, R. J. L. & Gilmore, M. S. *Enterococcus* Diversity, Origins in Nature,  
842 and Gut Colonization. in *Enterococci: From Commensals to Leading Causes of Drug*  
843 *Resistant Infection* (eds. Gilmore, M. S., Clewell, D. B., Ike, Y. & Shankar, N.)  
844 (Massachusetts Eye and Ear Infirmary, 2014).
- 845 83. Lebreton, F. *et al.* Tracing the Enterococci from Paleozoic Origins to the Hospital. *Cell* **169**,

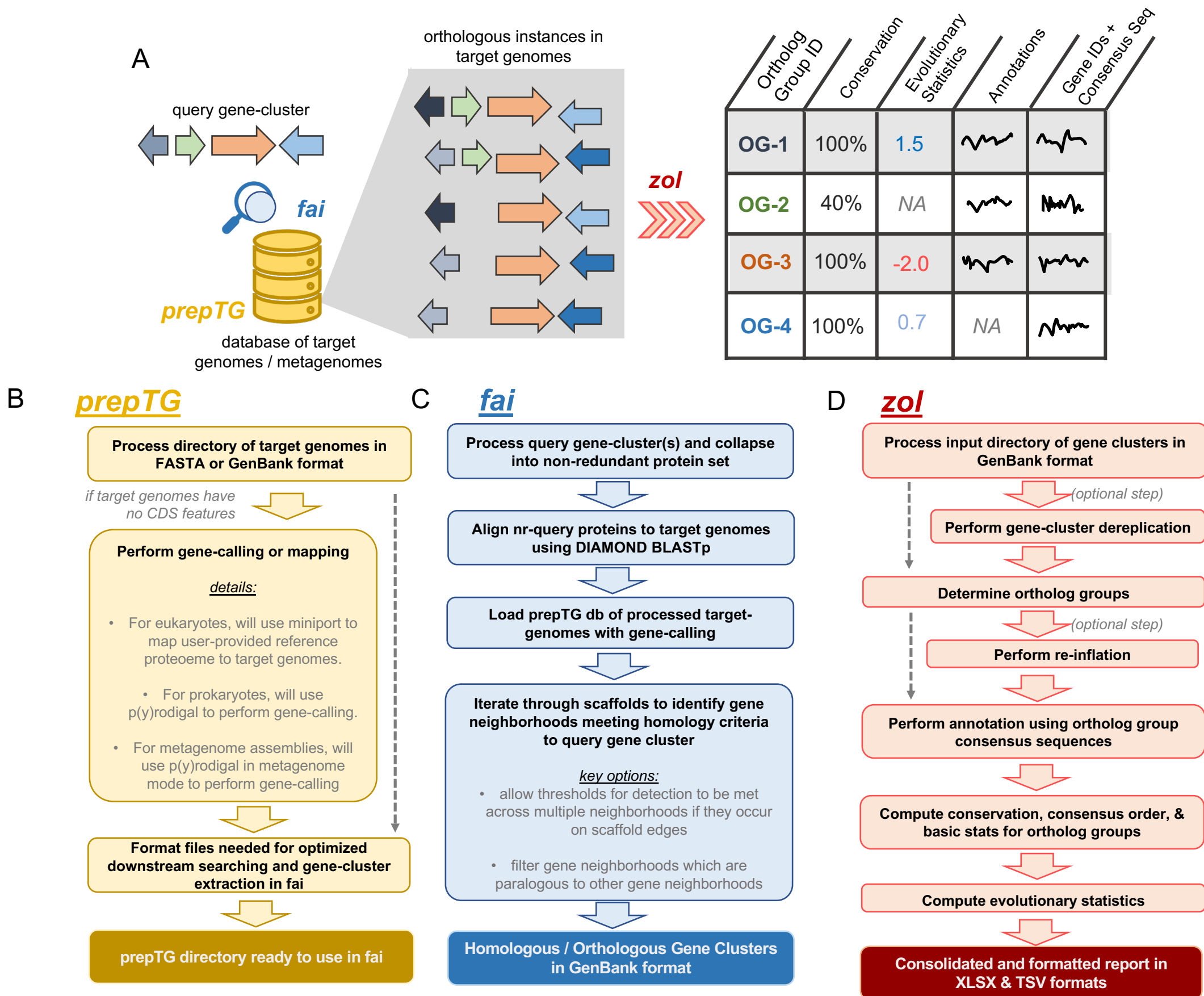
- 846 849-861.e13 (2017).
- 847 84. Schwartzman, J. A. *et al.* Global diversity of enterococci and description of 18 novel  
848 species. *bioRxiv* 2023.05.18.540996 (2023) doi:10.1101/2023.05.18.540996.
- 849 85. Parks, D. H. *et al.* GTDB: an ongoing census of bacterial and archaeal diversity through a  
850 phylogenetically consistent, rank normalized and complete genome-based taxonomy.  
851 *Nucleic Acids Res.* (2021) doi:10.1093/nar/gkab776.
- 852 86. Reid, K. C., Cockerill, F. R., III & Patel, R. Clinical and epidemiological features of  
853 *Enterococcus casseliflavus/flavescens* and *Enterococcus gallinarum* bacteremia: a report of  
854 20 cases. *Clin. Infect. Dis.* **32**, 1540–1546 (2001).
