

Beyond observation: genomic traits and machine learning algorithms for predicting fungal lifestyles

Yanpeng Chen

School of Life Science and Technology, Center for Informational Biology, University of Electronic Science and Technology of China

Pengwei Su

School of Life Science and Technology, Center for Informational Biology, University of Electronic Science and Technology of China

Marc Stadler

Rong Xiang

Kevin D. Hyde

Wenhui Tian

Sajeewa S. N. Maharachchikumbura (✉ sajeewa83@yahoo.com)

Research Article

Keywords: FCWDEs, Genomics, machine learning, PCWDEs, secretome, TEs

Posted Date: June 29th, 2023

DOI: <https://doi.org/10.21203/rs.3.rs-3118609/v1>

License:  This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

Additional Declarations:

Supplementary Figures and Tables are not available with this version.

1 **Beyond observation: genomic traits and machine learning algorithms for predicting fungal lifestyles**

2
3 Yanpeng Chen¹, Pengwei Su¹, Marc Stadler^{2,3}, Rong Xiang⁴, Kevin D. Hyde^{5,6}, Wenhui Tian¹, Sajeewa S. N.
4 Maharachchikumbura^{1*}

5
6 ¹School of Life Science and Technology, Center for Informational Biology, University of Electronic Science and
7 Technology of China, Chengdu 610054, China

8 ²Helmholtz Centre for Infection Research GmbH, Department Microbial Drugs and German Centre for Infection
9 Research (DZIF), Partner Site Hannover-Braunschweig, 38124 Braunschweig, Germany

10 ³Institute of Microbiology, Technische Universität Braunschweig, Inhofenstraße 7, 38124 Braunschweig,
11 Germany

12 ⁴Precision Medicine Center, The Second Affiliated Hospital of Chongqing Medical University, Chongqing
13 404100, China

14 ⁵Center of Excellence in Fungal Research, Mae Fah Luang University, Chaing Rai 57100, Thailand

15 ⁶Innovative Institute for Plant Health, Zhongkai University of Agriculture and Engineering, Guangzhou 510225,
16 China

17
18 ***Corresponding author:**

19 Sajeewa S. N. Maharachchikumbura

20
21 **Correspondence:**

22 Sajeewa S. N. Maharachchikumbura: sajeewa83@yahoo.com or sajeewa@uestc.edu.cn

23
24 **Abstract**

25 Economically and agriculturally important fungal species have various lifestyles, and they may shift from
26 mutualistic or saprobic to pathogenic depending on the habitat, host tolerance, and resource availability.
27 Traditionally, the determination of fungal lifestyles has been based on observation at a particular host or habitat.
28 Therefore, potential fungal pathogens have been neglected until they cause devastating impacts on human health,
29 food security, and ecosystem stability. This study focused on the class Sordariomycetes to explore the genomic
30 traits that could be used to determine the lifestyles of fungi and the possibility of predicting fungal lifestyles using
31 machine learning algorithms. A total of 638 representative genomes covering five subclasses, 17 orders and 50
32 families were selected and annotated. Through an extensive literature survey, the lifestyles of 555 genomes were
33 determined, including plant pathogens, saprotrophs, entomopathogens, mycoparasites, endophytes, human
34 pathogens and nematophagous fungi. We evaluated the influence of sequencing technologies and concluded that
35 second sequencing technologies have no influence on genome completeness but tend to generate a reduced size
36 of transposable elements. We constructed three numerical matrices: a basic genomic feature matrix including 25
37 features; a functional protein matrix including 24 features; and a combined matrix. The most comprehensively
38 comparative analysis to date across multiple lifestyles was conducted based on these matrices. Results indicate
39 that basic genomic features reflect more on phylogeny rather than lifestyle, but the abundance of functional
40 proteins displays relatively high discrimination not only in differentiating taxonomic groups at the higher levels
41 but also in differentiating lifestyles. Genome size, GC content and gene number showed powerful discrimination
42 for differentiating higher ranks, especially at the subclass level. Plant pathogens have the largest secretome;
43 whereas entomopathogens have the smallest secretome; and the abundance of secretomes is a useful indicator to
44 clearly differentiate plant pathogens from entomopathogens, mycoparasites, saprotrophs and entomopathogens,

