

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

2

3 Rauf Salamzade^{1,2}, Patricia Q. Tran^{3,4}, Cody Martin^{2,3}, Abigail L. Manson⁵, Michael S.
4 Gilmore^{5,6,7}, Ashlee M. Earl⁵, Karthik Anantharaman³, Lindsay R. Kalan^{1,8,9}

5

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

8 ²Microbiology Doctoral Training Program, University of Wisconsin-Madison, Madison, WI, USA

9 ³Department of Bacteriology, University of Wisconsin-Madison, Madison, WI, USA

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

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

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

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

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

16 ⁹M.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 instances of a query
29 gene cluster of interest amongst a database of target genomes. Subsequently, zol enables
30 reliable, context-specific inference of ortholog groups for individual protein-encoding genes
31 across gene cluster instances. In addition, zol performs functional annotation and computes a
32 variety of evolutionary statistics for each inferred ortholog group. Importantly, in comparison to
33 tools for visual exploration of homologous relationships between gene clusters, zol can scale to
34 thousands of gene cluster instances and produce detailed reports that are easy to digest. To
35 showcase fai and zol, we apply them for: (i) longitudinal tracking of a virus in metagenomes, (ii)
36 discovering novel population-level genetic insights of two common BGCs in the fungal species
37 *Aspergillus flavus*, and (iii) uncovering large-scale evolutionary trends of a virulence-associated
38 gene cluster across thousands of genomes from a diverse bacterial genus.

39

40 **Background**

41

42 *De novo* ortholog grouping typically involves searching for reciprocal best hits of proteins
43 between pairs of genomes, indicative of orthology, and subsequently clustering pairs of inferred
44 orthologs and in-paralogs across multiple genomes¹⁻⁴. Initial methods for orthology inference
45 were designed to be able to identify orthologs between distinct species but limited in the number
46 of genomes they could process¹⁻³. This limitation is largely due to the all-vs-all alignment of

47 proteomes, core to most methods for *de novo* ortholog grouping, which is an $O(n^2)$ operation
48 and a major computational bottleneck. Approaches to overcome this procedure include limiting
49 proteome comparisons by using a guiding-phylogeny^{5,6}, adapting alignment searching
50 parameters and heuristics to further boost speeds^{7,8}, or preliminary aggressive clustering of
51 proteins into coarse homolog groups⁹. Recently, graph-based and iterative-clustering
52 approaches have also allowed vast scalability to thousands of bacterial genomes, but are
53 primarily designed for application to a single species¹⁰⁻¹³.

54 Available orthology inference methods struggle to infer ortholog groups across large
55 datasets of taxonomically diverse genomes, potentially representing thousands of species, such
56 as a set of metagenome-assembled genomes (MAGs) related to a common microbiome. While
57 multiple methods exist to identify instances of previously established ortholog groups within the
58 predicted proteome of a metagenome¹⁴⁻¹⁷, these are unable to account for proteins not
59 represented in their database. Recently, independent advancements in methods to collapse
60 large protein sets based on sequence similarity have enabled rapid clustering of millions of
61 sequences¹⁸⁻²⁰. These approaches have even been used on massive protein datasets gathered
62 from across multiple metagenomic datasets²¹; however, more resolute delineation of functionally
63 analogous ortholog groups across thousands of genomes from multiple species remains difficult
64 to perform *de novo*.

65 Of relevance, within bacterial genomes, genes are often co-located within smaller,
66 discrete, multi-gene units, which we will broadly refer to as gene clusters. Examples of gene
67 clusters include operons^{22,23}, phages²⁴, metabolic gene clusters²⁵, biosynthetic gene clusters
68 (BGCs)²⁶⁻²⁹, and pathogenicity islands^{30,31}. Although less common, eukaryotic genomes can
69 also contain genes aggregated within discrete clusters³²⁻³⁴. Sometimes gene clusters are highly
70 conserved, encoding for products essential to the survival of the organism³⁵. In other cases, a
71 single gene cluster can exhibit variability in gene carriage and order across different strains or
72 species³⁶⁻³⁸. This is often the case for BGCs encoding specialized metabolites or virulence-
73 associated gene clusters, where evolution of gene content and sequence divergence can
74 influence fitness and contribute to adaptation within a changing ecosystem³⁹⁻⁴¹.

75 Syntenic conservation has been used to assist *de novo* identification of homologous
76 instances of a gene cluster of interest in diverse target genomes⁴²⁻⁴⁵. Homologous gene cluster
77 instances can then be comprehensively investigated to delineate homolog or ortholog groups of
78 the proteins found across them^{44,46}. While such targeted approaches can alleviate time and
79 computational resources by avoiding more comprehensive identification of orthologs at genome-
80 wide scales, currently available methods are mostly designed for specific types of gene clusters,
81 such as BGCs^{42,44,45}. Many of the software implementing such approaches also do not provide
82 support for uniform annotation of coding sequences in target genomes, which can decrease
83 sensitivity for gene cluster detection. In addition, most methods do not account for gene cluster
84 paralogy, which has been observed for BGCs in bacterial³⁸ and fungal genomes³³, or provide
85 specialized capabilities for finding gene clusters across fragmented genomes or metagenomic
86 assemblies³⁸.

87 Following identification of homologous gene clusters in target genomes, software to
88 understand the evolutionary relationships between gene cluster instances and infer protein
89 ortholog groups have largely applied coarse protein clustering and aimed to provide
90 visualization based exploration to users^{44,46-48}. Visual assessment of related gene clusters and

91 manual refinement of ortholog groups work well at smaller scales but become impractical when
92 dealing with hundreds to thousands of gene cluster instances. Scalability challenges are due to
93 both computational costs needed to render visuals as well as the figures becoming convoluted
94 and difficult to interpret. An effective solution to ease the identification of evolutionary trends
95 amongst homologous gene clusters is to first identify ortholog groups⁴⁴ and present information
96 pertaining to their conservation and sequence divergence within tabular reports^{10,38}. Such
97 tabular reports scale by the number of unique ortholog groups and can be organized by their
98 consensus order along gene cluster instances. We recently introduced construction of such
99 reports in a software suite for exploring microdiversity amongst homologous BGCs from a single
100 taxon³⁸; however, the functionality was difficult to use outside of the suite and reliant on
101 orthologous relationships between proteins of gene clusters being known in advance.

102 Here, we introduce the zol suite, providing functionalities for gene cluster detection and
103 subsequent inference and investigation of protein ortholog groups across homologous gene
104 clusters. The versatility and scalability of these programs is demonstrated through application to
105 three types of gene clusters within different genomic contexts including a virus within
106 environmental metagenomes, fungal secondary metabolite encoding biosynthetic gene clusters,
107 and a conserved polysaccharide antigen locus from the diverse bacterial genus of
108 *Enterococcus*.

109

110 Results

111

112 **fai and zol allow for the rapid inference of gene cluster orthologs across diverse** 113 **genomes**

114

115 The zol suite consists of three major programs: prepTG (*prepare target genomes*), fai
116 (*find additional instances*), and zol (*zoom on locus*) (**Figure 1A**). First, prepTG and fai can be
117 run to process a set of target genomes and rapidly search for a query gene cluster within them,
118 respectively. Afterwards, zol can perform reliable and efficient context-limited inference of
119 ortholog groups across homologous gene cluster instances identified using a flexible
120 InParanoid-type algorithm³. For each ortholog group, zol will further compute evolutionary
121 statistics, such as Tajima's D⁴⁹, and functional annotations, using several, diverse databases
122 suitable for a variety of gene clusters, including those specific to phages⁵⁰, virulence elements⁵¹,
123 and BGCs⁵². Ultimately, zol will summarize data in a table report where each row corresponds
124 to a distinct ortholog group. This report is automatically color formatted and provided as an
125 XLSX spreadsheet to allow for easy interpretation of the data, which can span thousands of
126 gene cluster instances.

127 To promote consistency in gene calling across target genomes, we have incorporated
128 computationally light-weight dependencies for *de novo* gene prediction in bacterial genomes^{53,54}
129 and protein-mapping in eukaryotic genomes⁵⁵ within prepTG, to prepare and format target
130 genomes for optimized gene cluster searching in fai (**Figure 1B**). prepTG also aims to provide a
131 convenient interface to transform genomic or metagenomic datasets into a format ready for
132 searching using fai. Options are available to download pre-built databases of distinct
133 representative genomes for 18 commonly studied bacterial taxa⁵⁶ or to build comprehensive

134 databases for any genus or species in the latest release of the Genome Taxonomy Database
135 (GTDB)⁵⁷.

136 `fai` features two key features which are absent in most existing methods for gene cluster
137 detection (**Figure 1C; Table S1; Supplementary Text**). First, it has an option to automatically
138 filter secondary instances of query gene clusters identified in target genomes, removing
139 potentially paralogous gene clusters from downstream investigations. Second, `fai` implements a
140 mode for searching for gene clusters in draft quality genomes, MAGs, or unbinned
141 metagenomic assemblies, where gene clusters might be fragmented across multiple scaffolds.
142 When this mode is activated, `fai` relaxes requirements for reporting a gene cluster as present in
143 a genome or metagenome if multiple homologous gene cluster regions are identified near
144 scaffold edges in a target genome and instead assesses whether reporting criteria are met in
145 unison across such instances (**Figure S1**). Similar to `prepTG`, `fai` also aims to provide
146 convenience for users and can accept query gene clusters in different formats to ease
147 searching for gene clusters and genomic islands cataloged in databases such as ICEberg⁵⁸,
148 MIBiG⁵², or IslandViewer⁵⁹. Query gene clusters can be provided as a coordinate along a
149 reference genome, in GenBank format, or as a set of proteins in FASTA format. In addition, to
150 simplify conservation and novelty assessment of a single isolate's BGCs, phages, and plasmids
151 relative to other genomes from the same genus or species, specialized wrapper programs of `fai`
152 are also provided within the `zol` suite (**Figure S2**).

153 `zol` will infer ortholog groups for proteins across homologous gene clusters and then
154 construct a tabular report with information on conservation, evolutionary trends, and annotation
155 for each individual ortholog group (**Figure 1D**). To make annotated reports generated by `zol`
156 more comprehensive for different types of gene clusters, several databases have been
157 included, such as VOGs⁵⁰, VFDB⁵¹, ISFinder⁶⁰, and CARD⁶¹. In addition, `zol` incorporates
158 HyPhy⁶² as a dependency and calculates various evolutionary statistics. Ultimately, beyond
159 high-throughput inference of ortholog groups across diverse genomic datasets, the rich tabular
160 report produced by `zol` provides complementary information to figures generated by
161 comparative visualization software such as `clinker`⁴⁶, `CORASON`⁴⁴, `gggenomes`⁶³, and `Easyfig`⁶⁴.

162 A key feature in `zol` is the ability to dereplicate gene clusters directly using `skani`⁶⁵, which
163 was recently shown to be more reliable at estimating average nucleotide identity (ANI) between
164 genomes of variable contiguity relative to comparative methods. Dereplication can allow for
165 more appropriate inference of evolutionary statistics to overcome availability or sampling biases
166 in genomic databases⁶⁶. It can also be used to subset distinct representative gene cluster
167 instances to make investigation using visualization software more tractable. Another important
168 ability of `zol` is a mode where users can provide a handful of known instances for a gene cluster
169 to estimate optimal parameters to search for additional instances of the gene cluster using `fai`.
170 We applied this functionality of `zol` on sets of homologous BGCs and phages to determine
171 distributions for search parameters in `fai` which users could consult as priors (**Figure S3;**
172 **Supplementary Text**).

173 Finally, `zol` allows for comparative investigations of gene clusters based on taxonomic or
174 ecological groupings⁶⁷⁻⁶⁹. For instance, users can designate a subset of gene clusters as
175 belonging to a specific population to allow `zol` to calculate ortholog group conservation across
176 just the focal set of gene clusters. In addition, `zol` will compute the fixation index⁷⁰, F_{ST} , for each
177 ortholog group to assess gene flow between the focal and complementary sets of gene clusters.

178

179 **Longitudinal tracking of a virus within lake metagenomic assemblies**

180

181 Metagenomic datasets represent a large reservoir of underexplored sequence
182 space^{71,72}. To demonstrate the ability of the zol suite to identify and investigate gene clusters in
183 metagenomes, we applied it to track a virus in a longitudinal metagenomic dataset profiling a
184 lake's microbiome over space and time⁷³.

185 We first identified large (≥ 20 kb) viruses, that were also predicted to represent circular
186 molecules, across a subset of the metagenomic assemblies corresponding to the earliest
187 sampling date⁷⁴. Afterwards, clustering based on the sequence and syntenic similarity of protein
188 domains led to the identification of a ~36kb highly conserved virus in two of the metagenomes
189 sampled from lower lake depths.

190 All 16 metagenomic assemblies, spanning five distinct sampling timepoints and four
191 distinct sampling depths, were processed through prepTG to identify coding sequences and
192 construct a database ready to search for gene clusters using fai. GenBank files with coding
193 sequence annotations for metagenomic assemblies generated by prepTG, amassing 27 Gb
194 total in size, were further provided as input for cblaster makedb, which serves a similar role to
195 prepTG in the cblaster suite to format genomic data for downstream gene cluster searches.
196 However, cblaster makedb does not feature the ability to perform *de novo* gene-calling for either
197 genomes or metagenomes and is not designed to accommodate the size of metagenomic
198 assemblies. During database construction, cblaster makedb required around 30 Gb of memory,
199 while prepTG needed less than 3 Gb of memory (**Figure S4A**).

200 Next, fai was used to perform a rapid, targeted search for this ~36 kb *Caudovirales* virus
201 across the full set of 16 metagenomes to identify additional instances of the virus. fai completed
202 its search of the metagenomes, featuring >20 million proteins and 10.7 million contigs, in less
203 than four minutes using 20 threads, performing similarly to cblaster, run using similar settings as
204 fai (**Figure S4B**). Of the 16 total metagenomes, the virus was found in ten metagenomes,
205 including all nine metagenomes surveying anoxic conditions ($p < 0.001$; one-sided Fisher's exact
206 test; **Figure 2A**). This is concordant with inferences for the host for the virus being *Rhodoferrax*,
207 which are purple bacterium featuring species classified as anaerobic photoheterotrophs^{73,75,76}. In
208 addition, *Rhodoferrax* classified MAGs from the metagenomic dataset were exclusively obtained
209 from anoxic conditions⁷³. To investigate how the gene repertoire of the virus evolved over time,
210 we next applied zol. zol-based analysis revealed that 45 (72.6%) of the 62 total distinct ortholog
211 groups were core to all instances of the virus across ten metagenomes with most completely
212 conserved in sequence over the course of 2.5 months (**Figure 2B; Table S2**). Furthermore, 15
213 of the 62 ortholog groups were not observed in the query viruses from the earliest sampling
214 date, suggesting the potential acquisition or duplication of genes in the virus during the span of
215 sampling at the lake.

216

217 **Investigating population-level and species-wide evolutionary trends of BGCs in the** 218 **eukaryotic species *Aspergillus flavus***

219

220 Low sensitivity for gene cluster detection in eukaryotic genome assemblies can arise
221 from their incompleteness, leading to gene clusters being fragmented across multiple

222 scaffolds^{77,78}, as well as challenges in *ab initio* gene prediction due to alternative splicing^{79,80}.
223 Therefore, many gene cluster detection software are either specific for bacterial genomes or
224 require coding sequence annotations for eukaryotic genomes to be provided by the user. To
225 overcome such challenges to user application, we integrated miniprot⁵⁵ into prepTG which
226 allows for mapping high-quality protein annotations from a reference genome to the remainder
227 of the genomes available for a species or genus. We showcase the ability of prepTG and fai to
228 simplify the reliable identification of gene clusters in eukaryotic genomes by using them to find
229 instances of two BGCs across genomes belonging to the fungal species *Aspergillus flavus*.