- 855 87. Monticelli, J., Knezevich, A., Luzzati, R. & Di Bella, S. Clinical management of non-faecium  
856 non-faecalis vancomycin-resistant enterococci infection. Focus on *Enterococcus gallinarum*  
857 and *Enterococcus casseliflavus/flavescens*. *J. Infect. Chemother.* **24**, 237–246 (2018).
- 858 88. Nevers, Y. *et al.* The Quest for Orthologs orthology benchmark service in 2022. *Nucleic*  
859 *Acids Res.* (2022).
- 860 89. Dale, J. L., Cagnazzo, J., Phan, C. Q., Barnes, A. M. T. & Dunny, G. M. Multiple Roles for  
861 *Enterococcus faecalis* Glycosyltransferases in Biofilm-Associated Antibiotic Resistance,  
862 Cell Envelope Integrity, and Conjugative Transfer. *Antimicrobial Agents and Chemotherapy*  
863 vol. 59 4094–4105 Preprint at <https://doi.org/10.1128/aac.00344-15> (2015).
- 864 90. Olm, M. R., Brown, C. T., Brooks, B. & Banfield, J. F. dRep: a tool for fast and accurate  
865 genomic comparisons that enables improved genome recovery from metagenomes through  
866 de-replication. *ISME J.* **11**, 2864–2868 (2017).
- 867 91. Jain, C., Rodriguez-R, L. M., Phillippy, A. M., Konstantinidis, K. T. & Aluru, S. High  
868 throughput ANI analysis of 90K prokaryotic genomes reveals clear species boundaries.  
869 *Nat. Commun.* **9**, 1–8 (2018).
- 870 92. van den Belt, M. *et al.* CAGECAT: The CompArative GEne Cluster Analysis Toolbox for  
871 rapid search and visualisation of homologous gene clusters. *BMC Bioinformatics* **24**, 181

- 872 (2023).
- 873 93. Bertelli, C. *et al.* Enabling genomic island prediction and comparison in multiple genomes to  
874 investigate bacterial evolution and outbreaks. *Microb. Genom.* **8**, (2022).
- 875 94. Zablocki, O., Jang, H. B., Bolduc, B. & Sullivan, M. B. VConTACT 2: A tool to automate  
876 genome-based prokaryotic viral taxonomy. in *Plant and Animal Genome XXVII*  
877 *Conference* (January 12- 16, 2019) (PAG, 2019).
- 878 95. Salamzade, R. *et al.* Inter-species geographic signatures for tracing horizontal gene  
879 transfer and long-term persistence of carbapenem resistance. *Genome Med.* **14**, 37 (2022).
- 880 96. Sheppard, A. E. *et al.* Nested Russian Doll-Like Genetic Mobility Drives Rapid  
881 Dissemination of the Carbapenem Resistance Gene blaKPC. *Antimicrob. Agents*  
882 *Chemother.* **60**, 3767–3778 (2016).
- 883 97. Groussin, M. *et al.* Elevated rates of horizontal gene transfer in the industrialized human  
884 microbiome. *Cell* **184**, 2053-2067.e18 (2021).
- 885 98. Crits-Christoph, A., Diamond, S., Butterfield, C. N., Thomas, B. C. & Banfield, J. F. Novel  
886 soil bacteria possess diverse genes for secondary metabolite biosynthesis. *Nature* **558**,  
887 440–444 (2018).
- 888 99. Bickhart, D. M. *et al.* Generation of lineage-resolved complete metagenome-assembled  
889 genomes by precision phasing. *bioRxiv* 2021.05.04.442591 (2021)  
890 doi:10.1101/2021.05.04.442591.
- 891 100. Chatterjee, A. *et al.* Bacteriophage Resistance Alters Antibiotic-Mediated Intestinal  
892 Expansion of Enterococci. *Infect. Immun.* **87**, (2019).
- 893 101. Chatterjee, A. *et al.* Parallel genomics uncover novel enterococcal-bacteriophage  
894 interactions. Preprint at <https://doi.org/10.1101/858506>.
- 895 102. Canfield, G. S. *et al.* Lytic bacteriophages facilitate antibiotic sensitization of *Enterococcus*  
896 *faecium*. Preprint at <https://doi.org/10.1101/2020.09.22.309401>.
- 897 103. van Tilburg Bernardes, E., Charron-Mazenod, L., Reading, D. J., Reckseidler-Zenteno, S.