45 and as well as differentiate entophytes from entomopathogens. Effectors have long been considered as disease
46 determinants, and we did observe that plant pathogens have more effectors than saprotrophs and entomopathogens.
47 However, we also observed a similar abundance of effectors in endophytes, suggesting that effectors maybe not a
48 reliable indicator for pathogenic fungi. Single functional protein could not differentiate all lifestyles, but
49 combinations of multiple numerical features of functional proteins result in accurate differentiation for most
50 lifestyles. Furthermore, models of six machine learning algorithms were trained, optimized and evaluated, and the
51 best-performance model was used to predict the lifestyle of 83 unlabeled genomes. Although the accuracy of the
52 best machine learning model was limited by the inadequate genome number of several lifestyles and the inaccurate
53 lifestyle assignments for some genomes, the predictive model still obtained a high degree of accuracy in
54 differentiating plant pathogens. The predictive model can be further optimized with more sequenced genomes in
55 the future, and provide a more reliable prediction. This can be used as an early warning system to identify
56 potentially devastating fungi and take appropriate measures to prevent their spread.

57

58 **Keywords:** FCWDEs, Genomics, machine learning, PCWDEs, secretome, TEs

59

60 Introduction

61 The class Sordariomycetes, established by Eriksson and Winka (Eriksson OE 1997), is the second-largest class of
62 the phylum Ascomycota (Hyde et al. 2020). Based on the latest outline of Wijayawardene et al. (2022), it
63 comprises 7 subclasses, 46 orders and 172 families. The perithecial ascomata and inoperculate unitunicate asci
64 are the main diagnostic features for distinguishing Sordariomycetes from other classes (Maharachchikumbura et
65 al. 2015). Most Sordariomycete species have been introduced based solely on either the anamorph or teleomorph,
66 and only a small number of them were characterized based on both anamorph and teleomorph (Wingfield et al.
67 2012; Maharachchikumbura et al. 2015; Réblová et al. 2016). Sordariomycete species are distributed worldwide
68 and have been found in almost every ecosystem (Wang et al. 2018; Luo et al. 2019; Kwon et al. 2021;
69 Maharachchikumbura et al. 2021). Although most Sordariomycetes are saprobic on organic matter from various
70 plants, the class also includes several notorious plant pathogens. For instance, *Pyricularia oryzae* (syn.
71 *Magnaporthe oryzae*; Magnaporthales, Pyriculariaceae), *Fusarium graminearum*, *F. oxysporum* (Hypocreales,
72 Nectriaceae), and *Colletotrichum* species (Glomerellales, Glomerellaceae), are listed in the top 10 fungal plant
73 pathogens (Dean et al. 2012). Moreover, several species, such as *Pyricularia grisea* and *Ophiostoma* spp, were
74 recognized as invasive plant pathogens, which altered the local natural ecosystems (Anderson et al. 2004; Solla
75 et al. 2005). Some species are related to human and animal diseases (Barros et al. 2011; Troy et al. 2013; Tortorano
76 et al. 2014; Řehulka et al. 2016; Jenks et al. 2018), while other species are of great importance to medicine,
77 agriculture, and industry (Crawford et al. 1952; Kaewchai et al. 2009; Xu et al. 2014).