230 The genus of *Aspergillus* is a source of several natural products, including aflatoxins, a
231 common and economically impactful contaminant of food⁸¹. The genus also contains species
232 that are model organisms for studying fungal secondary metabolism^{34,82,83}. Examination of the
233 secondary metabolome of *A. flavus* has revealed that different clades or populations can exhibit
234 variability in their metabolite production despite high conservation of core BGC genes encoding
235 enzymes for synthesis of these metabolites^{37,84}. For instance, population B *A. flavus* were
236 identified as producing a greater abundance of the insecticide leporin B relative to populations A
237 and C^{37,85}. We showcase zol's ability to aid comparative analysis of gene clusters from different
238 populations through application to the leporin BGC. We further show how zol can detect
239 variation in sequence conservation for different genes from the aflatoxin BGC and be inclusive
240 of genes present in target genome annotations but missing in the query gene cluster, allowing
241 for comprehensive profiling of BGC auxiliary content.

242 Based on read alignment to a reference genome, the leporin cluster was recently
243 identified to be a core component of the *A. flavus* genome³⁷. However, a restricting factor in the
244 direct prediction of gene clusters in *A. flavus* assemblies is the lack of gene annotations, with
245 only 11 (5.1%) of 216 genomes from the species in NCBI's GenBank database having coding
246 sequence predictions (**Figure 3A**). Therefore, we mapped high-quality protein predictions for a
247 reference *A. flavus* genome⁸⁶ to the remainder of the 216 genomes available for the species.
248 Running fai in "draft mode" led to the identification of the leporin BGC within 212 (98.1%)
249 assemblies, consistent with the prior read mapping-based investigation suggesting that the BGC
250 was core to the species³⁷. In comparison, the CAGECAT server⁸⁷, which runs cblaster⁴⁵, was
251 limited to genomes with protein coding annotations available on NCBI and thus unable to
252 assess the remaining 205 genomes for the presence of the leporin BGC (**Figure 3B**). We also
253 investigated the ability of non-targeted approaches for BGC detection to identify the leporin
254 BGC by applying antiSMASH followed by BiG-SCAPE for clustering related BGCs and matching
255 them to characterized BGCs in the MIBiG database. When this approach was applied using
256 GenBank files prepared by prepTG, the gene cluster clan corresponding containing the leporin
257 BGC was found in all *A. flavus* genomes provided as input. However, when antiSMASH was run
258 using *de novo* gene prediction in antiSMASH based on GlimmerHMM⁸⁸ with *Cryptococcus* gene
259 annotation models, recovery of the leporin BGC was limited (**Figure 3B**).

260 Of the 212 genomes with the leporin BGC identified by fai, 202 contained instances that
261 were high-quality and not near scaffold edges. This set of 202 instances of the gene cluster was
262 further investigated using zol with options to perform comparative investigation of BGC
263 instances from *A. flavus* population B genomes to instances from other populations. High
264 sequence conservation was observed for all genes in the leporin gene cluster as previously
265 reported³⁷ (**Table S3**). Further, alleles for genes in the BGC from population B genomes were

266 generally more similar to each other than to alleles from outside the population, as indicated by
267 high F_{ST} values (>0.85 for 9 of 10 genes) (**Figure 3C**; **Table S3**). While regulation of secondary
268 metabolites in *Aspergillus* is complex⁸⁹, zol analysis showed that the three essential genes for
269 leporin production⁸⁵ also had the lowest variation in the 100 bps upstream their exonic
270 coordinates (**Figure S5**). This suggests higher variability is occurring in the transcription of the
271 accessory *lep* genes within the species. This supports experimental evidence that has shown
272 gene knockouts depleting certain leporin species will still permit the production of others⁸⁵.
273 *fai* and *zol* were also applied to the BGC encoding aflatoxin across *A. flavus*⁹⁰ (Table
274 S4). Similar to the leporin BGC, the aflatoxin BGC was highly prevalent in the species and found
275 in 71.8% of genomes. However, in contrast to the leporin BGC, the aflatoxin BGC contained
276 several genes with positive Tajima's D values, indicating greater sequence variability for these
277 coding regions across the species (**Figure 3D**). One of the genes with a positive Tajima's D
278 value was *afIX*, which has been shown to influence conversion of the precursor veriscolorin A to
279 downstream intermediates in the aflatoxin biosynthesis pathway⁹¹ (**Figure 3E**). An abundance
280 of sites with mid-frequency alleles in the oxidoreductase encoding gene could represent
281 granular control for the amount of aflatoxin relative to intermediates produced. The polyketide
282 synthase gene *pksA* had the lowest Tajima's D value of -2.4, which suggests it is either highly
283 conserved or under purifying selection (**Figure 3F**). In addition, because the reference proteome
284 used to infer genomic coding regions was constructed recently⁸⁶, *fai* and *zol* detected several
285 highly conserved genes within the aflatoxin BGC that are not represented in the original
286 reference gene cluster input for *fai*⁵². This includes a gene annotated as a noranthrone
287 monooxygenase and recently characterized as contributing to aflatoxin biosynthesis^{92,93} (**Figure**
288 **3D**).

289

290 **Identification of the Enterococcal polysaccharide antigen and assessment of context** 291 **restricted orthology inference**

292

293 To demonstrate the ability of *zol* and *fai* to reliably identify ortholog groups across
294 multiple species and thousands of genomes, we used the tools to assess the distribution of the
295 enterococcal polysaccharide antigen (Epa) and its individual genes across the diverse genus of
296 *Enterococcus*. Because previous comparative genomic investigations have been performed
297 between *epa* loci from different species^{94,95}, we also showcase how such prior insight can be
298 used to tailor parameters in *fai* for searching for the locus across the full genus and how results
299 from *fai* can be assessed for appropriate selection of parameter values in *zol*.

300 The Epa is a signature component of the cellular envelope of multiple species within
301 *Enterococcus*⁹⁴⁻⁹⁷ and has mostly been characterized in the species *Enterococcus faecalis*^{96,98-}
302 ¹⁰¹. While molecular studies have provided evidence that the locus contributes to enterococcal
303 host colonization¹⁰⁰, evasion of immune systems¹⁰², and sensitivity to antibiotics¹⁰³ and
304 phages^{103,104}, it was only recently that the structure of Epa was resolved and a model for its
305 biosynthesis and localization formally proposed¹⁰¹. A homologous instance of the *epa* locus was
306 identified in the other prominent pathogenic species from the genus, *Enterococcus*
307 *faecium*^{94,95,105}; however, the prevalence and conservation of *epa* across the diverse genus of
308 *Enterococcus*¹⁰⁶⁻¹⁰⁸ remains poorly studied.

309 We first assessed the performance of fai and zol to identify *epa* loci across
310 representative genomes for each of the 92 species of *Enterococcus* in GTDB R214⁵⁷ and
311 subsequently delineate protein ortholog groups relative to other methods. Specifically, we
312 compared the runtime and ortholog group predictions of fai and zol to the combination of
313 cblaster and clinker as well as OrthoFinder, an established software for multi-species ortholog
314 group delineation, run on full genomes. For this comparison, the parameter settings for fai and
315 cblaster as well as zol and clinker were adapted to match each other more closely, with an
316 exception being to run fai in draft-mode, which lacks an analogous feature in cblaster. The
317 combination of fai and zol was the fastest of the three methods tested and able to identify
318 ortholog groups for the *epa* locus in approximately one minute (**Figure 4A, S6**). Orthology
319 inferences from fai and zol exhibited high overlap with orthology predictions by the alternate two
320 methods, finding 96.3% of ortholog protein pairs identified by at least two of the three methods
321 (**Figure 4B**). We also applied all three methods to determine *epa* locus orthologs across low
322 quality representative genomes for each species to demonstrate the convenience of fai's ability
323 to be run in "draft mode" and improve sensitivity for detecting fragmented gene clusters in
324 comparison to cblaster. fai identified 2.1-fold more exclusive ortholog pairs in common with
325 OrthoFinder, expected to be relatively robust to the effects of assembly fragmentation, than the
326 number of ortholog pairs shared exclusively by cblaster and clinker with OrthoFinder (**Figure**
327 **4C**). In addition, we performed evolutionary-simulation of the *epa* locus, allowing for sequence
328 gains and losses, and assessed context-limited orthology inference by zol, clinker and
329 OrthoFinder (**Figure S7; Supplementary Text**). zol was able to recover a high fraction of true
330 positive ortholog relations and was the best method at avoiding prediction of false positive
331 orthologs.

332 Next, to properly and comprehensively assess the distribution of *epa* across the entire
333 set of 5,291 genomes in GTDB classified as one of the 92 *Enterococcus* species⁵⁷, we applied
334 fai with more careful consideration of parameter values and requested more advanced features
335 for gene cluster detection. A sensitive searching criterium was selected based on prior
336 comparative genomics for the locus^{94,95} and its coordinates along the *E. faecalis* V583 genome
337 as a reference^{99,101}. For detection of *epa* orthologous regions, co-location of at least seven of
338 the 14 *epa* genes previously identified as conserved in both *E. faecalis* and *E. faecium* was
339 required. The default threshold for syntenic conservation of homologous instances to the query
340 gene cluster was disregarded to increase sensitivity for the detection of *epa* in enterococcal
341 species more distantly related to *E. faecalis*. In addition, key proteins were specified and the
342 length of the flanking context to include as part of the loci was expanded. Using these criteria,
343 5,085 of the genomes assessed were found to possess an *epa* locus, with phylogenomic
344 investigations further revealing that the locus is highly conserved in three of the four major
345 clades of *Enterococcus* (**Figure 4D; Table S5**).

346 Based on fai's reports, we realized that to achieve optimal clustering for ortholog groups
347 across the diverse set of *epa* loci identified, we needed to lower the default thresholds for
348 percent identity and coverage that protein pairs needed to exhibit for being considered as
349 orthologs (**Figure 4D; Table S5**). We ran zol on both the full set of 5,052 high-quality *epa* loci
350 and only loci from species representative genomes. For the comprehensive analysis, zol was
351 able to identify 14 ortholog groups as core or near-core, found in >90% of loci instances (**Table**
352 **S6**). When provided 30 threads, zol completed in 30.7 hours and had a maximum memory

353 usage of 101.3 GB. The more restricted analysis of zol to investigate *epa* instances from 65
354 species representative genomes was to allow for assessing the quality of ortholog group
355 predictions using phylogenetics (**Table S7**). After applying zol on *epa* from species
356 representative genomes, orthology predictions were assessed through construction of a
357 maximum-likelihood phylogeny of *epa* associated glycosyltransferases. Ortholog groups which
358 corresponded to glycosyltransferases from *E. faecalis* V583 were labelled on the phylogeny and
359 confirmed to match distinct phylogenetic clades, which suggests their appropriate delineation
360 (**Figure 4EF**). zol further identified several *epa* associated glycosyltransferase ortholog groups
361 that were absent in the *E. faecalis* representative genome and other representative genomes
362 from the *E. faecalis* clade (**Figure 4G**). These distinct glycosyltransferases might impact the
363 final structure or decoration of Epa in other *Enterococcus* species.

364

365 **zol identifies genetic diversity of *epaX*-like glycosyltransferases in *E. faecalis***

366

367 zol features several options related to the dereplication of input gene clusters to retain
368 only distinct representative instances for orthology inference and other downstream analytics
369 (**Figure S8**). Importantly, the application of these methods can substantially reduce zol's
370 runtime and impact some of the evolutionary statistics computed (**Figure S8, S9, S10,**
371 **Supplementary Text**). Whether dereplication is appropriate for a particular analysis should thus
372 be carefully considered by users depending on their research aims. In particular, dereplication
373 can impact investigations for highly sequenced bacterial taxa, including the opportunistic
374 pathogen *E. faecalis*. For such pathogens, certain lineages, such as those commonly isolated at
375 clinics, might be overrepresented in genomic databases, and the researcher may find it
376 beneficial for the analysis to apply dereplication.

377 To showcase the scalability of zol and its ability to expand knowledge for even well-
378 studied gene clusters, we applied it to high-quality, complete *epa* loci from 1,232 *E. faecalis*
379 genomes without dereplication. In accordance with prior studies^{94,101}, zol was able to distinguish
380 core and strain-variable patterns. The report from zol showed that one end of the locus
381 corresponds to genes which are highly conserved and core to *E. faecalis* (*epaA-epaR*), whereas
382 the other end contained strain-specific genes (**Figure 5A; Table S8**). Using zol, we further
383 found that variably conserved genes exhibit high sequence dissimilarity, as measured using
384 both Tajima's D and average sequence entropy, in comparison to the core genes of the locus
385 (**Figure 5BC**). These statistics were robust to the application of dereplication and thus unlikely
386 to be heavily impacted by well-sequenced lineages (**Figure S9, S10**).

387 One ortholog group, corresponding to the glycosyltransferase *epaX*, exhibited
388 substantially higher sequence variation than other *epa* associated glycosyltransferases (**Figure**
389 **5BD**). This finding was further validated through phylogenetic analysis of glycosyltransferases
390 from the species, which highlighted the breadth of diversity observed for the *epaX* ortholog
391 group relative to other *epa* associated glycosyltransferases (**Figure 5E**).

392

393 **Discussion**

394

395 Here fai and zol are introduced to enable large-scale evolutionary investigations of gene
396 clusters in diverse taxa. Together these tools overcome current bottlenecks in computational

397 biology to infer orthologous sets of genes at scale across thousands of diverse genomes and
398 large metagenomic assemblies.

399 The set of input gene clusters for zol does not need to be produced by fai. cblaster⁴⁵ is
400 another tool that can identify instances of a query gene cluster within a set of target genomes
401 and extract them in GenBank format for downstream investigations using zol. For those lacking
402 computational resources needed for fai analysis, cblaster offers remote searching of BGCs
403 using NCBI's BLAST infrastructure and non-redundant databases. More recently, CAGECAT⁸⁷,
404 a highly accessible web-application for running cblaster, was also developed and can similarly
405 be used to identify and extract gene cluster instances from genomes represented in NCBI
406 databases. In contrast to these tools, prepTG and fai feature algorithms and options for users
407 interested in: (i) identification of gene clusters in metagenomes, (ii) performing standardized
408 gene annotation across target genomes, (iii) improved sensitivity for gene cluster detection in
409 draft-quality assemblies, and (iv) automated filtering of secondary, or paralogous, matches to
410 query gene clusters. In addition, users can apply zol to further investigate homologous sets of
411 gene clusters identified from IslandCompare¹⁰⁹, BiG-SCAPE⁴⁴, or vConTACT2¹¹⁰ analyses,
412 which perform comprehensive clustering of predicted genomic islands, BGCs, or viruses.

413 The application of fai to identify gene clusters in metagenomes is demonstrated here
414 through rapid, targeted detection of a virus across lake metagenomic assemblies. We expect
415 that both fai and zol will gain greater relevance for metagenomic applications in the future as
416 long-read sequencing becomes cheaper. Importantly, the tools can be applied directly on
417 assemblies without the need for binning scaffolds into MAGs, avoiding complications associated
418 with binning¹¹¹. In addition to their application to viral tracking, fai and zol's application to
419 metagenomes could be useful for assessing the presence of concerning transposons carrying
420 antimicrobial resistance traits¹¹²⁻¹¹⁴ and identifying novel auxiliary genes within known BGCs
421 which may tailor the resulting specialized metabolites and expand chemical diversity^{115,116}.

