

1 **The evolution of the gliotoxin biosynthetic gene cluster in *Penicillium* fungi**

2

3 Charu Balamurugan^{1,2}, Jacob L. Steenwyk^{1,2,3,*}, Gustavo H. Goldman⁴, & Antonis Rokas^{1,2,*}

4

5 ¹ Vanderbilt University, Department of Biological Sciences, VU Station B #35-1634, Nashville,
6 TN 37235, United States of America

7 ² Vanderbilt Evolutionary Studies Initiative, Vanderbilt University, Nashville, TN, United States

8 ³ Howards Hughes Medical Institute and the Department of Molecular and Cell Biology,
9 University of California, Berkeley, Berkeley, CA, USA

10 ⁴ Faculdade de Ciências Farmacêuticas de Ribeirão Preto, Universidade de São Paulo, Ribeirão
11 Preto, São Paulo, Brazil

12

13 *Correspondence should be addressed to: jlsteenwyk@berkeley.edu and

14 antonis.rokas@vanderbilt.edu

15

16 Running title: Gliotoxin Evolution in *Penicillium* Fungi

17

18 Keywords: comparative genomics; evolutionary biology; secondary metabolic gene clusters;
19 duplication and loss; plant pathogen; secondary metabolism; specialized metabolism

20

21 **Abstract**

22 Fungi biosynthesize a diversity of secondary metabolites, small organic bioactive molecules that
23 play diverse roles in fungal ecology. Fungal secondary metabolites are often encoded by
24 physically clustered sets of genes known as biosynthetic gene clusters (BGCs). Fungi in the
25 genus *Penicillium* produce diverse secondary metabolites that have been both useful (e.g., the
26 antibiotic penicillin and the cholesterol-lowering drug mevastatin) and harmful (e.g., the
27 mycotoxin patulin and the immunosuppressant gliotoxin) to human affairs. BGCs often also
28 encode resistance genes that confer self-protection to the secondary metabolite-producing
29 fungus. Some *Penicillium* species, such as *Penicillium lilacinoechinulatum* and *Penicillium*
30 *decumbens*, are known to produce gliotoxin, a secondary metabolite with known
31 immunosuppressant activity; however, an evolutionary characterization of the BGC responsible
32 for gliotoxin biosynthesis among *Penicillium* species is lacking. Here, we examine the
33 conservation of genes involved in gliotoxin biosynthesis and resistance in 35 *Penicillium*
34 genomes from 23 species. We found homologous, less fragmented gliotoxin BGCs in 12
35 genomes, mostly fragmented remnants of the gliotoxin BGC in 21 genomes, whereas the
36 remaining two *Penicillium* genomes lacked the gliotoxin BGC altogether. In contrast, we
37 observed broad conservation of homologs of resistance genes that reside outside the BGC across
38 *Penicillium* genomes. Evolutionary rate analysis revealed that BGCs with higher numbers of
39 genes evolve slower than BGCs with few genes. Even though the gliotoxin BGC is fragmented
40 to varying degrees in nearly all genomes examined, ancestral state reconstruction suggests that
41 the ancestor of *Penicillium* species possessed the gliotoxin BGC. Our analyses suggest that genes
42 that are part of BGCs can be retained in genomes long after the loss of secondary metabolite
43 biosynthesis.

44

45 **Introduction**

46 Gliotoxin is a secondary metabolite produced by certain fungi, including the major opportunistic
47 human pathogen *Aspergillus fumigatus* (Raffa and Keller 2019). Secondary metabolites are
48 bioactive molecules of low molecular weight that are not required for the organism's growth but
49 aid survival in harsh environments (Raffa and Keller 2019). Genes that participate in the
50 biosynthesis of secondary metabolites, including gliotoxin, typically reside next to each other in
51 fungal genomes and form biosynthetic gene clusters (BGCs) (Rokas *et al.* 2020). The gliotoxin
52 BGC is implicated in human pathogenicity because gliotoxin suppresses the immune response of
53 the mammalian host through diverse mechanisms, including by inhibiting protein complexes
54 necessary for the generation of antimicrobial reactive oxygen species, decreasing cytotoxic
55 activities of T lymphocytes, and preventing integrin activation (Dolan *et al.* 2015; Raffa and
56 Keller 2019). Gliotoxin's role in modulating host biology suggests that it is a virulence factor
57 (Raffa and Keller 2019). For example, virulence is attenuated in certain animal models of disease
58 when *gliP*, the non-ribosomal peptide synthetase gene involved in gliotoxin biosynthesis, is
59 deleted (Sugui *et al.* 2007).

60

61 Fungi that produce gliotoxin need to be resistant to the toxin. Several genes contribute to
62 resistance, such as the thioredoxin reductase gene *gliT*, located within the gliotoxin BGC
63 (Schrettl *et al.* 2010). *gliT* deletion strains of *A. fumigatus* exhibit resistance to gliotoxin
64 oxidation and unchecked methylation (Owens *et al.* 2015). As a result, *gliT*-deficient *A.*
65 *fumigatus* are hypersensitive to gliotoxin (Owens *et al.* 2015). Other resistance genes encoding
66 transcription factors, transporters, and oxidoreductases, reside outside the BGC and – like *gliT* –

67 are found in both gliotoxin-producing and non-producing species (Castro *et al.* 2022). For
68 example, the transcription factor RglT is the primary regulator of *gliT* (Ries *et al.* 2020). Seven
69 other genes are known to be regulated by *rglT* and contribute to gliotoxin resistance: *gtmA*
70 (encodes a *bis*-thiomethyltransferase, AFUA_2G11120), *kojR* (transcription factor,
71 AFUA_5G06800), *abcCI* (ABC-transporter, AN7879/AFUA_1G10390), *mtrA*
72 (methyltransferase, AN3717/AFUA_6G12780), AN9051 (oxidoreductase, AFUA_7G00700),
73 AN1472 (MFS transporter, AFUA_8G04630), and AN9531 (NmrA-transcription factor,
74 AFUA_7G06920) (Castro *et al.* 2022).

75
76 Though progress has been made in understanding the mechanisms and functions of the gliotoxin
77 biosynthetic pathway, several questions remain, especially concerning the evolutionary and
78 ecological significance of this BGC in lineages that contain a mix of biotechnologically and
79 medically relevant fungi, such as *Penicillium* (Steenwyk *et al.* 2019). For example, *Penicillium*
80 *camemberti* and *Penicillium roqueforti* contribute to cheese production (Nelson 1970; Lessard *et*
81 *al.* 2012), whereas *Penicillium expansum*, *Penicillium digitatum*, and *Penicillium italicum* are
82 postharvest pathogens of citrus fruits, stored grains, and other cereal crops (Marcet-Houben *et al.*
83 2012; Ballester *et al.* 2015; Li *et al.* 2015). Examination of the gliotoxin BGC in the genomes of
84 *Penicillium* species will shed light on the evolution of the gliotoxin BGC within Aspergillaceae,
85 the family encompassing both *Aspergillus* and *Penicillium* species.