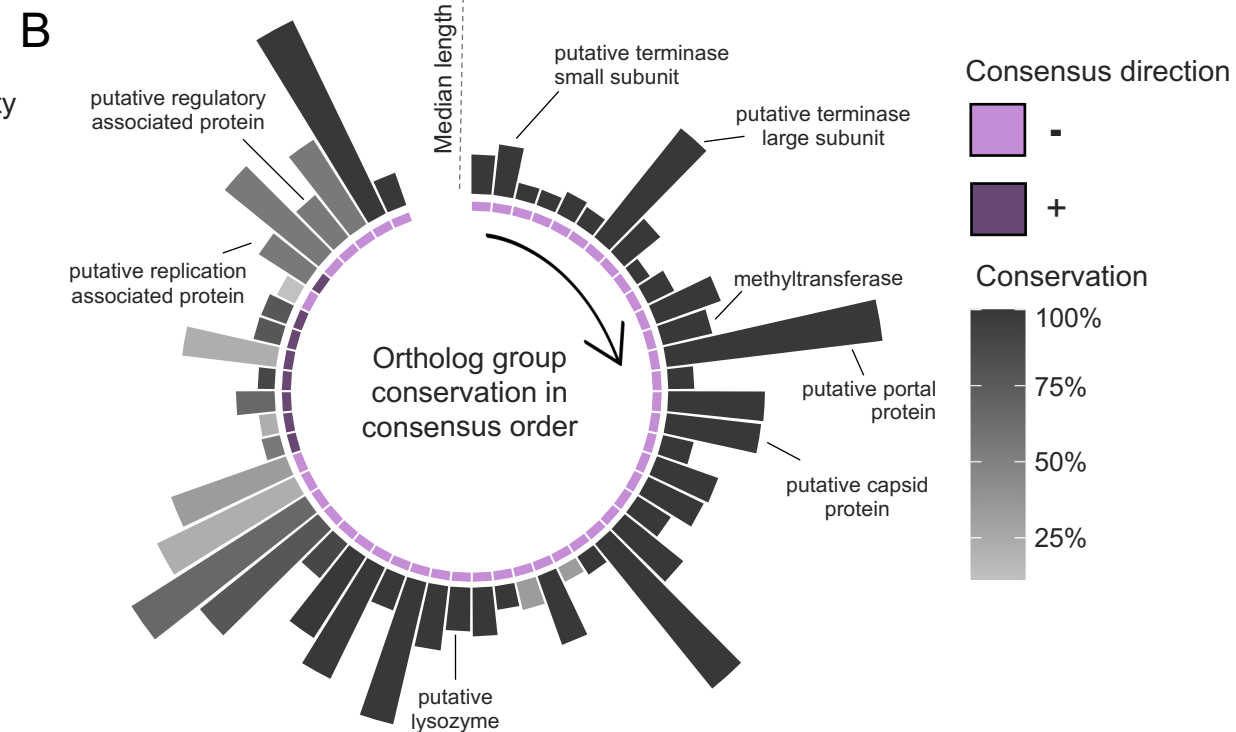
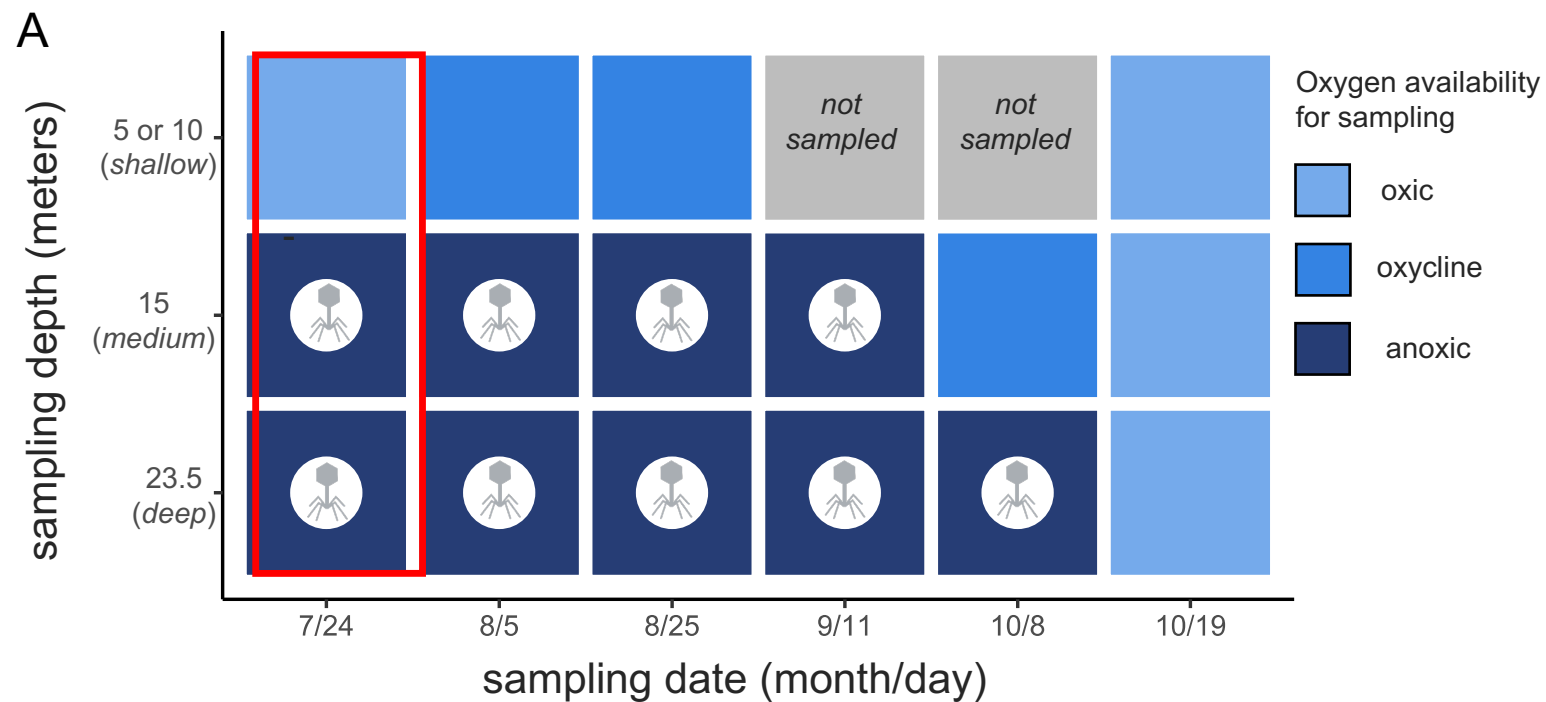


- 898 L. & Lewenza, S. Exopolysaccharide-repressing small molecules with antibiofilm and  
899 antivirulence activity against *Pseudomonas aeruginosa*. *Antimicrob. Agents Chemother.* **61**,  
900 (2017).