78
79 Diverse lifestyles, including saprotrophic, necrotrophic, hemibiotrophic and biotrophic are present in
80 Sordariomycetes, all of which represent distinct survival strategies evolved by fungi during their interactions with
81 their hosts, companions and associated environments (Presti et al. 2015; Boddy 2016; Rai and Agarkar 2016).
82 Saprotrophs live and feed on non-living organic matter from other organisms, contributing to the global carbon
83 cycle by breaking down complex organic matter into simpler substances (Hobbie and Horton 2007; Mäkelä et al.
84 2014). Fungi of necrotrophic, hemibiotrophic or biotrophic lifestyles are important plant pathogens that pose a
85 serious threat to economically important crops and are responsible for serious losses in quality and yield
86 (Mapuranga et al. 2022). Necrotrophic fungi have a broad host range, and commonly produce diverse toxic
87 molecules (e.g., lytic enzymes, metabolites) to kill host cells and subsequently derive nutrients from dead or dying
88 tissues for growth (van Kan 2006; Mengiste 2012; Singh et al. 2014; Ismaiel and Papenbrock 2015; Newman and
89 Derbyshire 2020). Biotrophic fungi are obligate parasites, which are completely dependent on the living host to
90 complete their life cycles and therefore have to maintain host viability (Glazebrook 2005; Delaye et al. 2013).
91 Hemibiotrophic fungi begin with an early biotrophic phase with their hosts, switching to a necrotrophic lifestyle
92 after killing the host cells (Mendgen and Hahn 2002; Lee and Rose 2010). Endophytic fungi absorb nutrients from
93 plant cells without causing visible symptoms of disease, sometimes in return benefiting plant growth via
94 enhancing the plant's tolerance to abiotic (e.g., drought and salt) and biotic stresses (e.g., insects and other fungal
95 pathogens) (Jia et al. 2016; Phurailatpam and Mishra 2020; Fontana et al. 2021; Wu et al. 2021). In accordance
96 with differences in hosts and substrates, Sordariomycetes are also characterized as plant pathogens, animal
97 pathogens, insect pathogens and mycoparasites. Some fungi are capable of switching between lifestyles.
98 Transitions from the endophytic lifestyle to the pathogenic lifestyle and vice versa have been observed in some
99 important fungal plant pathogens (O'Connell et al. 2012; Rai and Agarkar 2016; Liu et al. 2022).

100

101 Lifestyle-associated genomic traits are a particularly interesting area of research, as pathogenic transitions are
102 highly relevant to gene gain and loss (Friesen et al. 2006; Spanu et al. 2010). *Pyrenophora tritici-repentis*
103 (Pleosporaceae, Pleosporales, Dothideomycetes) becomes highly pathogenic on wheat (*Triticum aestivum*) by

104 obtaining the proteinaceous host-specific toxin *ToxA* from *Stagonospora nodorum* (Phaeosphaeriaceae,
105 Pleosporales, Dothideomycetes), demonstrating that the transfer of the virulence gene is an essential source for
106 the emergence of new pathogens (Friesen et al. 2006). An exclusively biotrophic lifestyle is related to gene losses
107 of primary and secondary metabolic enzymes (Spanu et al. 2010). The convergent losses of decay-related genes
108 and the expansion of symbiosis-related genes are the genetic bases for the evolution of mycorrhizal habits (Kohler
109 et al. 2015). Transposable elements (TEs), also known as “jumping genes,” are crucial genetic factors in both
110 eukaryotic and prokaryotic genomes that shape the evolution of fungal genomes by altering genome plasticity and
111 architecture, interrupting functional genes, generating novel genes or mediating horizontal gene transfer (Lorrain
112 et al. 2021). TEs are critical contributors to fungal pathogenicity by facilitating the diversification of effector genes
113 and even generating novel effector genes (Fouché et al. 2019). In addition, plant symbionts tend to have more TEs
114 than animal parasites (Muszewska et al. 2017a).