422 Reidentifying gene clusters in eukaryotic genomes remains difficult due to technical
423 challenges in gene prediction owing to the presence of alternative splicing. The ability of fai and
424 zol to perform population-level genetics on BGCs from the eukaryotic species *A. flavus* was
425 demonstrated. While there are over 200 genomes of *A. flavus* in NCBI, only 5.1% have coding-
426 sequence information readily available. We used miniprot⁵⁵ to map high quality gene coordinate
427 predictions from a representative genome in the species⁸⁶ to the remainder of genomic
428 assemblies with prepTG which enabled high sensitivity detection of BGCs with fai. Our analysis
429 provides additional support that the leporin BGC is conserved across the species³⁷ using an
430 assembly-based approach.

431 The ability of zol to identify ortholog groups across 5,052 gene cluster instances from 71
432 distinct species using limited computational resources was demonstrated through investigation
433 of the *epa* locus across *Enterococcus*. While such large-scale investigations will be largely
434 limited to those with access to a server, we expect datasets to often feature some degree of
435 species level redundancy. For instance, 80.2% of the 5,052 *epa* instances were from only two
436 species, *E. faecalis* and *E. faecium*. Thus, to alleviate computational costs, we have included
437 functions for dereplication of gene clusters and reinflation of ortholog groups in zol. Applying
438 these features to the comprehensive set of *epa* loci using 30 threads, reduced runtime from
439 30.7 to 3.5 hours and maximum memory usage from 101.3 GB to 83.2 GB (**Table S9**).

440 We further assessed the quality of ortholog group predictions by fai and zol using
441 phylogenetic investigations and comparisons with other software for homology inference.
442 Specifically, we compared orthology inference results from fai and zol to predictions obtained
443 from the combination of cblaster and clinker as well as OrthoFinder¹¹⁷, which was used to
444 detect ortholog groups at the genome-wide scale. Notably, clinker⁴⁶, which is developed by the
445 authors of cblaster, is primarily designed to produce interactive visualizations showing
446 relationships between related gene cluster instances. clinker's application of single-linkage
447 clustering to determine related sets of genes and to color matching genes in figures is expected
448 to produce relatively coarse ortholog groups. OrthoFinder was chosen as a representative
449 method for standard multi-species orthology inference because it has been shown to perform
450 well for several criteria in prior benchmarking studies^{117,118}. Through application to identification
451 of ortholog groups for diverse *epa* loci from multiple distinct species and evolutionary simulation
452 of the locus from *E. faecalis*, we found zol produces reliable orthology predictions that are
453 mostly in accordance with alternate orthology inference methods while exhibiting restraint for
454 over clustering. In the future, we are considering further improving the algorithm for ortholog
455 group classifications within zol. Specifically, we might take a similar approach to OrthoFinder in
456 which coarse ortholog groups are first identified and later refined using phylogenetics.

457 Our investigation of *epa* loci from multiple species revealed the presence of a multitude
458 of glycosyltransferases associated with production or decoration of the polysaccharide,
459 including some that are absent in the representative *E. faecalis* genome, the species in which
460 the polysaccharide has been most extensively characterized. Through population-genetic
461 investigations of the locus in *E. faecalis* using zol, we further determined that an ortholog group
462 containing *epaX*-like glycosyltransferases possessed high sequence divergence relative to other
463 glycosyltransferases associated with the locus. In addition to influencing the ability of *E. faecalis*
464 to colonize hosts¹⁰⁰, mutations in *epaX* and other genes from the ortholog group have also been
465 shown to impact susceptibility to phage predation^{119–122}. Therefore, we hypothesize that
466 extensive evolution of the *epaX* ortholog group is a result of contrasting selective forces,
467 pressuring *E. faecalis* to retain or (re-)acquire the glycosyltransferase to gain a fitness
468 advantage within hosts but also lose the gene to escape phage predation.

469 **Conclusions**

470
471
472 Practically, zol presents a comprehensive analysis tool for comparative genetics of
473 related gene clusters to facilitate detection of evolutionary patterns that might be less apparent
474 from visual analysis. Fundamentally, the algorithms presented within fai and zol enable the
475 reliable detection of orthologous gene clusters, and subsequently orthologous proteins, across
476 multi-species datasets spanning thousands of genomes and help overcome a key barrier in
477 scalability for comparative genomics.

478 **Methods**

479 **Software availability**

480
481
482

483 zol is provided as an open-source software suite, developed primarily in Python3 on GitHub at:
484 <https://github.com/Kalan-Lab/zol>. Docker and Bioconda¹²³ based installations of the suite are
485 supported. For the analyses presented in this manuscript, we used v1.4.1 of the zol software
486 package¹²⁴. Version information for major dependencies of the zol suite^{53,55,62,65,125–132} and other
487 software used^{44,74,133} for analyses in this study is provided in Table S10. Code and input files for
488 generation of figures in this manuscript are provided separately on GitHub at:
489 https://github.com/Kalan-Lab/Salamzade_etal_zol.

490

491 **Availability of data and materials**

492

493 Genomes and metagenomes used to showcase the application of fai and zol are listed with
494 GenBank accession identifiers in Table S11. Total metagenomes and their associated
495 information from Lake Mendota microbiome samplings were originally described in Tran *et al.*
496 2023⁷³ and deposited in NCBI under BioProject PRJNA758276. Genomic assemblies available
497 for *A. flavus* in NCBI's GenBank database on Jan 31st, 2023 were downloaded in FASTA
498 format using ncbi-genome-download (<https://github.com/kblin/ncbi-genome-download>).
499 Genomic assemblies for *Enterococcus* that met quality and taxonomic criteria for belonging to
500 the genus or related genera (e.g. *Enterococcus_A*, *Enterococcus_B*, etc.) in GTDB⁵⁷ release
501 R207 were similarly downloaded from NCBI's GenBank database using ncbi-genome-download
502 in FASTA format.

503

504 Assessment of compute time, memory usage, and disk space: The UNIX *time* command was
505 applied to measure the runtime and memory usage of programs. Specifically, the “Elapsed (wall
506 clock) time” was regarded as the runtime and the “Maximum resident set size (kbytes)” as the
507 maximum memory usage. The UNIX *du* command was used to measure the final disk space
508 used by various programs. All analyses were computed on the same server running Ubuntu
509 18.04.06 LTS with AMD EPYC 7451 24-Core processors, 472 GB of 288-Pin DDR4 random-
510 access memory, and a Samsung 970 Pro solid disk drive.

511

512 **Overview of tools and algorithms**

513

514 prepTG - processing and preparing target genomes for searching with fai: prepTG allows users
515 to create a database of target genomes that can be searched for homologous instances of
516 query gene clusters with fai. In addition to formatting and producing files for optimizing fai
517 searches, prepTG integrates pyrodigal⁵³, prodigal⁵⁴, and miniprot⁵⁵ for gene-calling or protein-
518 mapping in prokaryotic and eukaryotic genomes as well as metagenomes to aid consistency in
519 fai's performance and limit bias due to potential differences in gene-calling methods. For
520 miniprot-based protein-mapping, coding sequence predictions are required to exhibit an identity
521 of at least 80% to the reference protein and instances of overlapping mRNA and exon features
522 are resolved by retaining only the highest scoring mappings.

523 prepTG also features options to download pre-built databases for select bacterial taxa
524 that are commonly studied⁵⁶, such as ESKAPE pathogens, or to download all genomes
525 belonging to any genus or species in GTDB R214⁵⁷ and subsequently construct a database *ab*
526 *initio*.

527

528 *fai - automated identification of homologous instances of gene clusters:* *fai* allows for rapid
529 detection of gene clusters in target genomes. It accepts a target genomes database prepared
530 by prepTG and query gene cluster(s). Query gene cluster(s) can be provided in one of three
531 formats: (i) GenBank file(s) with CDS features, (ii) a coordinate along a reference genome, or
532 (iii) a set of proteins. When using coordinates along a reference genome to define a gene
533 cluster, *fai* re-performs gene-calling along the reference using pyrodigal⁵³ and extracts a local
534 GenBank file corresponding to the specified region.

535 *zol* implements HMM-based and CDS separation-based approaches for determining
536 homologous gene cluster instances in target genomes, which can further be combined in a
537 hybrid approach. For both approaches, homologs of proteins from query gene clusters are first
538 searched for in predicted proteomes of target genomes using DIAMOND alignment¹³⁰. Then, in
539 “Gene-Clumper” mode, which is the default, scaffolds with homologs of query proteins are
540 dynamically assessed for whether homologs are within a maximum number of CDS predictions
541 to be regarded as belonging to the same gene cluster. In “HMM” mode, scaffolds of target
542 genomes are instead scanned gene-by-gene using an HMM and neighborhoods or sets of
543 genes are regarded as being in a state of homology to the query gene cluster if several
544 individual genes depict homology to the proteins from the query gene cluster(s). The algorithm
545 is similar to *IsaBGC-Expansion*³⁸, however, it is not dependent on a preliminary genome-wide
546 orthology grouping analysis and thus features a different set of filters to still enable high-
547 throughput automated detection of homologous gene cluster segments as a result. *IsaBGC-*
548 *Expansion* is reliant on a preliminary orthology analysis to identify BGC-specific genes that
549 could be used to differentiate true homologous instances of BGCs and customize weighting of
550 HMM emission probabilities for distinct genes. It further requires the length of genes within
551 putative homologous regions to be within a certain deviation from the median length of known
552 gene instances. In contrast, *fai* has preconfigured emission probabilities which can be
553 customized by users and has no length requirement for potential homologous instances of
554 genes. *fai* further allows the “HMM-based” approach to be run with the parameter for
555 aggregating CDS predictions for the “Gene-Clumper” mode, whereby, gene cluster segments
556 detected by the HMM can be joined with other such segments if they are within a certain
557 number of CDS features from each other. Similar to *IsaBGC-Expansion*, syntenic similarity
558 between candidate and query gene cluster segments can also be used to filter candidate
559 segments using a gene cluster-wide correlation metric³⁸.

560 By default, *fai* requires filters pertaining to the number of genes from query gene clusters
561 to be met for each homologous gene cluster candidate segment. However, in “draft mode”,
562 thresholds for detection of gene clusters within target genomes are assessed in aggregate for
563 putative gene cluster segments found near scaffold edges (< 2,000 bp). Visual reports produced
564 by *fai* showcasing the sequence similarity of target genome proteins to the query protein(s) can
565 then be manually investigated by users to assess the validity of fragmented gene cluster
566 instances. In addition, *fai* features an option to filter for paralogous, overlapping candidate
567 segments of a gene cluster in target genomes and offers an intuitive visualization of gene
568 cluster segments, if requested, to allow users to assess their quality, including proximity of
569 candidate segments to scaffold edges. Together, these options enable the large-scale

570 identification of orthologous gene clusters across genomes which can then be leveraged by zol
571 to perform context-specific inference of protein ortholog groups.

572 In addition to a directory of homologous gene clusters in GenBank format, to serve as
573 input for zol analysis, and a small set of visual PDF files, fai generates an in-depth report on
574 which target genomes have the query gene cluster as an XLSX spreadsheet. This spreadsheet
575 includes information such as the average amino acid identity (AAI), syntenic similarity, and
576 number of conserved genes for gene clusters from target genomes relative to the query gene
577 cluster. The spreadsheet allows for easy sorting of various columns to assist identification of
578 which target genomes feature a gene cluster to the desired degree of similarity for the user.

579

580 *zol - computes a variety of evolutionary statistics and can perform gene cluster specific*

581 *dereplication*: The zol workflow begins by processing the input directory of gene cluster
582 GenBank files to assess validity and perform filtering of gene clusters or individual proteins.
583 Filtering can be performed at the gene cluster level by requesting filtering of draft-quality gene
584 clusters, those marked as being near scaffold edges, or low-quality gene clusters, those with
585 $\geq 10\%$ missing base-pairs (e.g. Ns) in their sequence. Filtering of individual proteins which are
586 near scaffold edges can also be performed if fai was used to identify the input gene cluster set,
587 because fai marks these proteins with a special feature tag in the resulting gene cluster
588 GenBank files.

589 Next, zol will perform dereplication of gene clusters, if requested by users, with skani⁶⁵
590 by clustering gene clusters which depict some user-defined coverage and identity thresholds
591 using single linkage clustering or more resolved MCL-based clustering, for which the inflation
592 parameter can be adjusted. Representative gene clusters are selected from each cluster as part
593 of the dereplication based on maximum length and, if comparative analysis is requested,
594 whether the representative gene cluster is part of the focal or focal-complement set of gene
595 cluster instances specified by the user.

596 The input set of gene clusters or set of dereplicated representative gene clusters is then
597 used to identify protein ortholog groups with an InParanoid-type approach³. Briefly,
598 DIAMOND¹³⁰ is used to perform all vs. all pairwise alignment between proteins from the set of
599 gene clusters after which the alignments are processed to identify reciprocal best hits (RBH)
600 between pairs of gene clusters. In-paralogs are identified within each gene cluster based on
601 whether two coding sequences depict more similarity to each other than one does to an RBH
602 with a different gene cluster. Bitscores, standardized through division by reflexive bitscore
603 values for query proteins, are used to assess homology. Specifically, the average normalized
604 bitscore between each pair of orthologs and in-paralogs is recorded. Afterwards, bitscores
605 between such protein pairs are further standardized through dividing them with the average
606 values between pairs of gene clusters to aid proper clustering of proteins downstream. This is
607 akin to the genome-wide normalization procedure recommended in OrthoMCL, owing to the
608 realization that orthologs between distantly related species are also more likely to exhibit lower
609 sequence similarity, which should be corrected for prior to MCL clustering². This information is
610 input into MCL with the inflation parameter set to 1.5, similar to other orthology inference
611 methods^{7,117}. The inflation parameter and minimum identity and coverage cutoffs to consider
612 valid pairs of in-paralogs and orthologs are adjustable by users.

613 Reinflation can also be requested by users to expand ortholog groups to include proteins
614 from the full input set of gene clusters if gene cluster dereplication was requested¹⁰. Reinflation
615 of ortholog groups is performed by first performing comprehensive and granular clustering of
616 proteins from all input gene clusters using CD-HIT¹²⁸, requiring proteins to depict >98%
617 sequence similarity and > 95% bi-directional coverage to the representative sequences of
618 clusters. Proteins in CD-HIT clusters are then mapped to ortholog groups if they co-cluster with
619 proteins from dereplicated gene clusters which are already assigned to ortholog groups.
620 Dereplication and reinflation are not recommended if sequence redundancy amongst the set of
621 input gene clusters is low. Stringent cutoffs used for CD-HIT clustering during reinflation assume
622 that dereplication was also run with stringent parameters to only collapse highly similar gene
623 clusters. Otherwise, reinflation could miss more distant instances of ortholog groups, resulting in
624 an underestimation of ortholog group conservation amongst gene clusters.

625 Next, zol will partition protein and nucleotide sequences from gene clusters according to
626 ortholog groups, perform protein alignment using MUSCLE¹³², and create codon alignments
627 using PAL2NAL¹³⁴. We also offer an option to use reference proteins to refine and filter
628 sequences based on multiple sequence alignment using MUSCLE¹³², which might be useful to
629 further filter intronic sequences in eukaryotic ORFs. Codon alignments are filtered for regions
630 with high ambiguity ($\geq 10\%$ gaps) using trimAL¹²⁶ which are then used downstream for
631 calculation of evolutionary statistics and to construct approximate maximum-likelihood
632 phylogenies using FastTree 2¹²⁷ for each ortholog group. Consensus protein sequences for
633 each ortholog group are finally constructed using HMMER3¹²⁹.

634 Using protein consensus sequences of each ortholog group, zol is next able to linearize
635 annotation of ortholog groups with various annotation databases including KOfam¹⁴, the PGAP
636 database¹³⁵, VFDB⁵¹, CARD⁶¹, MIBiG⁵², ISfinder⁶⁰, the PaperBLAST database¹³⁶, and Pfam¹³⁷.
637 A custom FASTA file can also be provided by users to annotate ortholog groups. The best hit
638 per ortholog group for each annotation database is selected by score, if annotation is HMM
639 based¹³⁸, or bitscore, if it is DIAMOND alignment based¹³⁰, and a default E-value cutoff of 1e-5.
640 The E-value of the alignment is provided in the zol report for each putative annotation except
641 Pfam domains. However, for Pfam annotations, only domains meeting trusted thresholds are
642 reported.