- 901 104. Shen, Y. & Loessner, M. J. Beyond antibacterials - exploring bacteriophages as  
902 antivirulence agents. *Curr. Opin. Biotechnol.* **68**, 166–173 (2021).
- 903 105. Shankar-Sinha, S. *et al.* The *Klebsiella pneumoniae* O antigen contributes to bacteremia  
904 and lethality during murine pneumonia. *Infect. Immun.* **72**, 1423–1430 (2004).
- 905 106. Zhao, S. *et al.* Adaptive evolution within gut microbiomes of healthy people. *Cell Host*  
906 *Microbe* **25**, 656-667.e8 (2019).
- 907 107. Grüning, B. *et al.* Bioconda: sustainable and comprehensive software distribution for the life  
908 sciences. *Nat. Methods* **15**, 475–476 (2018).
- 909 108. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular  
910 biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 911 109. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated  
912 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973  
913 (2009).
- 914 110. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood  
915 trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 916 111. Huang, Y., Niu, B., Gao, Y., Fu, L. & Li, W. CD-HIT Suite: a web server for clustering and  
917 comparing biological sequences. *Bioinformatics* **26**, 680–682 (2010).
- 918 112. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 919 113. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND.  
920 *Nat. Methods* **12**, 59–60 (2014).
- 921 114. Schreiber, J. Pomegranate: fast and flexible probabilistic modeling in python. *J. Mach.*  
922 *Learn. Res.* (2017).
- 923 115. Edgar, R. C. Muscle5: High-accuracy alignment ensembles enable unbiased assessments

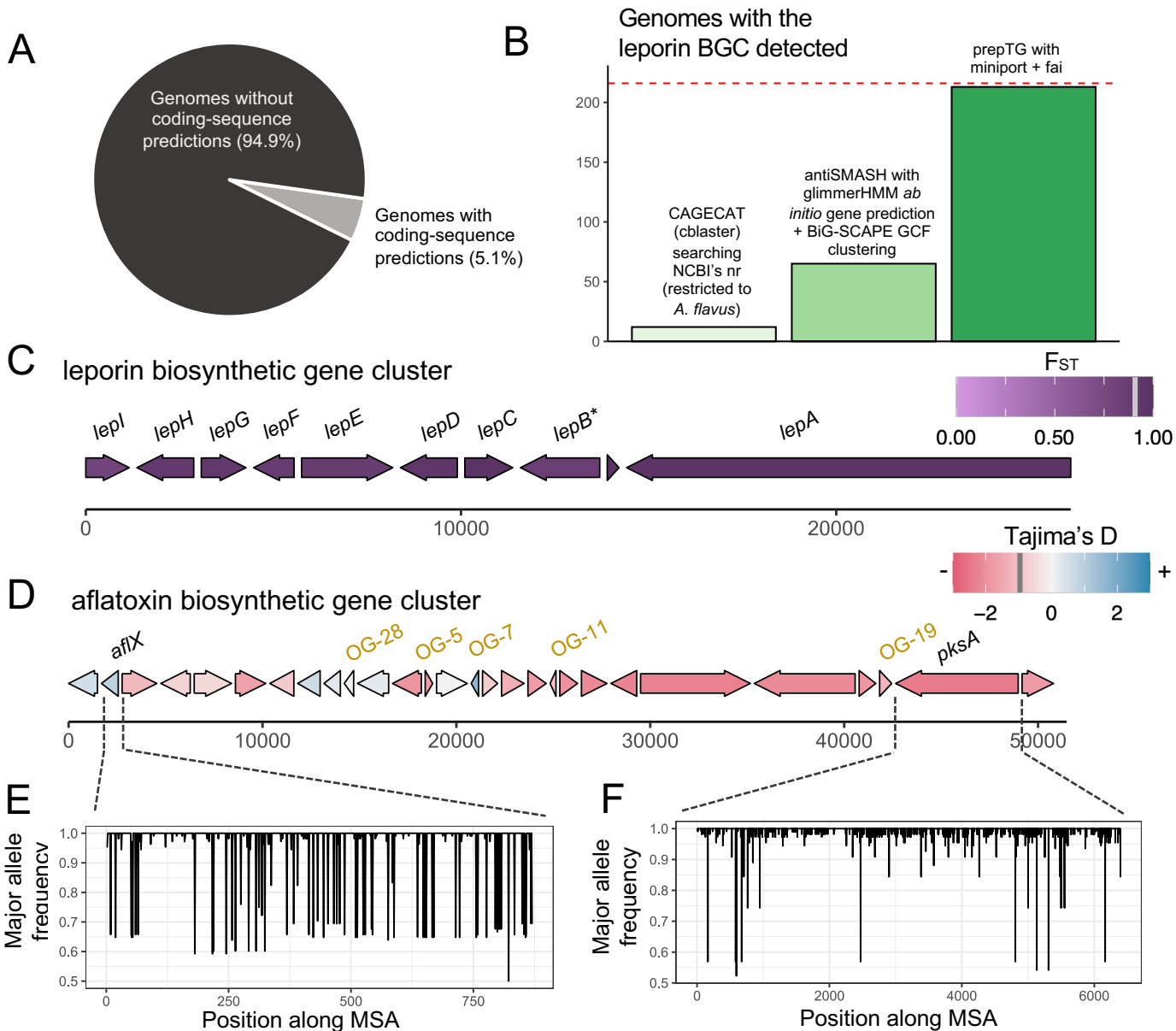
- 924 of sequence homology and phylogeny. *Nat. Commun.* **13**, 6968 (2022).