86
87 Considering the close relatedness of *Penicillium* and *Aspergillus*, it is interesting that evidence of
88 gliotoxin production is scant within the former. To fill this gap, we employed a genome-scale
89 approach to infer the evolutionary history of the gliotoxin BGC among 35 strains of 23

90 *Penicillium* species. We found that most *Penicillium* genomes examined contained fragmented
91 gliotoxin BGCs and two lacked a BGC. However, some *P. expansum* strains had two
92 homologous gliotoxin BGCs. Codon optimization analysis reveals that genes in *Penicillium*
93 BGCs are lowly optimized, whereas genes in *Aspergillus* gliotoxin BGCs are highly optimized.
94
95 In contrast, gliotoxin resistance genes in *Penicillium* and *Aspergillus* fungi have similar degrees
96 of codon optimization, suggesting that *Penicillium* species encounter exogenous gliotoxin in
97 their environments. Examination of evolutionary rates revealed that genes from highly
98 fragmented gliotoxin BGCs evolved at significantly higher rates than genes from less fragmented
99 BGCs, suggesting that less fragmented BGCs have been experiencing relaxation of selective
100 constraints for longer. Ancestral state reconstructions indicate that the *Penicillium* ancestor
101 possessed a less fragmented gliotoxin BGC, followed by distinct trajectories of duplication and
102 loss, highlighting the diverse evolutionary pathways of the gliotoxin BGC in *Penicillium* species.

103

104 **Materials and Methods**

105 **I. Data collection and quality assessment**

106 We retrieved the genomes and gene annotations of 35 *Penicillium* strains from 23 species as well
107 as of two outgroups (*Aspergillus fumigatus* and *Aspergillus fischeri*) from NCBI
108 (<https://www.ncbi.nlm.nih.gov/>) (Table S1).

109

110 Genome assembly and annotation quality were examined to evaluate whether the dataset is
111 sufficient for comparative genomics. The quality and characteristics of the genomes (N50, L50,
112 assembly size, number of scaffolds, and gene count) were evaluated using BioKIT (v0.1.0)

113 (Steenwyk *et al.* 2022) (Figure S1). The average N50 value was 1,850,972.1 bases, where 46%
114 of proteomes consisted of N50 values greater than 1 Mb, and the lowest N50 value was 31,119
115 bases for *P. expansum* CMP 1. Gene annotation completeness was assessed using BUSCO
116 (v5.0.0) (Waterhouse *et al.* 2018) (Figure S2). BUSCO uses a predetermined set of near-
117 universally conserved single-copy genes (or BUSCO genes) to identify their presence in a query
118 proteome (characterized as single-copy, duplicated, or fragmented) or absence. We used the
119 4,181 BUSCO genes from the Eurotiales OrthoDB dataset (Manni *et al.* 2021; Zdobnov *et al.*
120 2021). Nearly all the genomes have high BUSCO gene coverage (average: 95.9% \pm 3.1%), with
121 the lowest percentages being for *P. coprophilum* (87.9%) and *P. decumbens* (85.3%).
122

123 II. Identification and characterization of gliotoxin BGC and resistance genes

124 a. Identification of gliotoxin GBC and resistance genes

125 The representative gliotoxin BGC (BGC0000361, Download date: April 2022) from the
126 *Aspergillus fumigatus* Af293 reference strain was downloaded from the Minimum Information
127 about a Biosynthetic Gene Cluster (MiBIG) database (Kautsar *et al.* 2019). Command-line NCBI
128 BLASTP (Camacho *et al.* 2009) searches for the Af293 gliotoxin BGC against the proteome of
129 each species were executed. Highly similar sequences were identified using an expectation value
130 threshold of 1e-4 and a query coverage of 50%. The resulting BLAST outputs were then cross-
131 referenced with the NCBI feature table file, which contains genome location information for
132 each gene, and parsed to identify clusters of homologs. Less fragmented BGCs are defined as
133 having at least 7 / 13 genes from the query gliotoxin BGC present, including *gliP*, encoding the
134 core nonribosomal peptide synthetase (Castro *et al.* 2022); mostly fragmented clusters are
135 defined as having at least 3 / 13 genes from the gliotoxin BGC without a requirement for this

136 cluster to include *gliP*. When identifying BGCs, we allowed up to four genes between each pair
137 of adjacent homologs using the *A. fumigatus* Af293 BGC from the MiBIG database (Kautsar *et*
138 *al.* 2019) as reference (Castro *et al.* 2022).

139

140 To rule out gene annotation errors in cases where we infer genes to be absent, we conducted
141 command-line NCBI tBLASTn searches for the Af293 gliotoxin BGC against the genome
142 sequences. Highly similar sequences were identified using an expectation value threshold of 1e-
143 10. The resulting outputs were analyzed, and no new presence/absence information was found.

144

145 Sequence similarity searches were also conducted for eight gliotoxin resistance genes
146 (*abcC1*/AN7879, *mtrA*/AN3717, AN9051, AN1472, AN9531, *rglT*, *gtmA*/AFU2G11120,
147 *kojR*/AFUA_5G06800), three of which were transcription factors (AN9531, *rglT*, *kojR*). We
148 used an expectation value threshold of 1e-3 and a query coverage threshold of 50%; we used a
149 lower query coverage threshold of 40% for the three transcription factors.

150

151 **b. Codon bias**

152 To estimate the potential functional significance of the partial gliotoxin BGCs present in
153 *Penicillium* genomes, mean gene-wise relative synonymous codon usage (gRSCU) was
154 determined for each clustered *gli* gene across all proteomes using BioKIT (Steenwyk *et al.*
155 2022). This provides insight into how codon usage bias influences the expression level of a
156 particular gene. The percentile rankings of each of the present and clustered *gli* genes were
157 calculated using the R package *dplyr* (v1.0.9) (Wickham *et al.* 2022), and these values, for each
158 species, were then plotted using the R package *ggplot2* (Wickham 2016).

159

160 **c. Synteny Analysis**

161 Alignments of representative *Penicillium* genomes with less and more fragmented gliotoxin
162 BGCs were generated using a GenomeDiagram in Biopython (Cock *et al.* 2009). Five genomes
163 (*A. fumigatus* Af293, *P. flavigenum* IBT 14082, *P. roqueforti* FM164, *P. nordicum* DAOMC
164 185683, and *P. expansum* CMP1) with the largest number of different, homologous *gli* cluster
165 genes above seven, and including *gliP*, were chosen to visualize the conservation of synteny of
166 less fragmented gliotoxin BGCs across the phylogeny. Similarly, the five genomes (*P. steckii*
167 IBT 24891, *P. vulpinum* IBT 29486, *P. rubens* 43M1, *P. camemberti* FM 013, and *P. italicum*
168 PHI 1) with the greatest number of different, homologous *gli* cluster genes above three and
169 below seven, and not needing to include *gliP*, were chosen to visualize synteny of mostly
170 fragmented BGCs across the phylogeny.