643 Next, zol will compute basic statistics per ortholog group including the consensus order,
644 consensus directionality, whether proteins are single-copy across gene clusters, the median
645 length of ortholog group sequences, their median GC% percentage, and GC skew values. The
646 consensus order and directionality are performed similarly to *IsaBGC-PopGene*³⁸. Afterwards, in
647 the sixth step, zol will calculate evolutionary statistics for each ortholog group including Tajima's
648 D^{49} , the proportion of filtered codon alignments which correspond to segregating sites, the
649 average sequence entropy of the filtered codon alignment and the 100 upstream region, and the
650 median and maximum Beta-RDgc. Beta-RDgc is a statistic that is derived from the Beta-RD
651 statistic which we described in *IsaBGC*³⁸ and measures the divergence of a pair of protein
652 sequences based on the expected divergence between the gene clusters. Values below one
653 suggest that protein divergence is larger for the pair than expected based on other shared
654 proteins between the two gene clusters; conversely, the opposite trend might suggest high
655 conservation of the particular protein between the gene clusters and potentially gene-specific
656 horizontal gene transfer. Finally, we perform site-specific selection analyses using the FUBAR¹³⁹

657 and GARD¹⁴⁰ methods offered in the HyPhy suite. While highly scalable relative to comparable
658 methods¹³⁹, these analyses can still take considerable time and are turned off by default.
659 Importantly, GARD recombination detection¹⁴⁰ and partitioning of input alignments for ortholog
660 groups can also be used for alternate HyPhy analyses with HyPhy Vision⁶², to extend beyond
661 the site-specific selection analyses using FUBAR¹³⁹ supported directly in zol.

662 Prior to the generation of a final report, zol allows users to perform an optional
663 comparative analysis between user-defined set(s) of focal and complementary or alternate gene
664 cluster instances. In these comparative analyses, the conservation and fixation index⁷⁰ is
665 calculated for each ortholog group.

666 Finally, we generate a consensus report and a spreadsheet in XLSX format where each
667 row corresponds to an ortholog group and columns correspond to basic statistics, evolutionary
668 statistics, and annotation information. Quantitative fields are automatically colored to make
669 visual detection of patterns easier for users. A basic heatmap showing the presence of ortholog
670 groups across gene clusters is also produced.

671 zol additionally features two alternate modes that can be triggered via specific
672 arguments. First, the “only-orthologs” argument will invoke zol to only compute ortholog groups
673 and exist after determining them. Second, the “select_fai_params_mode” argument allows
674 users to provide a handful of known instances for a gene cluster and determine appropriate
675 thresholds for searching for additional instances of the gene cluster using fai. This mode
676 assumes that the known instances provided are representative of the breadth of diversity
677 expected for the gene cluster amongst the target genomes being searched.

678
679 *abon, atpoc, and apos – tools for assessing novelty and conservation of BGCs, phages, and*
680 *plasmids from a single strain*: The zol suite features three small wrapper programs called abon,
681 atpoc, and apos which assess the conservation and novelty of a single genome’s BGC-ome,
682 phage-ome, and plasmid-ome, respectively, relative to a target genome database constructed
683 by prepTG. The target genomes database could be all other genomes belonging to the focal
684 genome’s species or genus. The three programs are wrappers of fai but also offer a simple
685 BLAST search alternative, to more thoroughly check for whether individual genes from BGCs,
686 phages, and plasmids are present in the target genomes being searched. These tools accept
687 results from standard software for annotation of BGCs^{133,141}, phages^{74,142,143}, and plasmids^{143,144}
688 but do not integrate them within the suite. Similar to fai and zol they produce auto-formatted
689 XLSX spreadsheets as primary results.

690

691 **Application of fai and zol to track a virus within lake metagenomes**

692

693 VIBRANT was used to identify viral contigs or sub-contigs in the three total metagenomes from
694 Tran *et al.* 2023⁷³ sampled on the earliest date of 07/24. Afterwards, predicted circular contigs
695 were clustered using BiG-SCAPE⁴⁴ which revealed a ~36 kb virus was found in two of the three
696 metagenomes.

697

698 prepTG was run on all 16 total metagenomic assemblies from the Tran *et al.* 2023 study,
699 performing gene calling with pyrodigal in metagenomics mode⁵³ to prepare for comprehensive
700 targeted searching of the virus with fai. fai was run with largely default settings, with filtering of

701 secondary instances of the virus requested to retain only the best matching scaffold or scaffold
702 segment resembling the queries. In addition, the syntenic correlation requirement of hits to the
703 query gene clusters was turned off to account for the circular nature of the virus, which the
704 assessment is not designed for. To assess the performance of cblaster for preparing the target
705 metagenomes database and subsequently searching for the virus, we provided GenBank files
706 with CDS features produced by prepTG as input for cblaster makedb and adjusted searching
707 parameters for cblaster search to more closely match what we used for fai.

708

709 **Microevolutionary investigations of leporin and aflatoxin BGCs in *Aspergillus flavus***

710

711 Genomic assemblies downloaded from NCBI GenBank were processed using prepTG. Of the
712 217 genomic assemblies downloaded, one, GCA_000006275.3, was dropped from the analysis
713 because the original GenBank file had multiple CDS features with the same name, leading to
714 difficulties in performing BGC prediction with antiSMASH¹³³, and because alternate assemblies
715 were available for the isolate. prepTG was run on all assemblies with miniprot⁵⁵ based gene-
716 mapping of the high-quality gene coordinate predictions available for *A. flavus* NRRL 3357
717 (GCA_009017415.1)⁸⁶ requested. Target genomes were then searched for the leporin
718 (BGC0001445) and aflatoxin (BGC0000008) BGCs using GenBank files downloaded from
719 MIBiGv3⁵² as queries. For leporin, AFLA_066840, as represented in the MIBiG database, was
720 treated as a key protein required for detection of the BGC. Similarly, for aflatoxin, PksA
721 (AAS90022.1), as represented in the MIBiG database, was treated as a key protein required for
722 detection of the BGC. Draft-mode and filtering of paralogous segments was requested. For both
723 analyses, ortholog groups found in fewer than 5% of gene cluster instances were disregarded.

724

725 We reidentified population B as previously delineated³⁷ using k-mer based ANI estimation¹⁴⁵ and
726 neighbor-joining tree construction¹⁴⁶. A discrete clade (n=81) in the tree was validated to feature
727 all isolates previously determined as part of population B³⁷ and thus regarded as such.

728

729 For comprehensive and *de novo* BGC prediction, antiSMASH was run on the 216 genomic
730 assemblies with 'glimmerhmm' requested for the option '--genefinding-tool'. Similarly,
731 antiSMASH was also run on full GenBank files for genomes generated by prepTG from
732 reference proteome-mapping via miniprot. For one genome, antiSMASH was unable to process
733 the full GenBank created by prepTG due to an error related to "inconsistent exon ordering".
734 BGCs from each set of genome annotations were independently clustered using BiG-SCAPE
735 with "mix" clustering analysis and MIBiG reference BGC integration requested. The gene cluster
736 family and clan matching the reference leporin BGC in MIBiG (BGC0001445) were regarded as
737 the leporin BGC. For remote cblaster⁴⁵ analysis, CAGECAT⁸⁷ was used to search NCBI's nr
738 database with proteins from the leporin BGC representative (BGC0001445) provided as a
739 query. Only 13 scaffolds, belonging to 12 assemblies (including GCA_000006275.3), were
740 identified.

741

742 **Evolutionary investigations of the *epa* locus across *Enterococcus***

743

744 All *Enterococcus* genomes represented in GTDB R207⁵⁷ (n=5,291) were downloaded using
745 ncbi-genome-download⁵³. The same query for *epa* was used for all analyses. Specifically,
746 coordinates extending from 2,071,671 to 2,115,174 along the *E. faecalis* V583 chromosome,
747 corresponding to genes EF2164 to EF2200, were used as a query for the *epa* locus in fai to
748 identify homologous instances in target genomes^{99,101}.

749

750 Comparing orthology/homology inferences between fai & zol, cblaster & clinker, and
751 OrthoFinder. Representative genome assemblies were selected for each of the 92 species of
752 *Enterococcus* in GTDB R214⁵⁷ based on the N50 metric. One set of species representative
753 genomes corresponded to those with the largest N50 values and the other set was comprised of
754 genomes with the lowest N50 values. The two sets of species representative genomes were
755 processed and investigated identically but independently. Gene calling was first performed for
756 genomes using prepTG with pyrodigal⁵³. To generate the input for OrthoFinder, proteins from
757 prepTG's genome-wide GenBank files were extracted in FASTA format. After, OrthoFinder was
758 run with default settings. Phylogenetic hierarchical orthogroups inferred by OrthoFinder were
759 used for comparisons. To perform gene cluster specific homology prediction with cblaster and
760 clinker, we first used cblaster makedb to convert the genome-wide GenBank files from prepTG
761 into a database that could be searched with cblaster search. cblaster search was run using the
762 criteria: (i) DIAMOND alignment sensitivity mode set to very-sensitive, (ii) the percentage of
763 query genes required to be present in a cluster set to 25%, (iii) 1e-10 as the maximum E-value
764 for protein hits to be considered, (iv) 0% as the minimum coverage for protein hits to be
765 considered, (v) 0% as the minimum identity for protein hits to be considered, (vi) the maximum
766 flanking context for the gene cluster to gather set to 0 bp, (vii) request for intergenic proteins to
767 be included, and (viii) a maximum of 4620 bp allowed to separate protein hits for them to be
768 considered as part of the same gene cluster, which should approximately correspond to the
769 aggregate length of 5 bacterial genes on average¹⁴⁷. Next, cblaster extract_clusters was used to
770 extract gene clusters found in target genomes by cblaster in GenBank format and provide them
771 as input for clinker. clinker was run using default settings but with only an output and matrix
772 output file requested to cut time needed to render an interactive figure, its primary intended
773 result file. To aid appropriate comparisons in orthology prediction, fai was largely run using
774 similar criteria as cblaster search: (i) DIAMOND alignment sensitivity mode set to very-sensitive,
775 (ii) the percentage of query genes required to be present in a cluster set to 25%, (iii) 1e-10 as
776 the maximum E-value for protein hits to be considered, (iv) the maximum flanking context for the
777 gene cluster to gather set to 0 bp, (v) a maximum of 5 proteins allowed to separate hits for
778 them to be considered as part of the same gene cluster, and (vi) syntenic similarity assessment
779 between target gene clusters and the query gene cluster turned off. However, draft-mode was
780 enabled in fai, which is not available in cblaster, to showcase the program's ability to improve
781 sensitivity for draft-quality assemblies. zol was applied with mostly default settings but with the
782 flags "only-orthologs", to stop after it determined ortholog groups, and "allow_edge_cds", to
783 allow usage of CDS features marked by fai to be near scaffold edges. All three methods were
784 provided 20 threads wherever possible.

785

786 Comprehensive and tailored usages of fai and zol for finding epa in Enterococcus: Based on
787 prior comparative analyses that had shown that gene conservation and gene order can be

788 slightly variable between *epa* loci from *E. faecalis* and *E. faecium*^{94,95}, we relaxed the syntenic
789 similarity requirement of candidate gene cluster matches in target genomes to the query in fai
790 from 0.6 to 0.0. In addition, we relaxed the minimum percentage of query proteins needed to
791 report a homologous instance of the *epa* locus to 10%. Instead, we required the presence of
792 50% of key *epa* proteins found in both *E. faecalis* and *E. faecium*, defined as
793 *epaABCDEFGHIJLMOPQR*, for the identification of valid homologous instances of the *epa* locus.
794 The E-value cutoff to determine presence for the key *epa* proteins was lowered from 1e-20 to
795 1e-10 to be inclusive of shorter genes and allow for higher levels of sequence divergence
796 across the *Enterococcus* genus. To gather auxiliary genes flanking the core *epa* region in target
797 genomes, we further requested the inclusion of CDS features found within 20 kb of the
798 boundary genes in detected instances of the *epa* locus within the resulting GenBank files
799 produced by fai. A phylogenetic heatmap was constructed for the presence of the *epa* locus
800 across a species tree using species representative genomes, selected based on largest
801 assembly N50, where the values of the heatmap corresponded to the maximum percent identity
802 of a query protein to their best match in target genomes. Because EF2173 and EF2185 are
803 identical transposases, they were shown as one column in the heatmap. The species tree was
804 constructed using GToTree¹⁴⁸ using HMMs for proteins regarded as largely single-copy core to
805 the phylum Bacillota. The phylogenetic heatmap visual was created using iTol¹⁴⁹.

806 From inspection of fai's resulting XLSX spreadsheet, zol's parameters were adjusted to
807 relax identity and coverage thresholds for assessing protein pairs for orthology prior to MCL
808 clustering to 20% and 25%, respectively. Identical processing was performed for the full set of
809 *epa* loci and *epa* loci from only species representative genomes. During the comprehensive
810 processing of all high-quality *epa* loci identified, one instance was dropped during zol analysis
811 despite meeting requirements because all CDS features in it were found near scaffold edges
812 and, by default, such features are not used in zol to aid more accurate inference of ortholog
813 groups and assessment of their sequence variation. A third run of zol was performed using
814 identical settings and all the gene cluster instances but leveraging the dereplication and
815 reinflation options to showcase how the combination of the options can reduce the runtime
816 needed for comprehensive processing. For dereplication of gene clusters, alignment fraction
817 was increased from the default of 95% to 99% and MCL was used for clustering to gather more
818 resolute representative gene clusters. Major ortholog groups determined between the
819 comprehensive and the dereplication + reinflation runs were found to be similarly conserved
820 based on matching to known *epa* genes.

821
822 Phylogenetic assessment of glycosyltransferase orthology predictions: Proteins from ortholog
823 groups determined by zol analysis of species representative genomes were extracted based on
824 whether the ortholog group was annotated as featuring the keywords: "glycosyl" and
825 "transferase" in Pfam protein domain annotations¹⁵⁰. Two additional ortholog groups were
826 included and featured the Pfam domain "Bacterial sugar transferase", including *epaR*, which is
827 also regarded as a glycosyltransferase¹⁰¹. The comprehensive set of glycosyltransferases were
828 next aligned using MUSCLE with the default align mode¹³². Filtering of the alignment was next
829 performed using trimal with options "-keepseqs -gt 0.9" to filter sites composed largely of gaps
830 and further filtered for sequences which were composed of >10% gaps or ambiguous
831 characters ("X"). IQ-TREE¹⁵¹ was used to construct a maximum-likelihood phylogeny with

832 ModelFinder limited to the WAG and LG substitution models. The phylogeny was visualized
833 using iTol¹⁴⁹ with classifications for ortholog groups most closely matching *E. faecalis* V583 *epa*
834 glycosyltransferases marked on leaves. Ortholog groups were assigned to specific *epa* gene
835 designations based on sequence alignment of their consensus sequences to *E. faecalis* V583
836 *epa*-associated proteins. Best matching ortholog groups for each *E. faecalis* V583 *epa*
837 glycosyltransferase were identified based on E-value.

838

839 **Large-scale evolutionary investigations of *epa* loci from *E. faecalis***

840

841 The full set of *epa* loci identified by fai in *E. faecalis* genomes were processed through zol
842 requesting for retention of only complete instances that were also distant from scaffold edges.
843 For projection of conservation, Tajima's D, and sequence entropy statistics onto genes for the
844 *epa* locus in *E. faecalis* V583, sequence alignment was used to identify the best matching
845 ortholog groups based on E-value. For the identical transposases, EF2173 and EF2185, data
846 from a common ortholog group was used for both.