115
116 To survive inside a host or a specific environment, fungi must be equipped with the necessary functional proteins
117 to absorb nutrients or to overcome physical and chemical barriers posed by hosts (de Jonge et al. 2011; McCotter
118 et al. 2016; Zeng et al. 2018). The term secretome refers to the complete secretory proteins of an organism, which
119 are released outside the cells to decay substrates and interact with microbes, plants, animals, insects, and other
120 fungi (Eastwood et al. 2011; Frey-Klett et al. 2011; Shang et al. 2015). The fungal secretome comprises various
121 functional groups of protein, including carbohydrate-active enzymes (CAZymes), proteases, lipases, small-
122 secreted proteins (SSPs) and other secretory proteins of unknown functions (Alfaro et al. 2014). Many
123 comparative genomic studies have focused on fungal CAZymes, searching for possible connections between
124 compositions of CAZymes and fungal lifestyles (Kubicek et al. 2014; Pellegrin et al. 2015; Kim et al. 2016; Knapp
125 et al. 2018; Chang et al. 2022). CAZymes include many plant cell wall-degrading enzymes (PCWDEs), and their
126 composition and abundance are often linked to a saprotrophic lifestyle, while this view has been challenged on
127 the grounds that the highest number of CAZymes have been observed in plant pathogenic fungi (Zhao et al. 2013;
128 Kubicek et al. 2014). Fungal effectors, also called virulence factors encoded by avirulence genes, are potent
129 weapons used by fungal pathogens against plant and animal immunity (Stergiopoulos and Wit 2009; Kale and
130 Tyler 2011). Most effectors are secreting cysteine-rich proteins and play an essential role in host-fungal
131 interactions by suppressing host defense responses for promoting host colonization (Lu and Edwards 2016; Wang
132 et al. 2020; Dasari et al. 2018). Some effectors are essential genetic factors in determining host species specificity,
133 which help identify potential pathogenic fungi to certain plants (Li et al. 2020). Effector repositories have been
134 considered to be potential markers for differentiating pathogenic and endophytic strains in the *Fusarium*
135 *oxysporum* species complex (Czislowski et al. 2021).

136
137 Machine learning is a branch of artificial intelligence that is commonly subclassified into unsupervised and
138 supervised methods (Deo 2015). The former is used to find naturally occurring connections or groupings within
139 observations based on little knowledge or even no background information regarding the outcome of the results
140 (Camacho et al. 2018). This is contrasted with the supervised method, which is the construction and optimization
141 of model-based and well-constructed training data with observations and corresponding results (Bzdok et al. 2018).
142 The model is then utilized to predict results of future instances. Both methods have been widely used for
143 unearthing hidden information in big data or complex biological data (Ma et al. 2014; Xu and Jackson 2019).
144 There are many applications of machine learning in species delimitation, such as in, successfully using
145 unsupervised machine learning methods to assign arachnid taxa into species (Derkarabetian et al. 2019),
146 developing a machine learning species identifier for the genus *Hebeloma* (Bartlett et al. 2022) and predicting
147 fungal lifestyles of Dothideomycetes (Haridas et al. 2020). Moreover, machine learning has been used to

148 characterize and classify images of clinically and agriculturally important fungi, which avoids potentially
149 subjective differences, reduces identification time, and lowers the costs (Zieliński et al. 2020; Tongcham et al.
150 2020).

151

152 To mine the association patterns of genomic traits and phylogeny and lifestyles, and further determine whether it
153 is possible to predict lifestyles using machine learning approaches, we carried out a systematic bioinformatic
154 analysis based on 638 Sordariomycete genomes. Firstly, we determined whether the sequencing technologies
155 significantly influence genome assemblies and TE abundance, which exists theoretically and practically but has
156 never been discussed in previous studies. Secondly, based on the study of Fijarczyk et al. (2022), we not only
157 compared the basic genomic traits across multiple lifestyles but also the functional protein groups. Furthermore,
158 we took the influence from phylogeny into account, and compared the difference of numerical genomic traits at
159 different taxonomic levels for determining lifestyle and phylogeny, which is the most important determinant in
160 shaping genomic traits. It is also an answer to resolve the long-standing controversy: whether differences in the
161 secreted proteins reflect phylogeny or pathogenicity (Pellegrin et al. 2015). Finally, we explored whether it is
162 possible to predict fungal lifestyles using machine learning algorithms.