171

172 **III. Phylogenetic Analysis**

173 **a. Species Tree Inference**

174 The evolutionary relationships of *Penicillium* species were obtained from a previous study
175 (Steenwyk *et al.* 2019) using treehouse (Steenwyk and Rokas 2019). For three species with
176 population-level data, within-species relationships were inferred using phylogenomics. To do so,
177 protein sequences of BUSCO genes were first aligned using MAFFT (v7.490) with the *--auto*
178 parameter (Kato and Standley 2013). Codon-based alignments were generated by threading the
179 corresponding DNA sequences onto the protein alignment with the *thread_dna* function in
180 PhyKIT (v1.11.2) (Steenwyk *et al.* 2021). The resulting nucleotide alignments were trimmed
181 using ClipKIT (v1.3.0) (Steenwyk *et al.* 2020) with default parameters. The resulting aligned and

182 trimmed sequences were concatenated into a supermatrix with 8,124,861 sites using the
183 *create_concat* function in PhyKIT. We then inputted the concatenated matrix into IQ-TREE 2
184 (v2.0.6), a software that implements a maximum likelihood framework for inferring phylogenies.
185 All other evolutionary relationships between species were constrained following the relationships
186 inferred in a previously published study (Steenwyk and Rokas 2019). The best-fitting
187 substitution model (GTR+F+I+G4) was determined using ModelFinder (Kalyaanamoorthy *et al.*
188 2017).

189

190 **b. Single-gene tree inference**

191 To infer the evolutionary history of genes in the gliotoxin BGCs, individual *gli* genes were
192 compiled and aligned with MAFFT (v7.490) using the *--auto* parameter (Kato and Standley
193 2013). The corresponding nucleotide sequences for each file were obtained from the CDS files
194 for each species, using the *faidx* function of BioKIT (v0.1.0) (Steenwyk *et al.* 2022). These
195 nucleotide sequences were then threaded onto the protein alignments using the *thread_dna*
196 function of PhyKIT (Steenwyk *et al.* 2021), resulting in a codon-based alignment. All individual
197 codon-based gene alignments were trimmed with ClipKIT (Steenwyk *et al.* 2020) with default
198 parameters. The trimmed alignments were used to construct a phylogeny using IQ-TREE 2
199 (Minh *et al.* 2020). The best-fitting substitution model was chosen for each *gli* gene using
200 Bayesian information criteria (BIC) implemented in ModelFinder (Kalyaanamoorthy *et al.* 2017)
201 from IQ-TREE 2. Branch support in each phylogenetic tree was assessed by 1000 bootstraps
202 using ultrafast bootstrapping approximation (Hoang *et al.* 2018). Tree visualization was carried
203 out using the R packages *ape* (v5.6.2) (Paradis and Schliep 2019) and *phytools* (v1.0.3) (Revell
204 2012).

205
206 To characterize variation in the evolution of individual genes of the gliotoxin BGC, the trimmed
207 alignments and maximum-likelihood trees from IQ-TREE 2 were used as input into the
208 *evolutionary_rate*, *total_tree_length*, and *pairwise_identity* functions of PhyKIT to estimate two
209 tree-based measures of evolutionary rate and one sequence-based measure. Evolutionary rate is
210 defined as the total tree length divided by the number of terminals (Telford *et al.* 2014; Steenwyk
211 *et al.* 2021). The total tree length is the sum of all branches (Steenwyk *et al.* 2021).

212

213 c. Ancestral state reconstructions

214 Ancestral state reconstruction for each gene of the gliotoxin BGC across three discrete characters
215 (“Presence clustered,” “Presence unclustered,” and “Absence”) was estimated using *phytools*
216 (v1.0.3) (Revell 2012). Presence generally indicates that a homolog of the particular gene was
217 identified. “Presence clustered” identifies an existing homolog of the specific gene within a
218 maximum distance of four genes from other homologs of the gliotoxin BGC. “Presence
219 unclustered” identifies an existing homolog of the particular gene without clustering. “Absence”
220 indicates that no homolog of the specific gene was identified. Estimation of ancestral character
221 states was done using the Dollo parsimony method. This method assumes that a complex
222 character lost during the evolution of a particular lineage cannot be regained (Rogozin *et al.*
223 2006). Count, a software package for the evolutionary analysis of homolog family sizes, was
224 used to generate these ancestral state reconstructions (Csűös 2010).

225

226 d. Tree Topology Testing

227 Tree topology testing was used to determine whether the duplication event resulting in the two
228 less fragmented, homologous gliotoxin BGCs in *P. expansum* strains MD 8 and d1 occurred
229 solely in the lineage of *P. expansum* or deeper in the tree. IQTREE 2 (Minh *et al.* 2020) was used
230 to compute log-likelihoods of a constrained tree (monophyly of *P. expansum gliP* homologs) and
231 the observed tree in which a polyphyly of *gliP* in both clusters is seen (inconsistent with the
232 known species tree). 1000 RELL replicates (Kishino *et al.* 1990) were performed. The AU test
233 results (Shimodaira 2002) was used for comparison.

234

235 **Results and Discussion**

236 I. **The gliotoxin BGC is fragmented in *Penicillium* species**

237 Presence / absence data of the 13 genes in the gliotoxin BGC among the 23 *Penicillium* species
238 analyzed reveals that the cluster is largely fragmented in the genus *Penicillium* (Figure 1). The
239 proteomes of 12 strains from 5 *Penicillium* species (*P. arizonense*, *P. flavigenum*, *P. roqueforti*,
240 *P. nordicum*, *P. expansum*), possessed less fragmented BGCs, and the proteomes of 23 strains
241 from 18 *Penicillium* species had mostly fragmented BGCs (Figure S3-S15). Two less fragmented
242 BGCs, which contained 10 / 13 genes and 7 / 13 genes, were identified in *P. expansum* strains d1
243 and MD 8, respectively. Regardless of the number of less fragmented BGCs found, to our
244 knowledge, none of the *Penicillium* species in question are known to produce gliotoxin, except
245 *P. decumbens* (Feng *et al.* 2018), suggesting that the absence of clustering in this species may be
246 due to strain heterogeneity and requires further exploration.