847

848 Investigation of glycosyltransferase phylogenetic diversity: A similar phylogeny of
849 glycosyltransferases was constructed for the *E. faecalis* analysis as was done for the
850 investigation of *epa* glycosyltransferases across species representatives of *Enterococcus*.
851 Glycosyltransferase ortholog groups were identified based on Pfam domains featuring the
852 keywords "glycosyl transferase" or because they matched *epa* genes regarded as
853 glycosyltransferases in prior studies¹⁰¹. To accommodate for the larger number of sequences: (i)
854 only ortholog groups found in >1% of *epa* loci instances were regarded, (ii) MUSCLE¹³² super5
855 mode was used for alignment, and (iii) FastTree 2¹²⁷ was used for approximate maximum-
856 likelihood phylogeny construction. After trimal based filtering of sites, only sequences which
857 featured greater than 20% gaps or ambiguous characters ("X") were filtered to retain *epaA*
858 in the final alignment prior to phylogeny construction.

859

860 **Abbreviations**

861

862 Biosynthetic gene cluster (BGC), mobile-genetic element (MGE), Enterococcal polysaccharide
863 antigen (Epa), coding sequence (CDS), average nucleotide identity (ANI), metagenome-
864 assembled genome (MAG).

865

866 **Declarations**

867

868 Ethics approval and consent to participate: Not applicable

869

870 Consent for publication: Not applicable

871

872 Availability of data and materials: All genomic and metagenomic datasets used for showcasing
873 the application of fai and zol are publicly available on NCBI with accessions provided in
874 Supplementary Table S11.

875

876 Competing interests: The authors declare that they have no competing interests.

877

878 Funding: This work was supported by grants from the National Institutes of Health awarded to
879 L.R.K (NIAID U19AI142720 and NIGMS R35GM137828) and the Broad Institute
880 (U19AI110818). The content is solely the responsibility of the authors and does not necessarily
881 represent the official views of the National Institutes of Health.

882

883 Authors' contributions: RS & LRK designed the bioinformatics toolkit. RS, KA, & LRK conceived
884 of showcase applications. RS developed the software and performed all analyses. PQT
885 performed assembly of metagenomes. RS & LRK performed the initial drafting of the
886 manuscript. RS, PQT, CM, ALM, MSG, AME, KA, and LRK performed interpretation of results
887 and revision of the manuscript. All authors read and approved the final manuscript.

888

889 Acknowledgments: The authors are grateful to James Kosmopoulos, Dr. Caitlin Pepperell, Dr.
890 Caitlin Sande, and Dr. Mary Hannah Swaney for feedback and assistance with data acquisition
891 as well as Dr. Devon Ryan and Dr. Robert A. Petit III for assistance with incorporation of the
892 suite into Bioconda.

893

894 **References**

895

- 896 1. Enright, A. J., Kunin, V. & Ouzounis, C. A. Protein families and TRIBES in genome
897 sequence space. *Nucleic Acids Res.* **31**, 4632–4638 (2003).
- 898 2. Li, L., Stoeckert, C. J., Jr & Roos, D. S. OrthoMCL: identification of ortholog groups for
899 eukaryotic genomes. *Genome Res.* **13**, 2178–2189 (2003).
- 900 3. Remm, M., Storm, C. E. & Sonnhammer, E. L. Automatic clustering of orthologs and in-
901 paralogs from pairwise species comparisons. *J. Mol. Biol.* **314**, 1041–1052 (2001).
- 902 4. van Dongen, S. M. Graph clustering by flow simulation. (2000).
- 903 5. Schreiber, F. & Sonnhammer, E. L. L. Hieranoid: hierarchical orthology inference. *J. Mol.*
904 *Biol.* **425**, 2072–2081 (2013).
- 905 6. Georgescu, C. H. *et al.* SynerClust: a highly scalable, synteny-aware orthologue clustering
906 tool. *Microb Genom* **4**, (2018).
- 907 7. Hu, X. & Friedberg, I. SwiftOrtho: A fast, memory-efficient, multiple genome orthology
908 classifier. *Gigascience* **8**, (2019).
- 909 8. Cosentino, S. & Iwasaki, W. SonicParanoid: fast, accurate and easy orthology inference.

- 910 *Bioinformatics* **35**, 149–151 (2019).
- 911 9. Ding, W., Baumdicker, F. & Neher, R. A. panX: pan-genome analysis and exploration.
912 *Nucleic Acids Res.* **46**, e5 (2018).
- 913 10. Page, A. J. *et al.* Roary: rapid large-scale prokaryote pan genome analysis. *Bioinformatics*
914 **31**, 3691–3693 (2015).
- 915 11. Bayliss, S. C., Thorpe, H. A., Coyle, N. M., Sheppard, S. K. & Feil, E. J. PIRATE: A fast and
916 scalable pangenomics toolbox for clustering diverged orthologues in bacteria. *bioRxiv*
917 (2019) doi:10.1101/598391.
- 918 12. Tonkin-Hill, G. *et al.* Producing polished prokaryotic pangenomes with the Panaroo
919 pipeline. *Genome Biol.* **21**, 180 (2020).
- 920 13. Gautreau, G. *et al.* PPanGGOLiN: Depicting microbial diversity via a partitioned
921 pangenome graph. *PLoS Comput. Biol.* **16**, e1007732 (2020).
- 922 14. Aramaki, T. *et al.* KofamKOALA: KEGG Ortholog assignment based on profile HMM and
923 adaptive score threshold. *Bioinformatics* **36**, 2251–2252 (2019).
- 924 15. Cantalapiedra, C. P., Hernández-Plaza, A., Letunic, I., Bork, P. & Huerta-Cepas, J.
925 eggNOG-mapper v2: Functional Annotation, Orthology Assignments, and Domain
926 Prediction at the Metagenomic Scale. *Mol. Biol. Evol.* **38**, 5825–5829 (2021).
- 927 16. Tatusov, R. L., Galperin, M. Y., Natale, D. A. & Koonin, E. V. The COG database: a tool for
928 genome-scale analysis of protein functions and evolution. *Nucleic Acids Res.* **28**, 33–36
929 (2000).
- 930 17. Melnyk, R. A., Hossain, S. S. & Haney, C. H. Convergent gain and loss of genomic islands
931 drive lifestyle changes in plant-associated *Pseudomonas*. *ISME J.* **13**, 1575–1588 (2019).
- 932 18. Steinegger, M. & Söding, J. MMseqs2 enables sensitive protein sequence searching for the
933 analysis of massive data sets. *Nat. Biotechnol.* **35**, 1026–1028 (2017).
- 934 19. Buchfink, B., Ashkenazy, H., Reuter, K., Kennedy, J. A. & Drost, H.-G. Sensitive clustering
935 of protein sequences at tree-of-life scale using DIAMOND DeepClust. *bioRxiv* (2023)

- 936 doi:10.1101/2023.01.24.525373.
- 937 20. Coelho, L. P. *et al.* Towards the biogeography of prokaryotic genes. *Nature* **601**, 252–256
938 (2022).
- 939 21. Steinegger, M. & Söding, J. Clustering huge protein sequence sets in linear time. *Nat.*
940 *Commun.* **9**, 2542 (2018).
- 941 22. Snyder, L., Henkin, T. M., Peters, J. E. & Champness, W. Molecular Genetics of Bacteria,
942 4th Edition. Preprint at <https://doi.org/10.1128/9781555817169> (2013).
- 943 23. Price, M. N., Arkin, A. P. & Alm, E. J. The life-cycle of operons. *PLoS Genet.* **2**, e96 (2006).
- 944 24. Ptashne, M. *A Genetic Switch: Gene Control and Phage. Lambda.* (Palo Alto, CA (US);
945 Blackwell Scientific Publications, 1986).
- 946 25. Andreu, V. P. *et al.* gutSMASH predicts specialized primary metabolic pathways from the
947 human gut microbiota. *Nature Biotechnology* Preprint at [https://doi.org/10.1038/s41587-](https://doi.org/10.1038/s41587-023-01675-1)
948 [023-01675-1](https://doi.org/10.1038/s41587-023-01675-1) (2023).
- 949 26. Cortes, J., Haydock, S. F., Roberts, G. A., Bevitt, D. J. & Leadlay, P. F. An unusually large
950 multifunctional polypeptide in the erythromycin-producing polyketide synthase of
951 *Saccharopolyspora erythraea*. *Nature* **348**, 176–178 (1990).
- 952 27. Donadio, S., Staver, M. J., McAlpine, J. B., Swanson, S. J. & Katz, L. Modular organization
953 of genes required for complex polyketide biosynthesis. *Science* **252**, 675–679 (1991).
- 954 28. Walsh, C. T. & Fischbach, M. A. Natural products version 2.0: connecting genes to
955 molecules. *J. Am. Chem. Soc.* **132**, 2469–2493 (2010).
- 956 29. Medema, M. H. *et al.* antiSMASH: rapid identification, annotation and analysis of secondary
957 metabolite biosynthesis gene clusters in bacterial and fungal genome sequences. *Nucleic*
958 *Acids Res.* **39**, W339-46 (2011).
- 959 30. Gal-Mor, O. & Finlay, B. B. Pathogenicity islands: a molecular toolbox for bacterial
960 virulence. *Cell. Microbiol.* **8**, 1707–1719 (2006).
- 961 31. Kaper, J. B., Nataro, J. P. & Mobley, H. L. Pathogenic *Escherichia coli*. *Nat. Rev. Microbiol.*

- 962 **2**, 123–140 (2004).
- 963 32. Bolwell, G. P. & Paul Bolwell, G. *Biochemistry & Molecular Biology of Plants*.
964 *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)
965 (2001).
- 966 33. Rokas, A., Mead, M. E., Steenwyk, J. L., Raja, H. A. & Oberlies, N. H. Biosynthetic gene
967 clusters and the evolution of fungal chemodiversity. *Nat. Prod. Rep.* **37**, 868–878 (2020).
- 968 34. Robey, M. T., Caesar, L. K., Drott, M. T., Keller, N. P. & Kelleher, N. L. An interpreted atlas
969 of biosynthetic gene clusters from 1,000 fungal genomes. *Proc. Natl. Acad. Sci. U. S. A.*
970 **118**, (2021).
- 971 35. Lindahl, L. & Zengel, J. M. Operon-specific regulation of ribosomal protein synthesis in
972 *Escherichia coli*. *Proc. Natl. Acad. Sci. U. S. A.* **76**, 6542–6546 (1979).
- 973 36. Cordero, O. X. & Polz, M. F. Explaining microbial genomic diversity in light of evolutionary
974 ecology. *Nat. Rev. Microbiol.* **12**, 263–273 (2014).
- 975 37. Drott, M. T. *et al.* Microevolution in the pansecondary metabolome of *Aspergillus flavus* and
976 its potential macroevolutionary implications for filamentous fungi. *Proc. Natl. Acad. Sci. U.*
977 *S. A.* **118**, (2021).
- 978 38. Salamzade, R. *et al.* Evolutionary investigations of the biosynthetic diversity in the skin
979 microbiome using IsaBGC. *Microb Genom* **9**, (2023).
- 980 39. Ziemert, N. *et al.* Diversity and evolution of secondary metabolism in the marine
981 actinomycete genus *Salinispora*. *Proc. Natl. Acad. Sci. U. S. A.* **111**, E1130-9 (2014).
- 982 40. van Bergeijk, D. A., Terlouw, B. R., Medema, M. H. & van Wezel, G. P. Ecology and
983 genomics of Actinobacteria: new concepts for natural product discovery. *Nat. Rev.*
984 *Microbiol.* **18**, 546–558 (2020).
- 985 41. Chevrette, M. G. *et al.* Evolutionary dynamics of natural product biosynthesis in bacteria.
986 *Nat. Prod. Rep.* **37**, 566–599 (2020).
- 987 42. Medema, M. H., Takano, E. & Breitling, R. Detecting sequence homology at the gene

- 988 cluster level with MultiGeneBlast. *Mol. Biol. Evol.* **30**, 1218–1223 (2013).
- 989 43. Abby, S. S., Néron, B., Ménager, H., Touchon, M. & Rocha, E. P. C. MacSyFinder: a
990 program to mine genomes for molecular systems with an application to CRISPR-Cas
991 systems. *PLoS One* **9**, e110726 (2014).
- 992 44. Navarro-Muñoz, J. C. *et al.* A computational framework to explore large-scale biosynthetic
993 diversity. *Nat. Chem. Biol.* **16**, 60–68 (2020).
- 994 45. Gilchrist, C. L. M. *et al.* Cblaster: A remote search tool for rapid identification and
995 visualization of homologous gene clusters. *Bioinformatics Advances* **1**, (2021).
- 996 46. Gilchrist, C. L. M. & Chooi, Y.-H. clinker & clustermap.js: automatic generation of gene
997 cluster comparison figures. *Bioinformatics* **37**, 2473–2475 (2021).
- 998 47. Hackl, T. & Ankenbrand, M. J. gggenomes: a grammar of graphics for comparative
999 genomics. *R package version 0.9*.
- 1000 48. moshi. *PyGenomeViz: A Genome Visualization Python Package for Comparative*
1001 *Genomics*. (Github).
- 1002 49. Tajima, F. Statistical method for testing the neutral mutation hypothesis by DNA
1003 polymorphism. *Genetics* **123**, 585–595 (1989).
- 1004 50. Graziotin, A. L., Koonin, E. V. & Kristensen, D. M. Prokaryotic Virus Orthologous Groups
1005 (pVOGs): a resource for comparative genomics and protein family annotation. *Nucleic*
1006 *Acids Res.* **45**, D491–D498 (2017).
- 1007 51. Liu, B., Zheng, D., Jin, Q., Chen, L. & Yang, J. VFDB 2019: a comparative pathogenomic
1008 platform with an interactive web interface. *Nucleic Acids Res.* **47**, D687–D692 (2019).
- 1009 52. Terlouw, B. R. *et al.* MIBiG 3.0: a community-driven effort to annotate experimentally
1010 validated biosynthetic gene clusters. *Nucleic Acids Res.* **51**, D603–D610 (2023).
- 1011 53. Larralde, M. Pyrodigal: Python bindings and interface to Prodigal, an efficient method for
1012 gene prediction in prokaryotes. *J. Open Source Softw.* **7**, 4296 (2022).
- 1013 54. Hyatt, D. *et al.* Prodigal: prokaryotic gene recognition and translation initiation site

- 1014 identification. *BMC Bioinformatics* **11**, 119 (2010).
- 1015 55. Li, H. Protein-to-genome alignment with miniprot. *Bioinformatics* **39**, (2023).
- 1016 56. Salamzade, R. & Kalan, L. R. skDER: microbial genome dereplication approaches for
1017 comparative and metagenomic applications. *bioRxiv.org* (2023)
1018 doi:10.1101/2023.09.27.559801.
- 1019 57. Parks, D. H. *et al.* GTDB: an ongoing census of bacterial and archaeal diversity through a
1020 phylogenetically consistent, rank normalized and complete genome-based taxonomy.
1021 *Nucleic Acids Res.* (2021) doi:10.1093/nar/gkab776.
- 1022 58. Liu, M. *et al.* ICEberg 2.0: an updated database of bacterial integrative and conjugative
1023 elements. *Nucleic Acids Res.* **47**, D660–D665 (2019).
- 1024 59. Bertelli, C. *et al.* IslandViewer 4: expanded prediction of genomic islands for larger-scale
1025 datasets. *Nucleic Acids Res.* **45**, W30–W35 (2017).
- 1026 60. Siguier, P., Perochon, J., Lestrade, L., Mahillon, J. & Chandler, M. ISfinder: the reference
1027 centre for bacterial insertion sequences. *Nucleic Acids Res.* **34**, D32–6 (2006).
- 1028 61. Alcock, B. P. *et al.* CARD 2023: expanded curation, support for machine learning, and
1029 resistome prediction at the Comprehensive Antibiotic Resistance Database. *Nucleic Acids*
1030 *Res.* **51**, D690–D699 (2023).
- 1031 62. Kosakovsky Pond, S. L. *et al.* HyPhy 2.5—A Customizable Platform for Evolutionary
1032 Hypothesis Testing Using Phylogenies. *Mol. Biol. Evol.* **37**, 295–299 (2019).
- 1033 63. Hackl, T., Duponchel, S., Barenhoff, K., Weinmann, A. & Fischer, M. G. Virophages and
1034 retrotransposons colonize the genomes of a heterotrophic flagellate. *Elife* **10**, (2021).
- 1035 64. Sullivan, M. J., Petty, N. K. & Beatson, S. A. Easyfig: a genome comparison visualizer.
1036 *Bioinformatics* **27**, 1009–1010 (2011).
- 1037 65. Shaw, J. & Yu, Y. W. Fast and robust metagenomic sequence comparison through sparse
1038 chaining with skani. *Nat. Methods* **20**, 1661–1665 (2023).
- 1039 66. Blackwell, G. *et al.* Exploring bacterial diversity via a curated and searchable snapshot of