- 925 116.Blin, K. *et al.* antiSMASH 6.0: improving cluster detection and comparison capabilities.
- 926 *Nucleic Acids Res.* **49**, W29–W35 (2021).
- 927 117.Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using
- 928 MinHash. *Genome Biol.* **17**, 132 (2016).
- 929 118.Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R
- 930 language. *Bioinformatics* **20**, 289–290 (2004).
- 931 119.Majoros, W. H., Pertea, M. & Salzberg, S. L. TigrScan and GlimmerHMM: two open source
- 932 ab initio eukaryotic gene-finders. *Bioinformatics* **20**, 2878–2879 (2004).
- 933 120.Olm, M. R. *et al.* Consistent Metagenome-Derived Metrics Verify and Delineate Bacterial
- 934 Species Boundaries. *mSystems* **5**, (2020).
- 935 121.Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–
- 936 D419 (2021).
- 937 122.Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new
- 938 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).
- 939



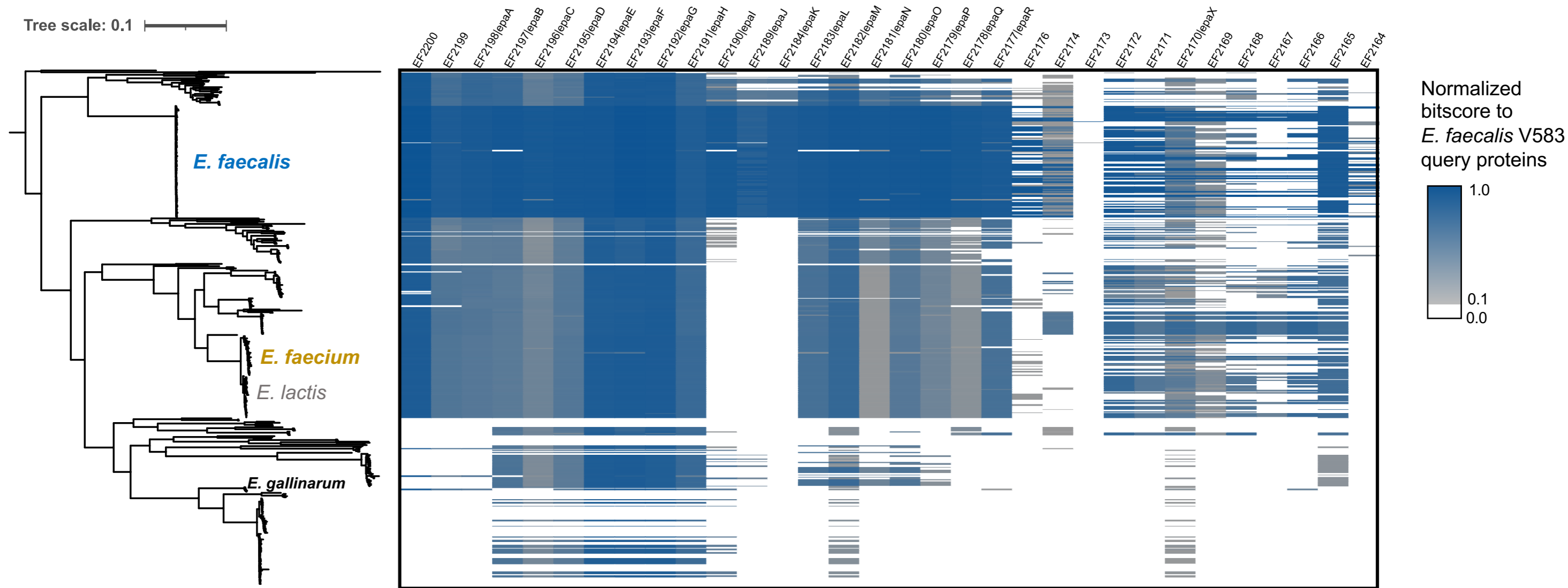
**Figure 1: Overviews of fai and zol.** **A)** A schematic of how prepTG, fai, and zol are integrated to perform evolutionary investigations by searching for gene-clusters. An overview of the prepTG (**B**), fai (**C**) and zol (**D**) algorithms and workflows.