247

248 II. **A complete gliotoxin BGC was present in the ancestor of *Penicillium* species**

249 Ancestral state reconstruction revealed the presence of all 13 genes in the gliotoxin BGC in the
250 ancestor of the *Penicillium* species used in our study (Figure 1). We infer that the first gene lost
251 was *gliH*, which is absent from 25 of the 35 *Penicillium* strains examined. The *gliH* gene
252 encodes an acetyltransferase that, when deleted, results in a loss of gliotoxin production in *A.*
253 *fumigatus* (Schrettl *et al.* 2010; Castro *et al.* 2022). Thus, the early loss of *gliH* in the genus
254 *Penicillium* may have been the key determinant of a lack of gliotoxin production. Further, the
255 synteny of genes in the BGC is mostly conserved and similar to the arrangement of the *A.*
256 *fumigatus* Af293 gliotoxin BGC across representative, less fragmented BGCs, such as *P.*
257 *flavigenum* IBT 14082 and *P. expansum* CMP 1 (Figure 2). In contrast, there is extensive
258 divergence in synteny conservation among mostly fragmented BGCs (Figure 2). To our
259 knowledge, none of the *Penicillium* species examined are known to produce gliotoxin, except *P.*
260 *decumbens* (Feng *et al.* 2018), suggesting that the absence of clustering in this species may be
261 due to strain heterogeneity and requires further exploration.

262

263 III. Resistance genes are broadly conserved

264 The presence/absence results of the eight resistance genes, portrayed in Figure 1, suggest that
265 their origins predate the *Aspergillus* and *Penicillium* genera (Figure S16-S23). All species
266 possessed *abcC1*, AN1472, AN9051, AN9531, and *kojR* homologs. In addition, only *Penicillium*
267 species with mostly fragmented gliotoxin BGCs lacked at least one resistance gene, such as
268 *gtmA*, *mtrA*, and *rglT*. *Penicillium chrysogenum* lacked both *rglT* and *gliT*, an observation
269 consistent with the transcriptional dependency of *gliT* to *rglT* (Ries *et al.* 2020).

270

271 IV. *Penicillium* species have experienced changes in gliotoxin BGC synteny over time

272 All genes of the gliotoxin BGC were broadly found within the genus *Penicillium*, except for
273 *gliH*, yet most were sparsely clustered (Figure 1). More specifically, 12 out of 35 *Penicillium*
274 species/strains were found to have a less fragmented, homologous BGC. Two strains of
275 *Penicillium expansum* (d1 and MD 8) were found to have two BGCs. Evidence of variation in
276 gene presence / absence is also evident within species. For example, *Penicillium roqueforti*
277 shows population variation in the presence of *gliZ*, a major transcriptional regulator of gliotoxin
278 biosynthesis (Bok *et al.* 2006); five strains of *P. roqueforti* lack *gliZ* whereas one strain has the
279 gene. As a result, we can conclude that the ancestor of *P. roqueforti* had a *gliZ* homolog, but the
280 gene was lost over time in most of the strains, highlighting the importance of population-level
281 sampling. Overall, we can see that the gliotoxin BGC has experienced relocations and
282 duplications of its genes, specifically in *Penicillium expansum* strains d1 and MD 8, as is
283 expected in the formation of most secondary metabolite-producing BGCs (Rokas *et al.* 2018).

284

285 V. Few *Penicillium* species contain codon-optimized gliotoxin BGCs

286 Compared to the two outgroup *Aspergillus* species, *A. fumigatus* and *A. fischeri*, *Penicillium*
287 species have much lower gRSCU value rankings (Figure 3). Specifically, the mean gRSCU
288 percentile rank of gliotoxin BGC genes among the *Aspergillus* outgroups is 0.81, while that
289 among the *Penicillium* species is 0.35; these scores suggest that *gli* genes from *Aspergillus* are
290 more codon-optimized than *gli* genes from *Penicillium*. Regardless of mean gRSCU values, *gliT*
291 and *gliA* homologs, when present, are ranked consistently in the top three to four clustered genes.
292 However, when considering resistance genes, the spread and range of their gRSCU values are
293 similar across all species. The mean gRSCU percentile rank of gliotoxin resistance genes among
294 the *Aspergillus* outgroups is 0.58, while that among the *Penicillium* species is 0.53. This allows

295 us to infer that these *Penicillium* species may ecologically encounter exogenous gliotoxin,
296 making *gliT*, encoding a gliotoxin-neutralizing enzyme, *gliA*, encoding a transporter that exports
297 gliotoxin, and non-TF resistance genes such as *abcC1*, encoding an ABC-transporter, rank in the
298 top percentiles among each of the species' gene sets.

299

300 **VI. *gli* genes in less fragmented clusters are evolving at a slower rate than mostly**
301 **fragmented clusters**

302 In the comparison of tree-based and sequence-based measures of evolutionary rate, *gli* genes
303 from less fragmented clusters are evolving at a significantly slower pace ($p < 0.0001$) than those
304 from mostly fragmented clusters across all three metrics, as seen by a two-way ANOVA with an
305 additive model (Figure 4, Figure S3-S15). This difference highlights a notable feature of many
306 BGCs, the fact that they are rapidly evolving, hinted at by their high variability and narrow
307 taxonomic range (Rokas *et al.* 2020).

308

309 **VII. A duplication of the gliotoxin BGC may have occurred before the divergence**
310 **between *P. flavigenum* and *P. roqueforti***

311 We conducted a tree topology test to infer when the gliotoxin BGCs found in *P. flavigenum*
312 occurred. The maximum likelihood phylogeny suggests that this duplication occurred before the
313 divergence between *P. flavigenum* and *P. roqueforti*. An alternative hypothesis is that
314 duplication occurred within *P. expansum*. This alternative hypothesis would be supported by
315 monophyly of *P. expansum* homologs of BGC genes. After conducting a tree topology test
316 comparing log likelihood values between the maximum likelihood phylogeny and an alternative
317 tree wherein *P. expansum* *gliP* homologs were constrained to be monophyletic, we found that the

318 constrained topology was significantly rejected (Approximately Unbiased test, $p = 7.34e-110$)
319 (Figure S24). In other words, it is unlikely duplication occurred within *P. expansum* lineage;
320 instead, duplication likely occurred more anciently, prior to the diversification of *P. expansum*.

321

322 **Conclusions**

323 The ancestor of *Penicillium* species likely possessed a complete gliotoxin BGC. A
324 duplication event of the BGC occurred in one lineage, likely prior to the divergence of *P.*
325 *flavigenum* and *P. roqueforti*. Also, the presence/absence results of the eight resistance genes
326 suggest that their origins predate the *Aspergillus* and *Penicillium* genera suggesting that
327 resistance has long been important among these species. The genes in *Penicillium* gliotoxin
328 BGCs are less codon optimized (gRSCU percentile rank mean: 0.35) compared to their
329 *Aspergillus* counterparts (gRSCU percentile rank mean: 0.81) suggesting that *gli* genes are much
330 more often expressed in *Aspergillus* species than in *Penicillium*. However, less fragmented
331 BGCs within *Penicillium* species are evolving at a slower rate than mostly fragmented clusters,
332 suggestive of potential functionality.