- 1040 archived DNA. *Access Microbiol.* **4**, (2022).
- 1041 67. Lebreton, F. *et al.* Emergence of epidemic multidrug-resistant *Enterococcus faecium* from
1042 animal and commensal strains. *MBio* **4**, (2013).
- 1043 68. Lieberman, T. D. *et al.* Genetic variation of a bacterial pathogen within individuals with
1044 cystic fibrosis provides a record of selective pressures. *Nat. Genet.* **46**, 82–87 (2014).
- 1045 69. Crits-Christoph, A., Olm, M. R., Diamond, S., Bouma-Gregson, K. & Banfield, J. F. Soil
1046 bacterial populations are shaped by recombination and gene-specific selection across a
1047 grassland meadow. *ISME J.* **14**, 1834–1846 (2020).
- 1048 70. Hudson, R. R., Slatkin, M. & Maddison, W. P. Estimation of levels of gene flow from DNA
1049 sequence data. *Genetics* **132**, 583–589 (1992).
- 1050 71. Pavlopoulos, G. A. *et al.* Unraveling the functional dark matter through global
1051 metagenomics. *Nature* **622**, 594–602 (2023).
- 1052 72. Vanni, C. *et al.* Unifying the known and unknown microbial coding sequence space. *Elife*
1053 **11**, (2022).
- 1054 73. Tran, P. Q. *et al.* Viral impacts on microbial activity and biogeochemical cycling in a
1055 seasonally anoxic freshwater lake. *bioRxiv* 2023.04.19.537559 (2023)
1056 doi:10.1101/2023.04.19.537559.
- 1057 74. Kieft, K., Zhou, Z. & Anantharaman, K. VIBRANT: automated recovery, annotation and
1058 curation of microbial viruses, and evaluation of viral community function from genomic
1059 sequences. *Microbiome* **8**, 90 (2020).
- 1060 75. Willems, A. The Family Comamonadaceae. in *The Prokaryotes* 777–851 (Springer Berlin
1061 Heidelberg, Berlin, Heidelberg, 2014).
- 1062 76. Roux, S. *et al.* iPHoP: An integrated machine learning framework to maximize host
1063 prediction for metagenome-derived viruses of archaea and bacteria. *PLoS Biol.* **21**,
1064 e3002083 (2023).
- 1065 77. Klassen, J. L. & Currie, C. R. Gene fragmentation in bacterial draft genomes: extent,

- 1066 consequences and mitigation. *BMC Genomics* **13**, 14 (2012).
- 1067 78. Thomma, B. P. H. J. *et al.* Mind the gap; seven reasons to close fragmented genome
1068 assemblies. *Fungal Genet. Biol.* **90**, 24–30 (2016).
- 1069 79. Drăgan, M.-A., Moghul, I., Priyam, A., Bustos, C. & Wurm, Y. GeneValidator: identify
1070 problems with protein-coding gene predictions. *Bioinformatics* **32**, 1559–1561 (2016).
- 1071 80. Scalzitti, N., Jeannin-Girardon, A., Collet, P., Poch, O. & Thompson, J. D. A benchmark
1072 study of ab initio gene prediction methods in diverse eukaryotic organisms. *BMC Genomics*
1073 **21**, 293 (2020).
- 1074 81. Jallow, A., Xie, H., Tang, X., Qi, Z. & Li, P. Worldwide aflatoxin contamination of agricultural
1075 products and foods: From occurrence to control. *Compr. Rev. Food Sci. Food Saf.* **20**,
1076 2332–2381 (2021).
- 1077 82. Bok, J. W. *et al.* Genomic mining for *Aspergillus* natural products. *Chem. Biol.* **13**, 31–37
1078 (2006).
- 1079 83. Vadlapudi, V. *et al.* *Aspergillus* Secondary Metabolite Database, a resource to understand
1080 the Secondary metabolome of *Aspergillus* genus. *Sci. Rep.* **7**, 7325 (2017).
- 1081 84. Hatmaker, E. A. *et al.* Genomic and Phenotypic Trait Variation of the Opportunistic Human
1082 Pathogen *Aspergillus flavus* and Its Close Relatives. *Microbiol Spectr* **10**, e0306922 (2022).
- 1083 85. Cary, J. W. *et al.* An *Aspergillus flavus* secondary metabolic gene cluster containing a
1084 hybrid PKS-NRPS is necessary for synthesis of the 2-pyridones, leporins. *Fungal Genet.*
1085 *Biol.* **81**, 88–97 (2015).
- 1086 86. Skerker, J. M. *et al.* Chromosome assembled and annotated genome sequence of
1087 *Aspergillus flavus* NRRL 3357. *G3* **11**, jkab213 (2021).
- 1088 87. van den Belt, M. *et al.* CAGECAT: The CompArative GEne Cluster Analysis Toolbox for
1089 rapid search and visualisation of homologous gene clusters. *BMC Bioinformatics* **24**, 181
1090 (2023).
- 1091 88. Majoros, W. H., Pertea, M. & Salzberg, S. L. TigrScan and GlimmerHMM: two open source

- 1092 ab initio eukaryotic gene-finders. *Bioinformatics* **20**, 2878–2879 (2004).
- 1093 89. Yang, K., Tian, J. & Keller, N. P. Post-translational modifications drive secondary
1094 metabolite biosynthesis in *Aspergillus*: a review. *Environ. Microbiol.* **24**, 2857–2881 (2022).
- 1095 90. Klich, M. A. *Aspergillus flavus*: the major producer of aflatoxin. *Mol. Plant Pathol.* **8**, 713–
1096 722 (2007).
- 1097 91. Cary, J. W., Ehrlich, K. C., Bland, J. M. & Montalbano, B. G. The Aflatoxin Biosynthesis
1098 Cluster Gene, *afIX*, Encodes an Oxidoreductase Involved in Conversion of Versicolorin A
1099 to Demethylsterigmatocystin. *Applied and Environmental Microbiology* vol. 72 1096–1101
1100 Preprint at <https://doi.org/10.1128/aem.72.2.1096-1101.2006> (2006).
- 1101 92. Cleveland, T. E. *et al.* Potential of *Aspergillus flavus* genomics for applications in
1102 biotechnology. *Trends Biotechnol.* **27**, 151–157 (2009).
- 1103 93. Ehrlich, K. C., Li, P., Scharfenstein, L. & Chang, P.-K. HypC, the anthrone oxidase involved
1104 in aflatoxin biosynthesis. *Appl. Environ. Microbiol.* **76**, 3374–3377 (2010).
- 1105 94. Palmer, K. L. *et al.* Comparative Genomics of Enterococci: Variation in *Enterococcus*
1106 *faecalis*, Clade Structure in *E. faecium*, and Defining Characteristics of *E. gallinarum* and *E.*
1107 *casseliflavus*. *mBio* vol. 3 Preprint at <https://doi.org/10.1128/mbio.00318-11> (2012).
- 1108 95. Qin, X. *et al.* Complete genome sequence of *Enterococcus faecium* strain TX16 and
1109 comparative genomic analysis of *Enterococcus faecium* genomes. *BMC Microbiol.* **12**, 135
1110 (2012).
- 1111 96. Xu, Y., Murray, B. E. & Weinstock, G. M. A cluster of genes involved in polysaccharide
1112 biosynthesis from *Enterococcus faecalis* OG1RF. *Infect. Immun.* **66**, 4313–4323 (1998).
- 1113 97. Hancock, L. E., Murray, B. E. & Sillanpää, J. Enterococcal Cell Wall Components and
1114 Structures. in *Enterococci: From Commensals to Leading Causes of Drug Resistant*
1115 *Infection* (eds. Gilmore, M. S., Clewell, D. B., Ike, Y. & Shankar, N.) (Massachusetts Eye
1116 and Ear Infirmary, Boston, 2014).
- 1117 98. Teng, F., Jacques-Palaz, K. D., Weinstock, G. M. & Murray, B. E. Evidence that the

- 1118 Enterococcal Polysaccharide Antigen Gene (*epa*) Cluster Is Widespread in *Enterococcus*
1119 *faecalis* and Influences Resistance to Phagocytic Killing of *E. faecalis*. *Infection and*
1120 *Immunity* vol. 70 2010–2015 Preprint at <https://doi.org/10.1128/iai.70.4.2010-2015.2002>
1121 (2002).
- 1122 99. Teng, F., Singh, K. V., Bourgogne, A., Zeng, J. & Murray, B. E. Further Characterization of
1123 the *epa* Gene Cluster and *Epa* Polysaccharides of *Enterococcus faecalis*. *Infection and*
1124 *Immunity* vol. 77 3759–3767 Preprint at <https://doi.org/10.1128/iai.00149-09> (2009).
- 1125 100. Rigottier-Gois, L. *et al.* The surface rhamnopolysaccharide *epa* of *Enterococcus faecalis* is
1126 a key determinant of intestinal colonization. *J. Infect. Dis.* **211**, 62–71 (2015).
- 1127 101. Guerardel, Y. *et al.* Complete structure of the enterococcal polysaccharide antigen (EPA) of
1128 vancomycin-resistant *Enterococcus faecalis* V583 reveals that EPA decorations are
1129 teichoic acids covalently linked to a rhamnopolysaccharide backbone. *MBio* **11**, (2020).
- 1130 102. Smith, R. E. *et al.* Decoration of the enterococcal polysaccharide antigen EPA is essential
1131 for virulence, cell surface charge and interaction with effectors of the innate immune
1132 system. *PLoS Pathog.* **15**, e1007730 (2019).
- 1133 103. Singh, K. V. & Murray, B. E. Loss of a Major Enterococcal Polysaccharide Antigen (*Epa*) by
1134 *Enterococcus faecalis* Is Associated with Increased Resistance to Ceftriaxone and
1135 Carbapenems. *Antimicrob. Agents Chemother.* **63**, (2019).
- 1136 104. Ho, K., Huo, W., Pas, S., Dao, R. & Palmer, K. L. Loss-of-Function Mutations in *epaR*
1137 Confer Resistance to \square NPV1 Infection in *Enterococcus faecalis* OG1RF. *Antimicrobial*
1138 *Agents and Chemotherapy* vol. 62 Preprint at <https://doi.org/10.1128/aac.00758-18> (2018).
- 1139 105. Fiore, E., Van Tyne, D. & Gilmore, M. S. Pathogenicity of Enterococci. *Microbiol Spectr* **7**,
1140 (2019).
- 1141 106. Lebreton, F., Willems, R. J. L. & Gilmore, M. S. *Enterococcus* Diversity, Origins in Nature,
1142 and Gut Colonization. in *Enterococci: From Commensals to Leading Causes of Drug*
1143 *Resistant Infection* (eds. Gilmore, M. S., Clewell, D. B., Ike, Y. & Shankar, N.)

- 1144 (Massachusetts Eye and Ear Infirmary, Boston, 2014).
- 1145 107. Lebreton, F. *et al.* Tracing the Enterococci from Paleozoic Origins to the Hospital. *Cell* **169**,
- 1146 849-861.e13 (2017).
- 1147 108. Schwartzman, J. A. *et al.* Global diversity of enterococci and description of 18 novel
- 1148 species. *bioRxiv* 2023.05.18.540996 (2023) doi:10.1101/2023.05.18.540996.
- 1149 109. Bertelli, C. *et al.* Enabling genomic island prediction and comparison in multiple genomes to
- 1150 investigate bacterial evolution and outbreaks. *Microb. Genom.* **8**, (2022).
- 1151 110. Bin Jang, H. *et al.* Taxonomic assignment of uncultivated prokaryotic virus genomes is
- 1152 enabled by gene-sharing networks. *Nat. Biotechnol.* **37**, 632–639 (2019).
- 1153 111. Meyer, F. *et al.* Critical Assessment of Metagenome Interpretation: the second round of
- 1154 challenges. *Nat. Methods* **19**, 429–440 (2022).
- 1155 112. Salamzade, R. *et al.* Inter-species geographic signatures for tracing horizontal gene
- 1156 transfer and long-term persistence of carbapenem resistance. *Genome Med.* **14**, 37 (2022).
- 1157 113. Sheppard, A. E. *et al.* Nested Russian Doll-Like Genetic Mobility Drives Rapid
- 1158 Dissemination of the Carbapenem Resistance Gene blaKPC. *Antimicrob. Agents*
- 1159 *Chemother.* **60**, 3767–3778 (2016).
- 1160 114. Groussin, M. *et al.* Elevated rates of horizontal gene transfer in the industrialized human
- 1161 microbiome. *Cell* **184**, 2053-2067.e18 (2021).
- 1162 115. Crits-Christoph, A., Diamond, S., Butterfield, C. N., Thomas, B. C. & Banfield, J. F. Novel
- 1163 soil bacteria possess diverse genes for secondary metabolite biosynthesis. *Nature* **558**,
- 1164 440–444 (2018).
- 1165 116. Bickhart, D. M. *et al.* Generation of lineage-resolved complete metagenome-assembled
- 1166 genomes by precision phasing. *bioRxiv* 2021.05.04.442591 (2021)
- 1167 doi:10.1101/2021.05.04.442591.
- 1168 117. Emms, D. M. & Kelly, S. OrthoFinder: phylogenetic orthology inference for comparative
- 1169 genomics. *Genome Biol.* **20**, 238 (2019).