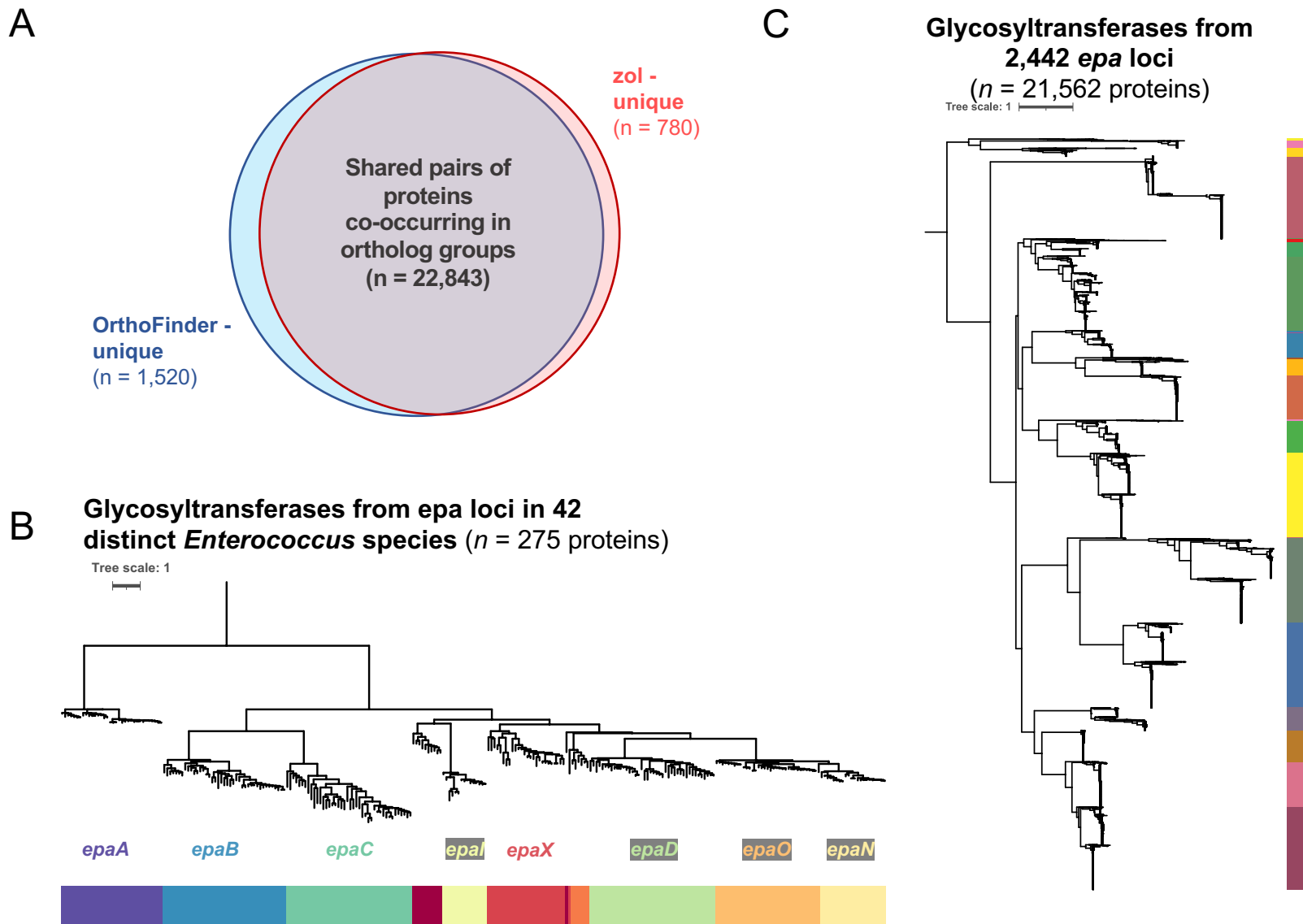
**Figure 2: Targeted viral detection in metagenomes using fai.** **A)** Total metagenomes from a single site in Lake Mendota across multiple depths and timepoints from Tran et al. 2023 were investigated using fai for the presence of a virus found in two of the three earliest microbiome samplings (red box). The presence of the virus is indicated by a phage icon. Metagenome samples are colored according to whether they corresponded to oxic, oxycline, or anoxic. The most shallow sampling depths varied for different dates and consolidated as a single row corresponding to a sampling depth of either 5 or 10 meters. **B)** The pangenome of the virus is shown based on the consensus order and directionality of coding sequences inferred by zol. Bar heights correspond to the median length of coding sequences and are colored based on the percentages of the nine metagenomes the virus was detected in. BioRender was used in generation of this figure.