163

164 **Materials and Methods**

165

166 **Genome collection**

167

168 The taxonomic scheme of the class Sordariomycetes has been updated continuously (Maharachchikumbura et al.
169 2015; Hyde et al. 2020; Wijayawardene et al. 2022), whereas the NCBI taxonomy database
170 (<https://www.ncbi.nlm.nih.gov/Taxonomy/Browser/wwwtax.cgi?id=147550>) does not keep up with the updates,
171 and some genomes were assigned incorrect lineage information (Shen et al. 2020; Liu et al. 2022). To ensure the
172 correctness of the taxonomic positions of selected genomes, a taxonomic framework table composed of all generic
173 names in Sordariomycetes and the parent lineage information, was prepared according to the taxonomic outline
174 (Wijayawardene et al. 2022), and some changes were added in keeping with the latest literature (Crous et al. 2021;
175 Sun et al. 2021, Magyar et al. 2022, Sugita & Tanaka 2022). We used the “Ascomycota” as the search term in
176 NCBI’s Genome Browser (<https://www.ncbi.nlm.nih.gov/data-hub/genome/?taxon=4890>, 12 August 2022) to
177 obtain all records of Ascomycota genomes, and then a table, including assembly accession, organism name, strain
178 identifiers, assemble level, and release date, was downloaded. Only records of the Sordariomycetes genome were
179 retained according to the generic names, and the lineage information of the genus was also integrated into the
180 table. These genomes were downloaded via NCBI command line tool datasets. Besides, we collected several
181 genomes from JGI MycoCosm (Grigoriev et al. 2013) with written permission. More details, such as lifestyles,
182 sources, and publication records, were determined by tracing the original literature, the sample details, and the
183 description of the corresponding BioProject records. Strains isolated from soil were marked as saprotrophs. If the
184 strains have two observed lifestyles, only one lifestyle was used as the training data, and the other lifestyle was
185 used to check the predictions. For a small number of strains from certain habits or undetermined sources, we noted
186 them according to the submitter’s description or as “Undetermined” in *Allantophomopsis lycopodina* ATCC 66958
187 (Leotiomyces) was selected as the outgroup.

188

189 **Assessment of genome completeness**

190

191 Assessment of genome quality is the primary step in genomic studies, which is important to recognize potential

192 issues in subsequent analysis (Smits 2019). Benchmarking Universal Single-Copy Orthologs (BUSCO) is an ideal
193 dataset for quantifying genome completeness (Simão et al. 2015) and conducting genome-scale phylogenetic
194 inference (Shen et al. 2018; Shen et al. 2020; Manni et al. 2021). Here, we used BUSCO version 5.2.2 (Manni et
195 al. 2021) with the ascomycota_odb10 database comprising 1,706 reference genes to assess the completeness of
196 the genome assemblies. Only genomes with BUSCO gene content larger than 80% were retained for subsequent
197 analyses.

198

199 **Phylogenetic inference**

200

201 The corresponding protein sequences of single-copy orthologs resulting from the BUSCO analysis were extracted
202 and assembled into a single-locus dataset to conduct phylogenetic analysis. Each locus dataset was aligned using
203 MAFFT version 7.310 (Kato et al. 2002) with options “--auto --maxiterate 1000” that allow the program
204 automatically to determine the approximate refinement strategy and conduct iterative refinement at most 1,000
205 times. Poorly aligned regions were removed using trimAl version 1.4 with the option “-gappycout”, and the
206 alignments with a length shorter than 100 were deleted. ModelFinder (Kalyaanamoorthy et al. 2017) implemented
207 in IQ-TREE2 (Minh et al. 2020) was used to choose the best-fit evolution model of each alignment based on the
208 Bayesian Information Criterion (BIC). All single-locus alignments were concatenated into a supermatrix using an
209 in-house python script. A single evolution model was determined by the occurrence and used in concatenation-
210 based phylogenetic analyses. Maximum-likelihood analysis was conducted using IQ-TREE2 with 1000 bootstrap
211 replicates of the SH-like approximate likelihood ratio test (SH-aLRT) (Guindon et al. 2010) and 1000 bootstrap
212 replicates of ultrafast bootstrap approximation (UFBoot) (Hoang et al. 2017) to estimate the reliability of each
213 internal branch. The strain *Allantophomopsis lycopodina* ATCC 66958 served as an outgroup to root the phylogeny.

214

215 **Identification and analysis of repetitive elements**

216

217 A *de novo* library of repeat consensus sequences was generated for each genome using RepeatModeler version
218 2.0.2 with search engine NCBI-RMBLAST version 2.11.0+. Next, repetitive sequences in genomes were
219 identified and soft-masked using RepeatMasker version 4.1.2 based on three repeat libraries including the *de novo*
220 library, Dfam 2.0 (Hubley et al. 2015), and the Repbase-derived library (20181026) (Bao et al. 2015). The
221 abundance of transposable element (TE) categories was summarized using an in-house python script, and further
222 visualized using the package ggplot2 in R.