- 1170 118. Nevers, Y. *et al.* The Quest for Orthologs orthology benchmark service in 2022. *Nucleic*
1171 *Acids Res.* (2022).
- 1172 119. Chatterjee, A. *et al.* Bacteriophage Resistance Alters Antibiotic-Mediated Intestinal
1173 Expansion of Enterococci. *Infect. Immun.* **87**, (2019).
- 1174 120. Chatterjee, A. *et al.* Parallel genomics uncover novel enterococcal-bacteriophage
1175 interactions. Preprint at <https://doi.org/10.1101/858506>.
- 1176 121. Canfield, G. S. *et al.* Lytic bacteriophages facilitate antibiotic sensitization of *Enterococcus*
1177 *faecium*. Preprint at <https://doi.org/10.1101/2020.09.22.309401>.
- 1178 122. Kirsch, J. M. *et al.* Targeted IS-element sequencing uncovers transposition dynamics
1179 during selective pressure in enterococci. *PLoS Pathog.* **19**, e1011424 (2023).
- 1180 123. Grüning, B. *et al.* Bioconda: sustainable and comprehensive software distribution for the life
1181 sciences. *Nat. Methods* **15**, 475–476 (2018).
- 1182 124. Salamzade, R. & Kalan, L. *Zol.* (Zenodo, 2024). doi:10.5281/ZENODO.10828137.
- 1183 125. Cock, P. J. A. *et al.* Biopython: freely available Python tools for computational molecular
1184 biology and bioinformatics. *Bioinformatics* **25**, 1422–1423 (2009).
- 1185 126. Capella-Gutiérrez, S., Silla-Martínez, J. M. & Gabaldón, T. trimAl: a tool for automated
1186 alignment trimming in large-scale phylogenetic analyses. *Bioinformatics* **25**, 1972–1973
1187 (2009).
- 1188 127. Price, M. N., Dehal, P. S. & Arkin, A. P. FastTree 2--approximately maximum-likelihood
1189 trees for large alignments. *PLoS One* **5**, e9490 (2010).
- 1190 128. Huang, Y., Niu, B., Gao, Y., Fu, L. & Li, W. CD-HIT Suite: a web server for clustering and
1191 comparing biological sequences. *Bioinformatics* **26**, 680–682 (2010).
- 1192 129. Eddy, S. R. Accelerated Profile HMM Searches. *PLoS Comput. Biol.* **7**, e1002195 (2011).
- 1193 130. Buchfink, B., Xie, C. & Huson, D. H. Fast and sensitive protein alignment using DIAMOND.
1194 *Nat. Methods* **12**, 59–60 (2014).
- 1195 131. Schreiber, J. Pomegranate: fast and flexible probabilistic modeling in python. *J. Mach.*

- 1196 *Learn. Res.* (2017).
- 1197 132. Edgar, R. C. Muscle5: High-accuracy alignment ensembles enable unbiased assessments
1198 of sequence homology and phylogeny. *Nat. Commun.* **13**, 6968 (2022).
- 1199 133. Blin, K. *et al.* antiSMASH 6.0: improving cluster detection and comparison capabilities.
1200 *Nucleic Acids Res.* **49**, W29–W35 (2021).
- 1201 134. Suyama, M., Torrents, D. & Bork, P. PAL2NAL: robust conversion of protein sequence
1202 alignments into the corresponding codon alignments. *Nucleic Acids Res.* **34**, W609-12
1203 (2006).
- 1204 135. Li, W. *et al.* RefSeq: expanding the Prokaryotic Genome Annotation Pipeline reach with
1205 protein family model curation. *Nucleic Acids Res.* **49**, D1020–D1028 (2021).
- 1206 136. Price, M. N. & Arkin, A. P. PaperBLAST: Text Mining Papers for Information about
1207 Homologs. *mSystems* **2**, (2017).
- 1208 137. Finn, R. D. *et al.* Pfam: the protein families database. *Nucleic Acids Res.* **42**, D222-30
1209 (2014).
- 1210 138. Larralde, M. & Zeller, G. PyHMMER: a Python library binding to HMMER for efficient
1211 sequence analysis. *Bioinformatics* **39**, (2023).
- 1212 139. Murrell, B. *et al.* FUBAR: a fast, unconstrained bayesian approximation for inferring
1213 selection. *Mol. Biol. Evol.* **30**, 1196–1205 (2013).
- 1214 140. Kosakovsky Pond, S. L., Posada, D., Gravenor, M. B., Woelk, C. H. & Frost, S. D. W.
1215 GARD: a genetic algorithm for recombination detection. *Bioinformatics* **22**, 3096–3098
1216 (2006).
- 1217 141. Carroll, L. M. *et al.* Accurate de novo identification of biosynthetic gene clusters with
1218 GECCO. *bioRxiv* 2021.05.03.442509 (2021) doi:10.1101/2021.05.03.442509.
- 1219 142. Akhter, S., Aziz, R. K. & Edwards, R. A. PhiSpy: a novel algorithm for finding prophages in
1220 bacterial genomes that combines similarity- and composition-based strategies. *Nucleic
1221 Acids Res.* **40**, e126 (2012).

- 1222 143. Camargo, A. P. *et al.* Identification of mobile genetic elements with geNomad. *Nat.*
1223 *Biotechnol.* (2023) doi:10.1038/s41587-023-01953-y.
- 1224 144. Robertson, J. & Nash, J. H. E. MOB-suite: software tools for clustering, reconstruction and
1225 typing of plasmids from draft assemblies. *Microb Genom* **4**, (2018).
- 1226 145. Ondov, B. D. *et al.* Mash: fast genome and metagenome distance estimation using
1227 MinHash. *Genome Biol.* **17**, 132 (2016).
- 1228 146. Paradis, E., Claude, J. & Strimmer, K. APE: Analyses of Phylogenetics and Evolution in R
1229 language. *Bioinformatics* **20**, 289–290 (2004).
- 1230 147. Xu, L. *et al.* Average gene length is highly conserved in prokaryotes and eukaryotes and
1231 diverges only between the two kingdoms. *Mol. Biol. Evol.* **23**, 1107–1108 (2006).
- 1232 148. Lee, M. D. GToTree: a user-friendly workflow for phylogenomics. *Bioinformatics* **35**, 4162–
1233 4164 (2019).
- 1234 149. Letunic, I. & Bork, P. Interactive Tree Of Life (iTOL) v4: recent updates and new
1235 developments. *Nucleic Acids Res.* **47**, W256–W259 (2019).
- 1236 150. Mistry, J. *et al.* Pfam: The protein families database in 2021. *Nucleic Acids Res.* **49**, D412–
1237 D419 (2021).
- 1238 151. Minh, B. Q. *et al.* IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference
1239 in the Genomic Era. *Mol. Biol. Evol.* **37**, 1530–1534 (2020).
- 1240

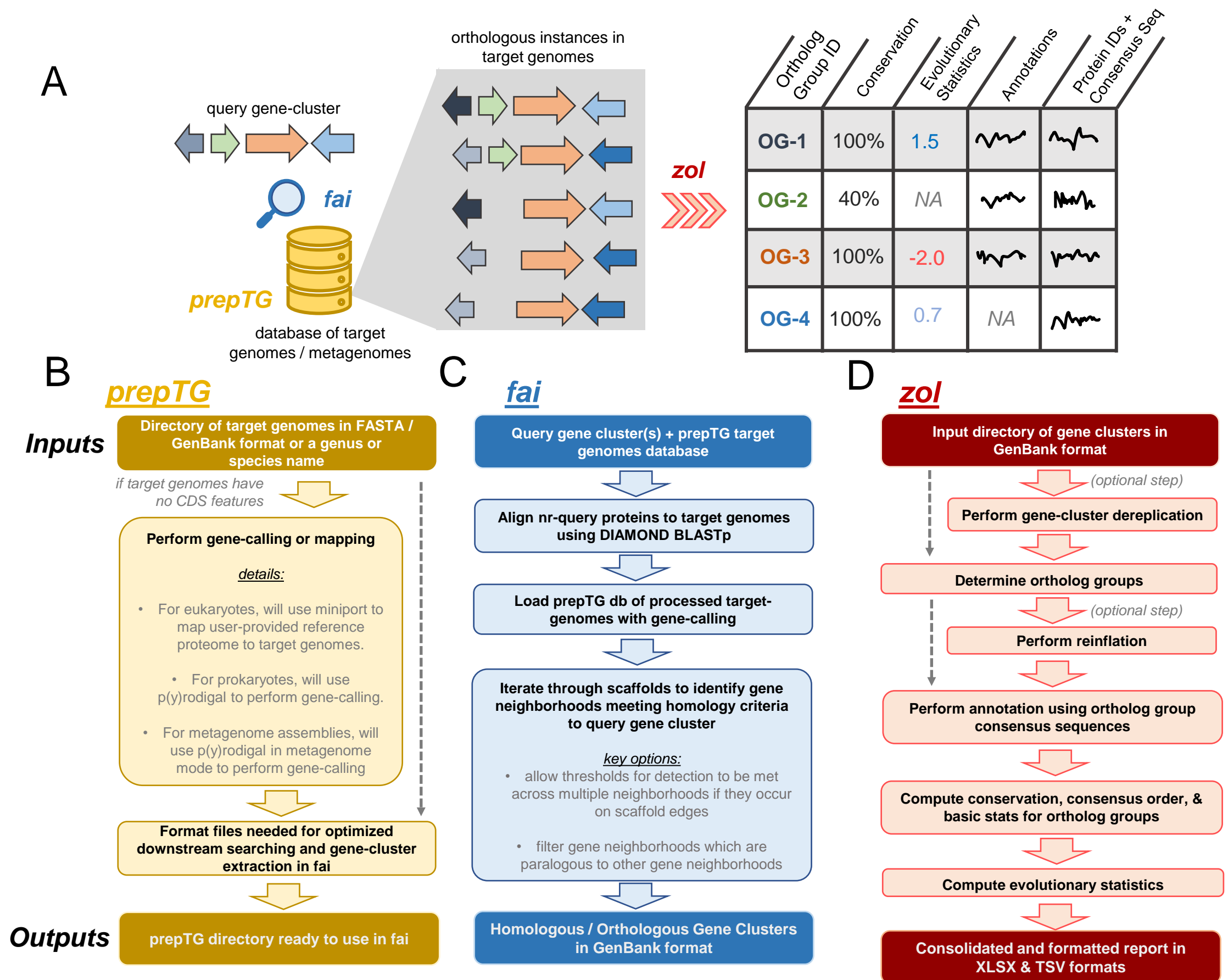
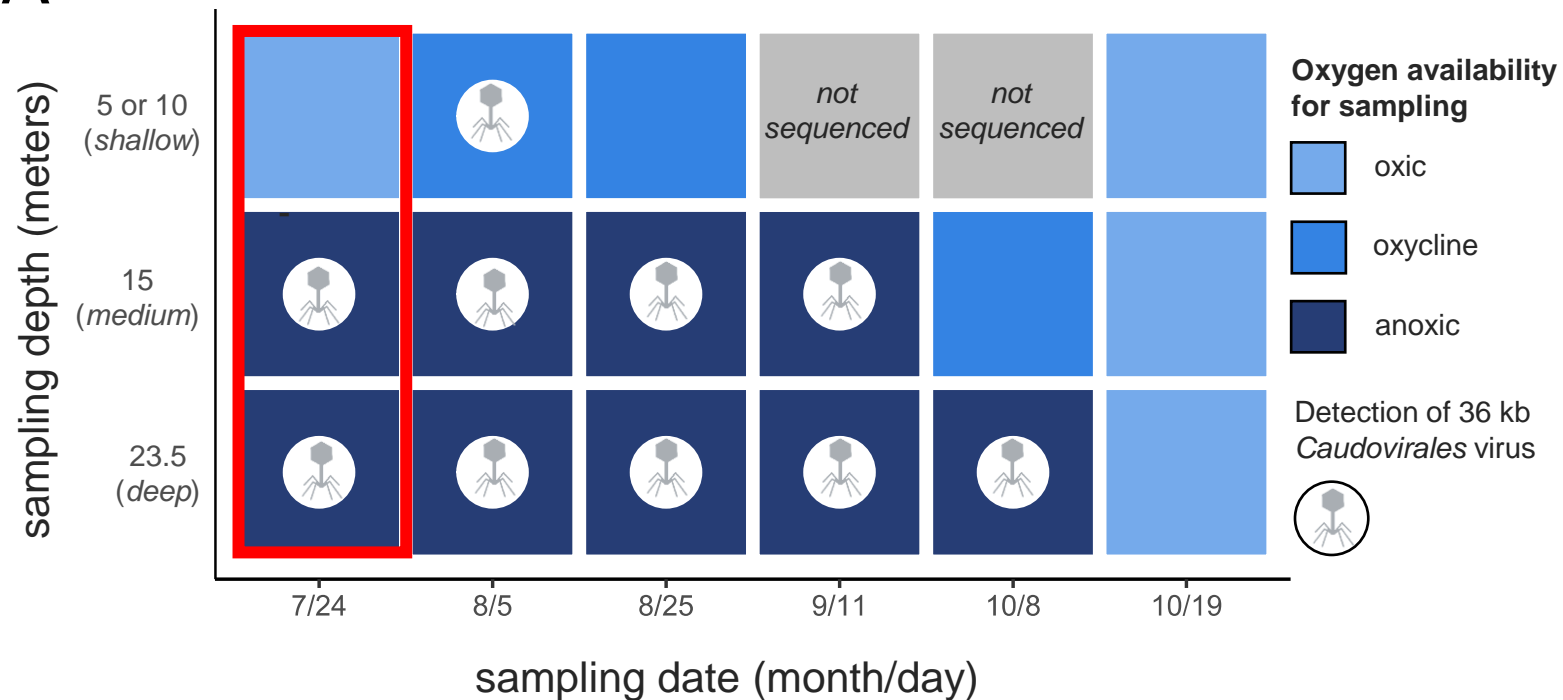


Figure 1: Overviews of fai and zol. **A)** A cartoon schematic of how prepTG, fai, and zol are integrated to perform evolutionary investigations by searching for gene-clusters. Certain statistics in the zol report will not be calculated if not enough instances of an ortholog group are identified, resulting in non-available (NA) values being reported. Squiggles correspond to arbitrary text pertaining to functional annotation information, etc. **B)** An overview of the prepTG, **C)** fai, and **D)** zol algorithms and workflows. Inputs and outputs for the programs are indicated with bolder coloring.

A



B

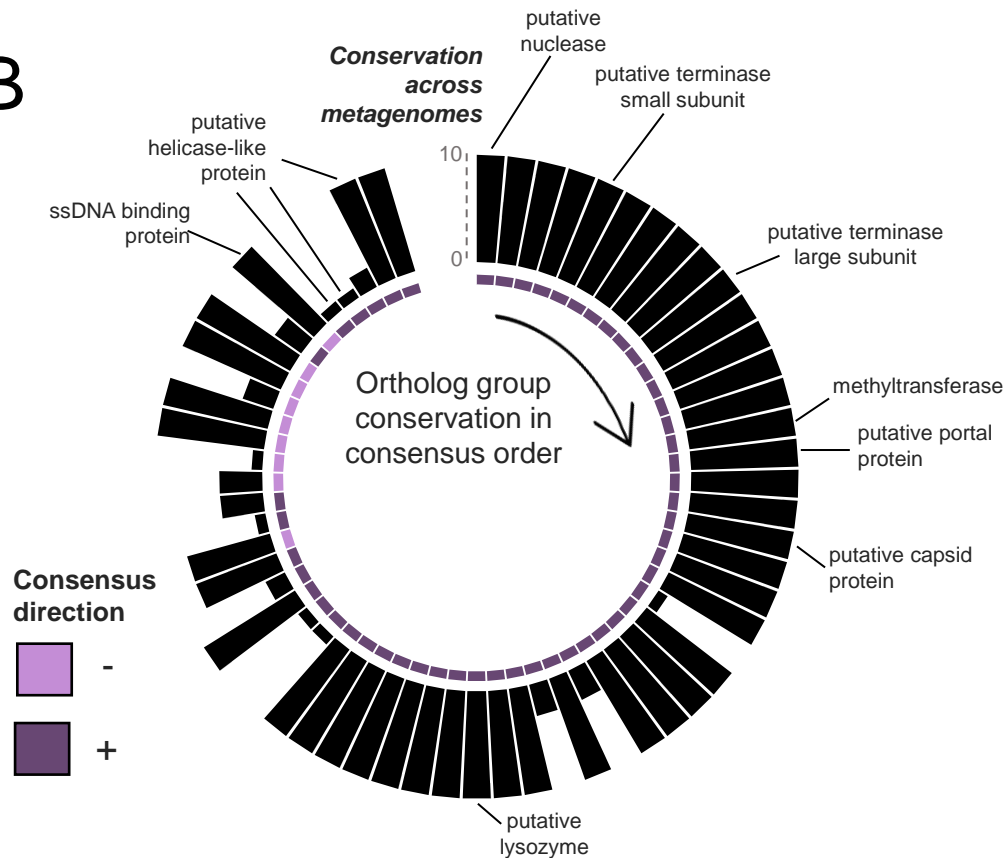


Figure 2: Targeted viral detection in metagenomes using fai. **A)** Total metagenomes from a single site in Lake Mendota across multiple depths and timepoints from Tran et al. 2023 were investigated using fai for the presence of a virus found in two of the three earliest microbiome samplings (red box). The presence of the virus is indicated by a virus 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. **D)** 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 conservation of the ortholog groups across the ten metagenomes the virus was detected in. BioRender was used in generation of this figure.