333 Although informative, this work only utilizes publicly available protein annotations of
334 biotechnologically and medically relevant *Penicillium* fungi, making it important to expand upon
335 the species/strains studied. Moreover, this same targeted gliotoxin analysis within a larger
336 phylogeny of *Aspergillus* species, for which there is greater evidence of the production of this
337 secondary metabolite, may be helpful. An analysis of gliotoxin BGCs encoded in all fungi would
338 also provide us with more insight into the evolutionary mechanisms that give rise to BGC
339 diversity. In addition, expanding on the causes of conservation of less fragmented gliotoxin
340 BGCs within a variety of *Penicillium* strains may be important, especially because evidence of

341 production is lacking. As a result, this exciting reality encourages further understanding of the
342 motivating hypothesis that individual secondary metabolites are “cards” of virulence in a larger
343 “hand” that fungi possess.

344

345 **Data availability**

346 The authors affirm that all data necessary for confirming the conclusions of the article are
347 present within the article, figures, tables, and supplemental material.

348

349 **Acknowledgements**

350 We thank members of the Rokas Laboratory at Vanderbilt University for support and
351 feedback on this work. We also thank the Vanderbilt Data Science Institute for their
352 undergraduate enrichment opportunities. This work was performed in part using resources
353 contained within the Advanced Computing Center for research and Education at Vanderbilt
354 University in Nashville, TN.

355

356 **Funding**

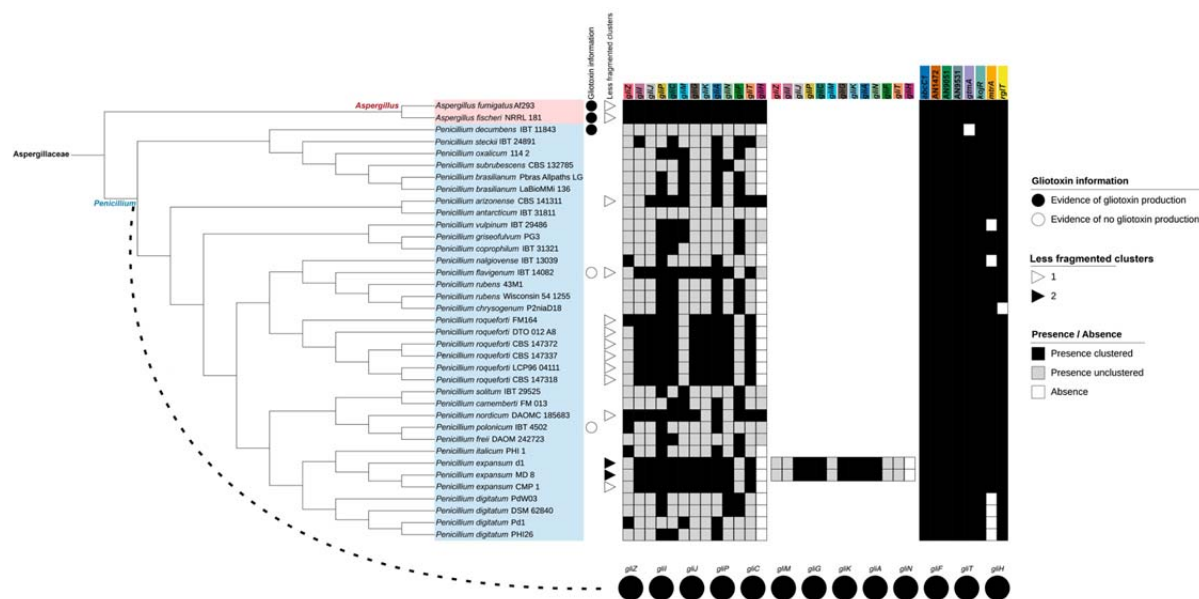
357 C.B. was supported by the Vanderbilt University Data Science Institute—Summer
358 Research Program. J.L.S. and A.R. were funded by the Howard Hughes Medical Institute
359 through the James H. Gilliam Fellowships for Advanced Study program. Research in A.R.’s lab
360 is supported by grants from the National Science Foundation (DEB-2110404), the National
361 Institutes of Health/National Institute of Allergy and Infectious Diseases (R01 AI153356), and
362 the Burroughs Wellcome Fund.

363

364 **Conflicts of interest**

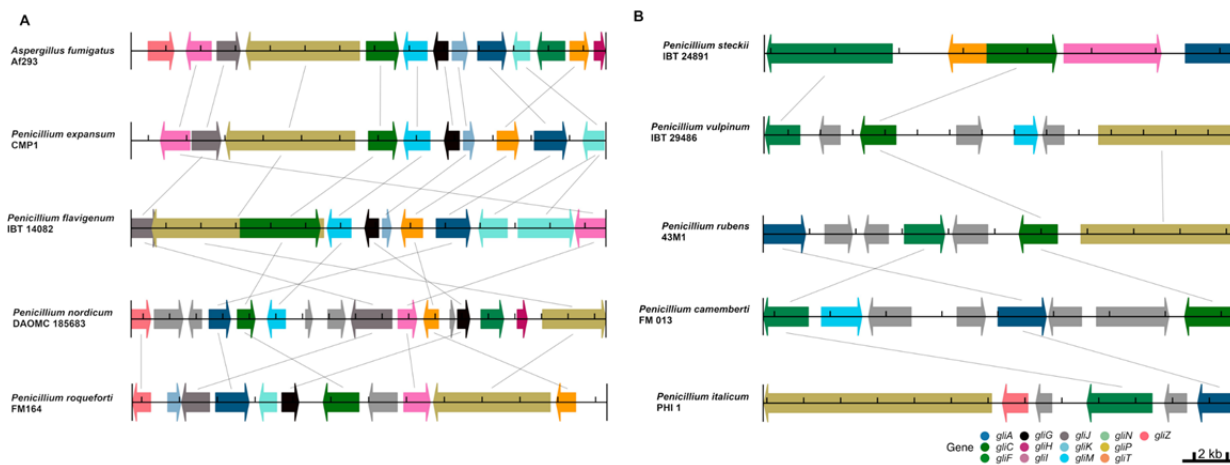
365 J.L.S. is a scientific consultant for Latch AI Inc. J.L.S. is a scientific advisor for WittGen

366 Biotechnologies. A.R. is a scientific consultant for LifeMine Therapeutics, Inc.



367 **Fig. 1 Phylogeny of *Penicillium* genomes.** Different genera are depicted using different-colored
 368 boxes. *Aspergillus* is shown in red and *Penicillium* in blue. Shaded circles next to species / strain
 369 names indicate gliotoxin production information from the literature, or lack thereof (Fischer *et al.*
 370 2000; Spikes *et al.* 2008; Knowles *et al.* 2020; Redrado *et al.* 2022). Shaded squares in the
 371 second column depict number of clusters identified. Remaining color strips depict gene presence
 372 clustered (black), presence unclustered (gray), and absence (white) according to the requirements
 373 outlined in the *Methods* section. Ancestral state reconstructions of each gene of the gliotoxin
 374 BGC (for the ancestor of *Penicillium* species) are presented in pie charts below the phylogeny.