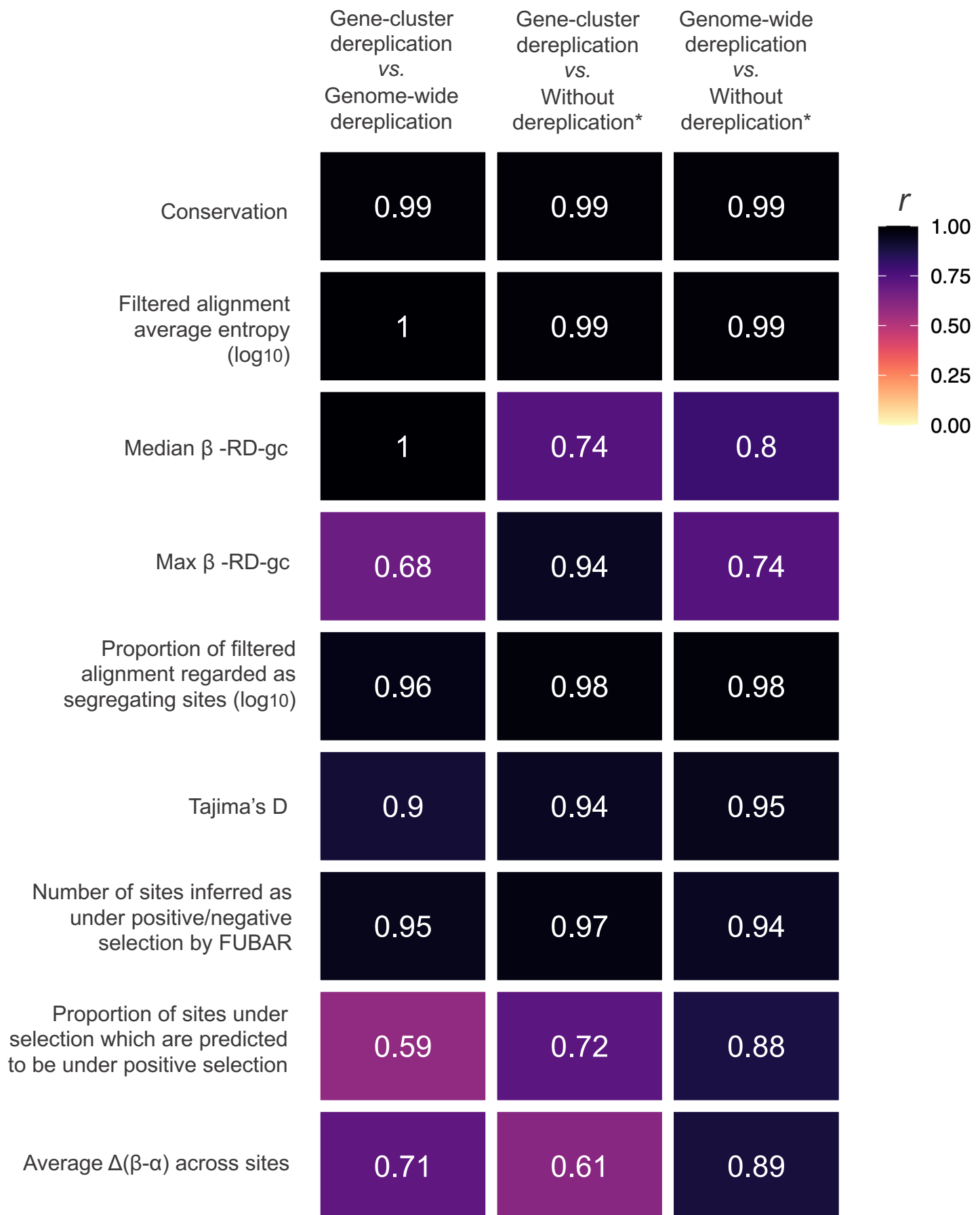
**Figure 3: Evolutionary trends of common BGCs in *A. flavus*.** **A)** The proportion of 216 *A. flavus* genomes from NCBI's GenBank database with coding-sequence predictions available. **B)** Comparison of the sensitivity of *fai* and alternate approaches based on assemblies for detecting the leporin BGC. The red-line indicates the total number of genomes ( $n=216$ ) assessed. A schematic of the **(C)** leporin and **(D)** aflatoxin BGCs is shown with genes present in  $\geq 10\%$  of samples shown in consensus order and relative directionality. Coloring of genes in **(C)** corresponds to  $F_{ST}$  values and in **(D)** to Tajima's  $D$  values, as calculated by *zol*. Grey bars in the legends, at **(C)** 0.92 and **(D)** -0.98, indicate the mean values for the statistics across genes in the BGC. \*For the leporin BGC, *lepB* corresponds to an updated open-reading frame (ORF) prediction by Skerker *et al.* 2021 which was the combination of AFLA\_066860 and AFLA\_066870 ORFs in the MIBiG entry BGC0001445 used as the query for *fai*. For the aflatoxin BGC, ORFs which were not represented in the MIBiG entry BGC0000008 but predicted to be within the aflatoxin BGC by mapping of gene-calls from *A. flavus* NRRL 3357 by Skerker *et al.* 2021 are shown in gold. The major allele frequency distributions are shown for **(E)** *afiX* and **(F)** *pksA*, which depict opposite trends in sequence conservation according to their respective Tajima's  $D$  calculations.



**Figure 4: The *epa* locus is conserved across most enterococcal species.** The distribution of the *epa* locus and associated genes, based on criteria used for running fai, is shown across 463 representative genomes across *Enterococcus*. Coloring of the heatmap corresponds to the normalized bitscore of the best alignment to coding sequences from *E. faecalis* V583.