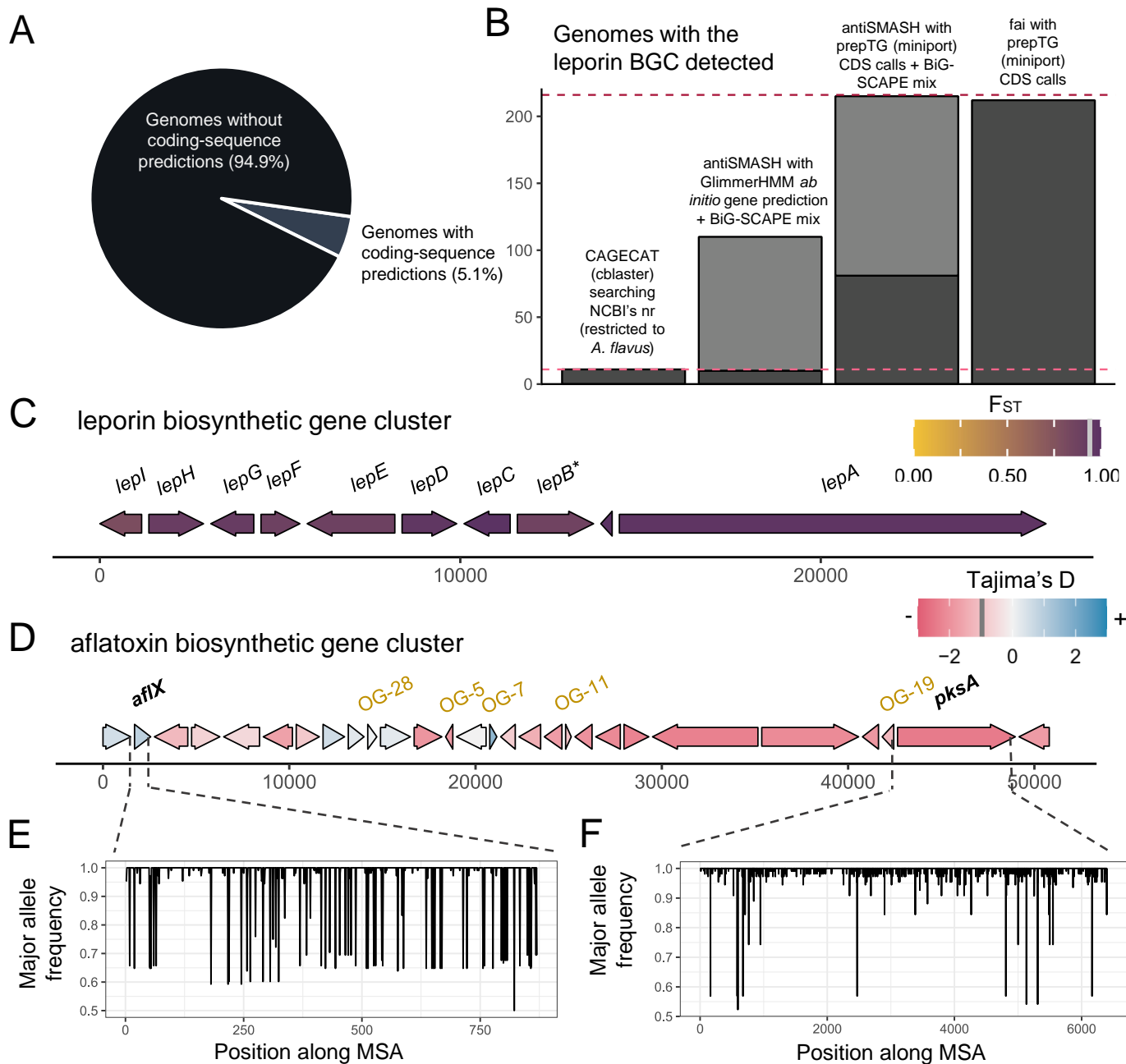
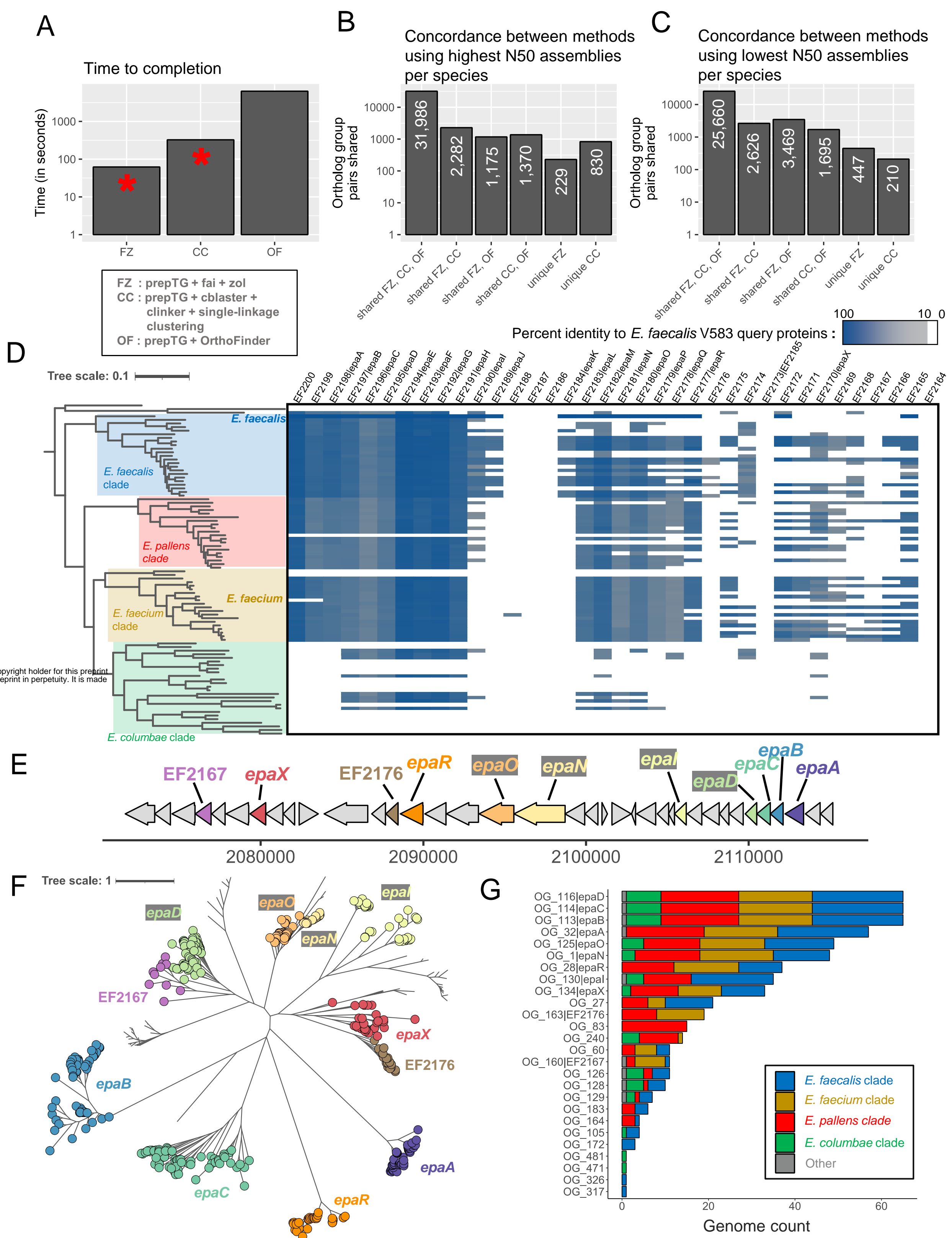


Figure 3: Evolutionary trends of common BGCs in *A. flavus*. **A)** The proportion of 216 *A. flavus* genomes from NCBI's GenBank database with coding-sequence predictions available. **B)** Comparison of the sensitivity of fai and alternate approaches based on assemblies for detecting the leporin BGC. The dashed violet line indicates the total number of genomes ($n=216$) assessed and the dashed pink line indicates the number of genomes with CDS features available on NCBI ($n=11$). Dark grey indicates instances identified by CAGECAT/cblaster or fai or as belonging to the same GCF as the reference leporin BGC from MIBiG by antiSMASH and BiG-SCAPE analysis. Lighter grey indicates the number of similar BGCs identified by BiG-SCAPE, belonging to the same clan but not to the same GCF as the reference leporin BGC. 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)** *aflX* and **(F)** *pksA*, which depict opposite trends in sequence conservation according to their respective Tajima's D calculations.



per 12, 2024. The copyright holder for this preprint
use to display the preprint in perpetuity. It is made
ense.

Figure 4: Searching for the *epa* locus across the diverse genus of *Enterococcus*. **A**) Overview of the time needed to run orthology/homology inference methods on the 92 genomes with the highest N50 for each distinct *Enterococcus* species. OrthoFinder was run at the genome-wide scale, while fai and cblaster were used to first identify genomic regions corresponding to the *epa* locus from *E. faecalis* V583 and subsequently zol and clinker were applied to determine ortholog groups, respectively. The red asterisks denote that manual assessment or filtering of homologous gene clusters identified by fai and cblaster is encouraged and thus additional time might be required for them. Counts showing the overlap in orthologous protein pair predictions by the three different methods are shown following their application to representative genomes from GTDB R214 with the **B**) highest N50 and **C**) lowest N50 for the 92 different species. **D**) The distribution of the *epa* locus, based on criteria used for running fai, is shown across a species phylogeny for 92 genomes representative of distinct *Enterococcus* species in GTDB R214. The coloring of the heatmap corresponds to the percent identity of the best matching protein from each genome to the query *epa* proteins from *E. faecalis* V583. **E**) A schematic of the *epa* gene cluster from *E. faecalis* V583 (from EF2164 to EF2200) with glycosyltransferase encoding genes shown in color. **F**) A maximum-likelihood phylogeny of zol-identified ortholog groups corresponding to glycosyltransferases in *epa* loci across *Enterococcus*. **G**) Distribution of different glycosyltransferase ortholog groups across the four major clades of *Enterococcus* are shown. For **D** and **F** the tree scales correspond to the number of amino acid substitutions along the alignments used for phylogeny construction.

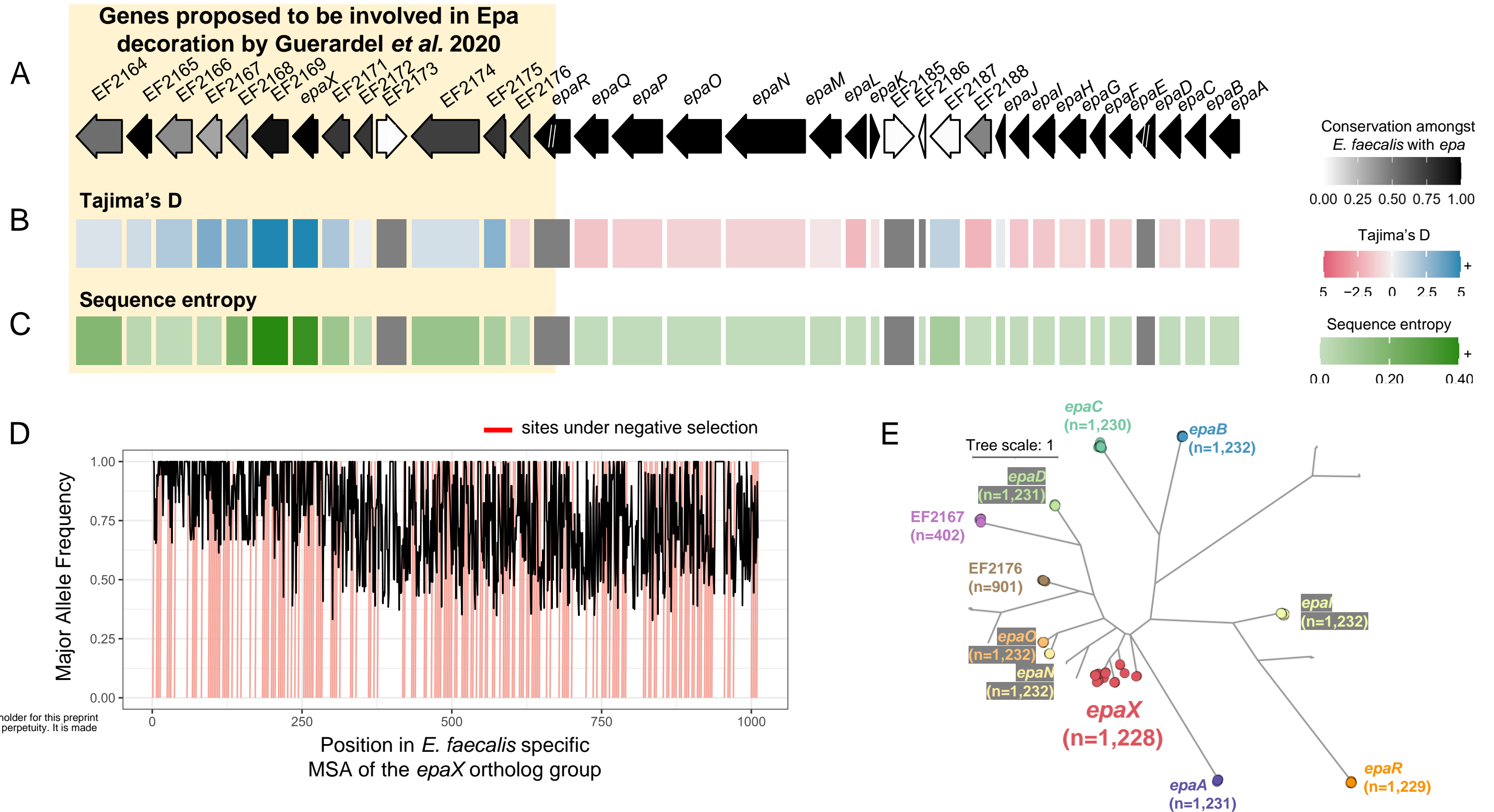


Figure 5: High sequence diversity of *epaX*-like glycosyltransferases amongst *E. faecalis*. A schematic of the *epa* locus from *E. faecalis* V583 with evolutionary statistics, **A**) conservation, **B**) Tajima's D and **C**) sequence entropy, gathered from the best corresponding ortholog group for each protein. Ortholog groups were inferred from zol investigation of 1,232 *epa* loci from the species. Genes upstream of and including *epaR* were recently proposed to be involved in Epa decoration by Guerardel *et al.* 2020. “//” indicates that the ortholog group was not single-copy in the context of the gene-cluster and calculation of evolutionary statistics for these genes was avoided (grey in panels B and C). Note, the same ortholog group was regarded for EF2173 and EF2185 which correspond to an identical *ISEf1* transposase. The length of proteins in the locus schematic are the median lengths of the corresponding ortholog groups. **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 of glycosyltransferase ortholog groups identified by zol which were found in >1% of *epa* instances. Ortholog groups identified by zol are indicated by colored circular nodes with names of *epa* genes from *E. faecalis* V583 noted where possible. The number of leaves/proteins for each clade is provided for labeled ortholog groups. The tree scale corresponds to the number of amino acid substitutions along the input protein alignment used for phylogeny construction.