375



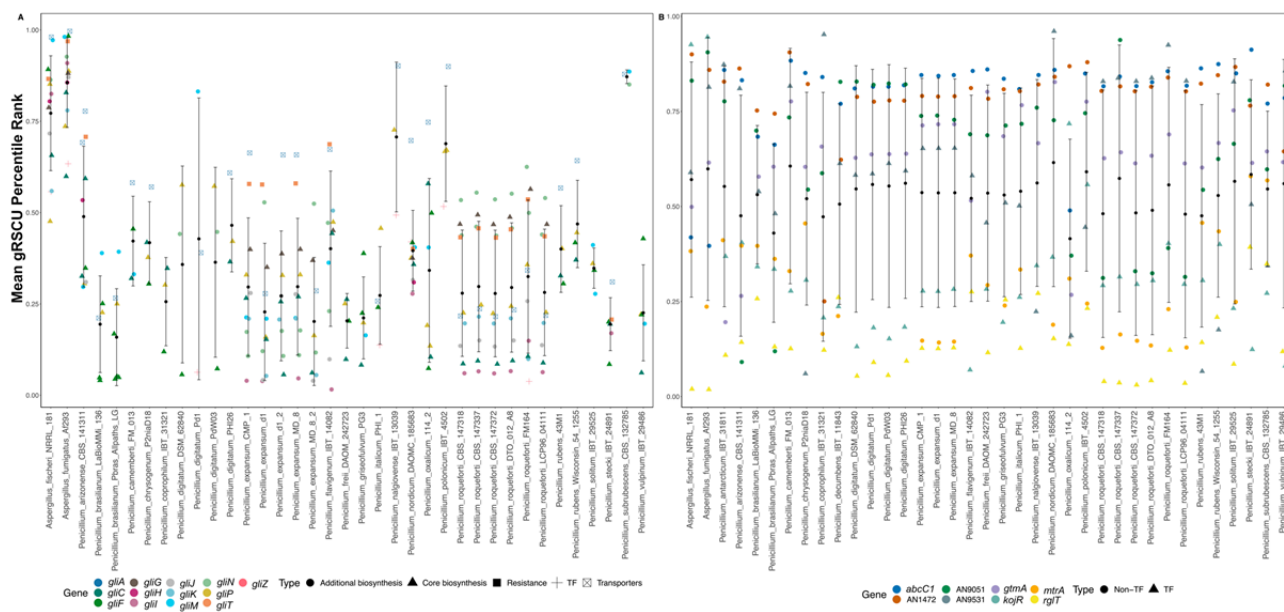
376

377 **Fig. 2 Conservation of gliotoxin BGC synteny for representative *Penicillium* species.**

378 Synteny analysis of representative genomes with less fragmented (A) and mostly fragmented (B)

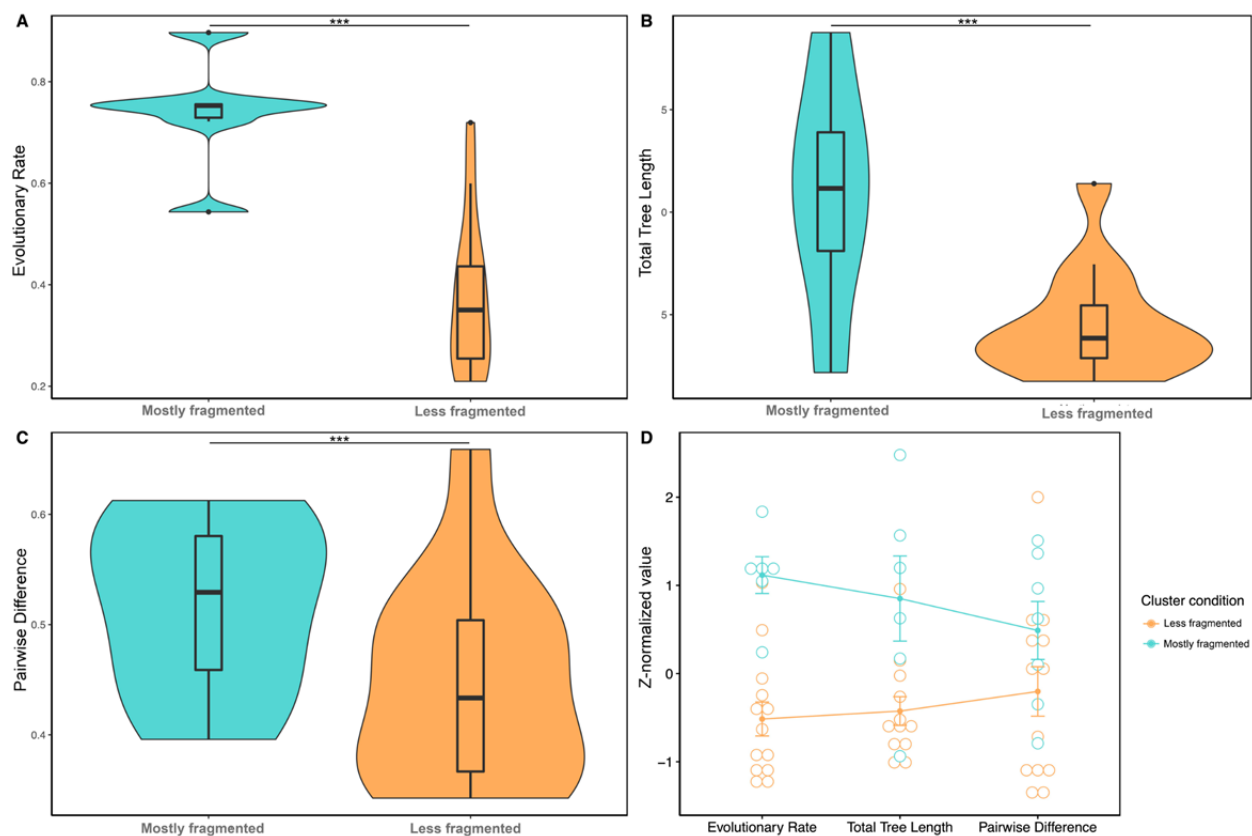
379 gliotoxin BGCs. Each interval along the track represents 2 kb.