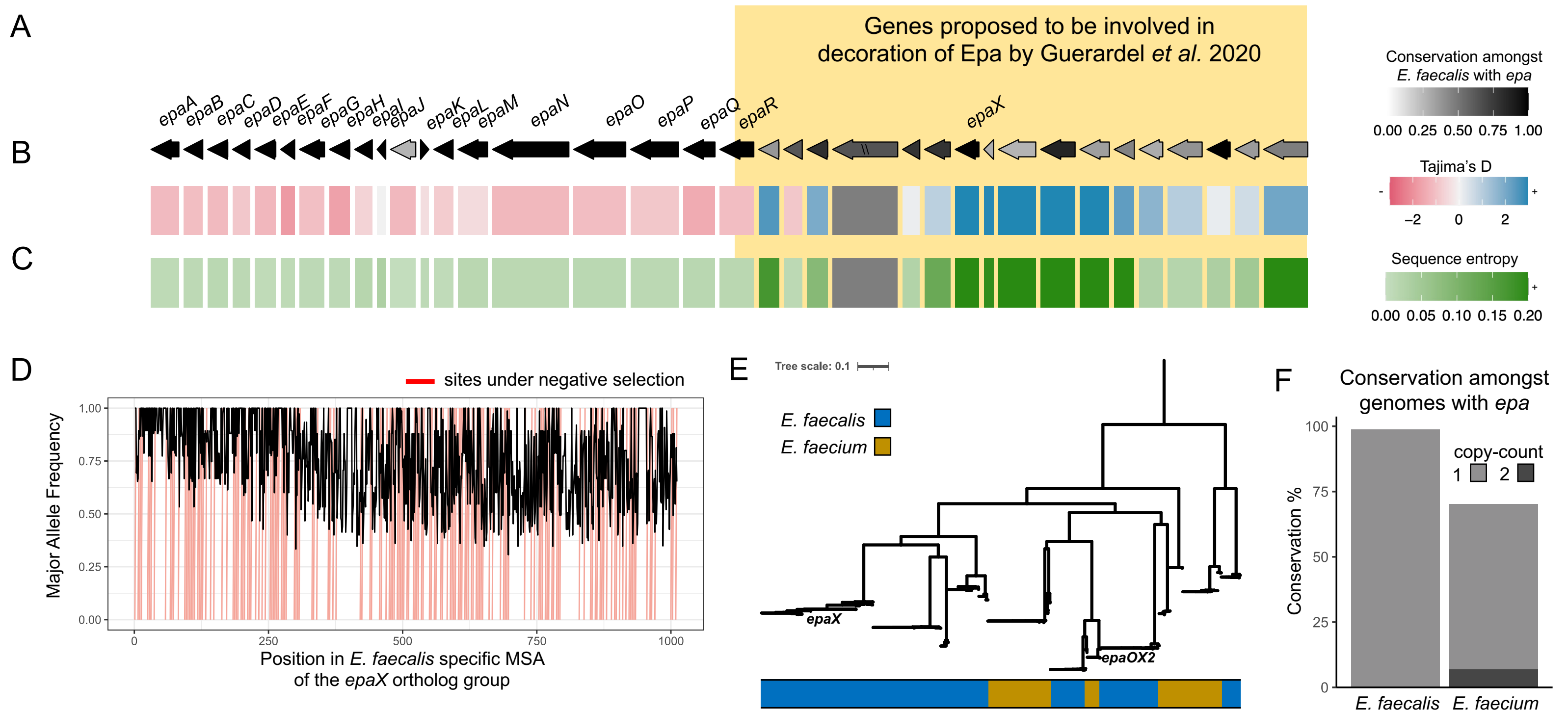


**Figure 5: Assessment of gene-cluster restricted ortholog grouping by *fai* and *zol*.** **A)** *zol* gene-cluster constricted ortholog group predictions for *epa* locus proteins from 42 distinct representative enterococcal species were compared to genome-wide predictions of ortholog groups by OrthoFinder. A phylogeny based on gap-filtered protein alignments of ortholog groups with domains featuring “glycosyl” and “transferase” as key words is shown from **(B)** *epa* loci in the 42 representative genomes and **(C)** a more comprehensive set of 2,442 *epa* loci. Each node represents a specific protein and coloring of the track corresponds to their ortholog group designations by *zol*. Note, (B) 2 (0.07%) and (C) 79 proteins (0.4%) were removed prior to phylogeny construction due to an abundance of gaps in the trimmed alignment.



**Figure 6: Effects of dereplication on the calculation of evolutionary statistics by zol.** The heatmap shows the correlation of values for analogous ortholog groups for various evolutionary statistics computed by zol when different approaches to dereplication are used. See Methods for further details. \*To simulate no dereplication, gene-cluster dereplication with re-inflation parameters were used in zol.





**Figure 7: Distribution of the *epa* locus and associated genes across the genus of *Enterococcus*.** **A)** A schematic is shown for the *epa* locus in *E. faecalis* for genes which were found in  $\geq 25\%$  of 83 representative genomes for the species presented in consensus order with consensus directionality as inferred by zol. The coloring corresponds to the conservation of individual genes. Genes upstream and/or including *epaR* were recently proposed to be involved in decoration of Epa by Guerardel *et al.* 2020. “//” indicates that the ortholog group was not single-copy in the context of the gene-cluster. The tracks below the gene showcase their sequence similarity across the *E. faecalis* genomes measured using **(B)** Tajima’s D and **(C)** the average sequence alignment entropy. **D)** The major allele frequency is depicted across the alignment for the ortholog group featuring *epaX*. Sites predicted to be under negative selection by FUBAR,  $\text{Prob}(\alpha > \beta) \geq 0.9$ , are marked in red. **E)** An approximate maximum-likelihood phylogeny based on gap-filtered codon alignments for the ortholog group corresponding to *epaX* and *epaX*-like proteins in the joint *E. faecalis* and *E. faecium* investigation of the *epa* locus using zol. **F)** Conservation of *epaX* is shown amongst *E. faecalis* and *E. faecium* genomes with a high-quality representation of the *epa* locus available. Coloring of the bars corresponds to the proportion of genomes with a certain copy-count of the *epaX*-like ortholog group.