380



381
 382
 383 **Fig. 3 Gene-wise relative synonymous codon usage (gRSCU) for gliotoxin BGC and**
 384 **resistance genes. (A)** Percentile rankings of gene-wise relative synonymous codon usage
 385 (gRSCU) among gliotoxin BGC genes, in comparison to all other genes. Types / functionality of
 386 each gene of the gliotoxin BGC is depicted by shape in the categories of “Core biosynthesis”,
 387 “Additional biosynthesis”, “Resistance”, “Transcription Factor”, “Transporter” **(B)** Percentile
 388 ranking of gene-wise relative synonymous codon usage (gRSCU) among gliotoxin resistance
 389 genes, in comparison to all other genes. Types / functionality of each resistance gene is depicted
 390 by shape in the categories of “Non-Transcription Factor and Transcription Factor”.

391



392

393 **Fig. 4 Evolutionary rate comparison across gliotoxin BGCs.** Multi-method comparison of
394 evolutionary rates between less fragmented and mostly fragmented gliotoxin BGCs. Less
395 fragmented clusters were required to contain a *gliP* ortholog and at least 7 different genes of the
396 cluster. Mostly fragmented clusters had no requirement to contain a *gliP* ortholog and only
397 needed to contain at least 3 different genes of the cluster. (A) Comparison of evolutionary rates,
398 as a function of total tree length divided by the number of taxa, between less fragmented and
399 mostly fragmented gliotoxin BGCs. (B) Comparison of total tree length between less fragmented
400 and mostly fragmented gliotoxin BGCs. (C) Comparison of pairwise identity between less
401 fragmented and mostly fragmented gliotoxin BGCs.

402

403 **Works Cited**

- 404 Ballester, A.-R., M. Marcet-Houben, E. Levin, N. Sela, C. Selma-Lázaro *et al.*, 2015 Genome,
405 Transcriptome, and Functional Analyses of *Penicillium expansum* Provide New Insights
406 Into Secondary Metabolism and Pathogenicity. *MPMI* 28: 232–248.
- 407 Bok, J. W., D. Chung, S. A. Balajee, K. A. Marr, D. Andes *et al.*, 2006 GliZ, a transcriptional
408 regulator of gliotoxin biosynthesis, contributes to *Aspergillus fumigatus* virulence. *Infect*
409 *Immun* 74: 6761–6768.
- 410 Camacho, C., G. Coulouris, V. Avagyan, N. Ma, J. Papadopoulos *et al.*, 2009 BLAST+:
411 architecture and applications. *BMC Bioinformatics* 10: 421.
- 412 Castro, P. A. de, A. C. Colabardini, M. Moraes, M. A. C. Horta, S. L. Knowles *et al.*, 2022
413 Regulation of gliotoxin biosynthesis and protection in *Aspergillus* species. *PLOS*
414 *Genetics* 18: e1009965.
- 415 Cock, P. J. A., T. Antao, J. T. Chang, B. A. Chapman, C. J. Cox *et al.*, 2009 Biopython: freely
416 available Python tools for computational molecular biology and bioinformatics.
417 *Bioinformatics* 25: 1422–1423.
- 418 Csűös, M., 2010 Count: evolutionary analysis of phylogenetic profiles with parsimony and
419 likelihood. *Bioinformatics* 26: 1910–1912.
- 420 Dolan, S. K., G. O’Keeffe, G. W. Jones, and S. Doyle, 2015 Resistance is not futile: gliotoxin
421 biosynthesis, functionality and utility. *Trends in Microbiology* 23: 419–428.
- 422 Feng, H., S. Liu, M. Su, E. L. Kim, J. Hong *et al.*, 2018 Gliotoxin is Antibacterial to Drug-
423 resistant Piscine Pathogens. *Nat Prod Sci* 24: 225–228.
- 424 Fischer, G., T. Müller, R. Schwalbe, R. Ostrowski, and W. Dott, 2000 Species-specific profiles
425 of mycotoxins produced in cultures and associated with conidia of airborne fungi derived

- 426 from biowaste. *International Journal of Hygiene and Environmental Health* 203: 105–
427 116.
- 428 Hoang, D. T., O. Chernomor, A. von Haeseler, B. Q. Minh, and L. S. Vinh, 2018 UFBoot2:
429 Improving the Ultrafast Bootstrap Approximation. *Mol Biol Evol* 35: 518–522.
- 430 Kalyaanamoorthy, S., B. Q. Minh, T. K. F. Wong, A. von Haeseler, and L. S. Jermin, 2017
431 ModelFinder: fast model selection for accurate phylogenetic estimates. *Nat Methods* 14:
432 587–589.
- 433 Katoh, K., and D. M. Standley, 2013 MAFFT Multiple Sequence Alignment Software Version 7:
434 Improvements in Performance and Usability. *Mol Biol Evol* 30: 772–780.
- 435 Kautsar, S. A., K. Blin, S. Shaw, J. C. Navarro-Muñoz, B. R. Terlouw *et al.*, 2019 MIBiG 2.0: a
436 repository for biosynthetic gene clusters of known function. *Nucleic Acids Research*
437 gkz882.
- 438 Kishino, H., T. Miyata, and M. Hasegawa, 1990 Maximum likelihood inference of protein
439 phylogeny and the origin of chloroplasts. *J Mol Evol* 31: 151–160.
- 440 Knowles, S. L., M. E. Mead, L. P. Silva, H. A. Raja, J. L. Steenwyk *et al.*, 2020 Gliotoxin, a
441 Known Virulence Factor in the Major Human Pathogen *Aspergillus fumigatus*, Is Also
442 Biosynthesized by Its Nonpathogenic Relative *Aspergillus fischeri*. *mBio* 11: e03361-19.
- 443 Lessard, M.-H., G. Bélanger, D. St-Gelais, and S. Labrie, 2012 The Composition of Camembert
444 Cheese-Ripening Cultures Modulates both Mycelial Growth and Appearance. *Applied*
445 *and Environmental Microbiology* 78: 1813–1819.
- 446 Li, B., Y. Zong, Z. Du, Y. Chen, Z. Zhang *et al.*, 2015 Genomic Characterization Reveals
447 Insights Into Patulin Biosynthesis and Pathogenicity in *Penicillium* Species. *MPMI* 28:
448 635–647.

- 449 Manni, M., M. R. Berkeley, M. Seppey, and E. M. Zdobnov, 2021 BUSCO: Assessing Genomic
450 Data Quality and Beyond. *Current Protocols* 1: e323.
- 451 Marcet-Houben, M., A.-R. Ballester, B. de la Fuente, E. Harries, J. F. Marcos *et al.*, 2012
452 Genome sequence of the necrotrophic fungus *Penicillium digitatum*, the main postharvest
453 pathogen of citrus. *BMC Genomics* 13: 646.
- 454 Minh, B. Q., H. A. Schmidt, O. Chernomor, D. Schrempf, M. D. Woodhams *et al.*, 2020 IQ-
455 TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic
456 Era. *Mol Biol Evol* 37: 1530–1534.
- 457 Nelson, J. Howard., 1970 Production of Blue cheese flavor via submerged fermentation by
458 *Penicillium roqueforti*. *J. Agric. Food Chem.* 18: 567–569.
- 459 Owens, R. A., G. O’Keeffe, E. B. Smith, S. K. Dolan, S. Hammel *et al.*, 2015 Interplay between
460 Gliotoxin Resistance, Secretion, and the Methyl/Methionine Cycle in *Aspergillus*
461 *fumigatus*. *Eukaryotic Cell* 14: 941–957.
- 462 Paradis, E., and K. Schliep, 2019 ape 5.0: an environment for modern phylogenetics and
463 evolutionary analyses in R. *Bioinformatics* 35: 526–528.
- 464 Raffa, N., and N. P. Keller, 2019 A call to arms: Mustering secondary metabolites for success
465 and survival of an opportunistic pathogen. *PLOS Pathogens* 15: e1007606.
- 466 Redrado, S., P. Esteban, M. P. Domingo, C. Lopez, A. Rezusta *et al.*, 2022 Integration of In
467 Silico and In Vitro Analysis of Gliotoxin Production Reveals a Narrow Range of
468 Producing Fungal Species. *Journal of Fungi* 8: 361.
- 469 Revell, L. J., 2012 phytools: an R package for phylogenetic comparative biology (and other
470 things). *Methods in Ecology and Evolution* 3: 217–223.

- 471 Ries, L. N. A., L. Pardeshi, Z. Dong, K. Tan, J. L. Steenwyk *et al.*, 2020 The *Aspergillus*
472 *fumigatus* transcription factor RglT is important for gliotoxin biosynthesis and self-
473 protection, and virulence. *PLoS Pathog* 16: e1008645.
- 474 Rogozin, I. B., Y. I. Wolf, V. N. Babenko, and E. V. Koonin, 2006 Dollo parsimony and the
475 reconstruction of genome evolution.
- 476 Rokas, A., M. E. Mead, J. L. Steenwyk, H. A. Raja, and N. H. Oberlies, 2020 Biosynthetic gene
477 clusters and the evolution of fungal chemodiversity. *Nat Prod Rep* 37: 868–878.
- 478 Rokas, A., J. H. Wisecaver, and A. L. Lind, 2018 The birth, evolution and death of metabolic
479 gene clusters in fungi. *Nat Rev Microbiol* 16: 731–744.
- 480 Schrettl, M., S. Carberry, K. Kavanagh, H. Haas, G. W. Jones *et al.*, 2010 Self-Protection against
481 Gliotoxin—A Component of the Gliotoxin Biosynthetic Cluster, GliT, Completely
482 Protects *Aspergillus fumigatus* Against Exogenous Gliotoxin. *PLOS Pathogens* 6:
483 e1000952.
- 484 Shimodaira, H., 2002 An Approximately Unbiased Test of Phylogenetic Tree Selection.
485 *Systematic Biology* 51: 492–508.
- 486 Spikes, S., R. Xu, C. K. Nguyen, G. Chamilos, D. P. Kontoyiannis *et al.*, 2008 Gliotoxin
487 Production in *Aspergillus fumigatus* Contributes to Host-Specific Differences in
488 Virulence. *The Journal of Infectious Diseases* 197: 479–486.
- 489 Steenwyk, J. L., T. J. Buida III, C. Gonçalves, D. C. Goltz, G. Morales *et al.*, 2022 BioKIT: a
490 versatile toolkit for processing and analyzing diverse types of sequence data. *Genetics*
491 221: iyac079.

- 492 Steenwyk, J. L., T. J. Buida III, A. L. Labella, Y. Li, X.-X. Shen *et al.*, 2021 PhyKIT: a broadly
493 applicable UNIX shell toolkit for processing and analyzing phylogenomic data.
494 *Bioinformatics* 37: 2325–2331.
- 495 Steenwyk, J. L., T. J. B. Iii, Y. Li, X.-X. Shen, and A. Rokas, 2020 ClipKIT: A multiple
496 sequence alignment trimming software for accurate phylogenomic inference. *PLOS*
497 *Biology* 18: e3001007.
- 498 Steenwyk, J. L., and A. Rokas, 2019 Treehouse: a user-friendly application to obtain subtrees
499 from large phylogenies. *BMC Research Notes* 12: 541.
- 500 Steenwyk, J. L., X.-X. Shen, A. L. Lind, G. H. Goldman, and A. Rokas, 2019 A Robust
501 Phylogenomic Time Tree for Biotechnologically and Medically Important Fungi in the
502 Genera *Aspergillus* and *Penicillium*. *mBio*.
- 503 Sugui, J. A., J. Pardo, Y. C. Chang, K. A. Zarembek, G. Nardone *et al.*, 2007 Gliotoxin is a
504 virulence factor of *Aspergillus fumigatus*: gliP deletion attenuates virulence in mice
505 immunosuppressed with hydrocortisone. *Eukaryot Cell* 6: 1562–1569.
- 506 Telford, M. J., C. J. Lowe, C. B. Cameron, O. Ortega-Martinez, J. Aronowicz *et al.*, 2014
507 Phylogenomic analysis of echinoderm class relationships supports Asterozoa.
508 *Proceedings of the Royal Society B: Biological Sciences* 281: 20140479.
- 509 Waterhouse, R. M., M. Seppey, F. A. Simão, M. Manni, P. Ioannidis *et al.*, 2018 BUSCO
510 Applications from Quality Assessments to Gene Prediction and Phylogenomics. *Mol Biol*
511 *Evol* 35: 543–548.
- 512 Wickham, H., 2016 *ggplot2: Elegant Graphics for Data Analysis*. Springer-Verlag New York.
- 513 Wickham, H., R. François, L. Henry, K. Müller, and RStudio, 2022 dplyr: A Grammar of Data
514 Manipulation.

515 Zdobnov, E. M., D. Kuznetsov, F. Tegenfeldt, M. Manni, M. Berkeley *et al.*, 2021 OrthoDB in
516 2020: evolutionary and functional annotations of orthologs. *Nucleic Acids Research* 49:
517 D389–D393.
518