223

224 **Recognition the influence of sequencing strategies**

225

226 In this study, the selected genomes mainly were generated from second- and third-generation sequencing
227 technologies. Given their differences in sequencing read length, we had to consider the impact of sequencing
228 technology on the genome, especially in the genome completeness and TE sizes. Therefore, we first excluded the
229 only one genome generated from the first-generation sequencing technology (Sanger sequencing), and divided
230 the other genomes into two groups according their sequencing strategies. If the genome was generated using only
231 the second-generation sequencing technologies, or with Sanger sequencing for improvement, we marked the
232 sequencing strategy of the genome as second-generation sequencing strategy. If the genome was generated using
233 only the third-generation sequencing technologies (Single-molecule real-time sequencing or Nanopore
234 sequencing), or with second-generation sequencing for improvement, we marked the sequencing strategy of the
235 genome as third-generation sequencing strategy. Comparative analyses of the completeness, continuity and TE

236 sizes of genomes generated from both different sequencing strategies, were conducted to figure out whether
237 sequencing strategies impact the number of genes and the abundance of TEs. We also took the taxonomic position
238 of the compared groups into consideration, to decrease the influence of phylogeny to the comparative results.

239

240 **Gene prediction and functional annotation**

241

242 Transfer RNA (tRNA) genes in each soft-masked genome were annotated using tRNAscan-SE version 2.0.9 with
243 default parameters (Chan et al. 2021). Models of protein-coding genes were predicted using the BRAKER2
244 pipeline (Brůna et al. 2021), which combines robust features of GeneMark-EP+ (Brůna et al. 2020) and
245 AUGUSTUS (Stanke et al. 2008). To improve gene prediction accuracy, fungal proteins with annotation scores
246 above 3 in UniProtKB (Consortium 2020) were downloaded and further reduced by removing redundant protein
247 sequences using CD-HIT version 4.8.1 (Fu et al. 2012). Sequence identity and alignment coverage were set to 0.8
248 to retain the representative sequences. Finally, a total of 95,251 protein sequences were used as external evidence
249 for gene structure prediction. Protein hints of homologous regions in each genome were produced using ProtHint
250 version 2.6.0 (Brůna et al. 2020) and further used in the BRAKER2 pipeline. Functional annotation, orthology
251 assignments and domain prediction of all predicted proteins were conducted using eggNOG-mapper version 2.1.3
252 (Cantalapiedra et al. 2021).

253

254 **Identification of secreted proteins and effectors**

255

256 Secreted proteins were identified using a widely used pipeline described previously (Pellegrin et al. 2015;
257 Miyauchi et al. 2020; Mesny et al. 2021). In brief, proteins with signal peptides were identified as candidate-
258 secreted proteins using SignalP version 4.1 with default parameters (Petersen et al. 2011). Then, membrane
259 proteins were removed using TMHMM version 2.0 (Melén et al. 2003) by detecting the presence of the
260 transmembrane helix. Glycosylphosphatidylinositol (GPI)-anchored proteins were removed using NetGPI version
261 1.1 (Gislason et al. 2021) online by detecting GPI-anchoring signals, and proteins residing in the endoplasmic
262 reticulum lumen were removed using PS-SCAN (Nielsen et al. 1997) by detecting KDEL motif (Lys-Asp-Glu-
263 Leu) in the C-terminal region. Two subcellular localization prediction tools, WoLF PSORT (Horton et al. 2007)
264 and TargetP version 2.0 (Emanuelsson et al. 2007), were used to confirm that only proteins assigned extracellular
265 tags were identified as secreted proteins.

266

267 Secreted CAZymes including auxiliary redox (AA) enzyme families were identified using run_dbCAN version
268 3.0.7 (Zhang et al. 2018). Proteases and lipases were identified by querying the MEROPS database (Rawlings et
269 al. 2017) and LED database release 3.0 (<http://www.led.uni-stuttgart.de>), respectively, using BLASTp with a cut-
270 off e-value of 1e-5. Other secreted proteins shorter than 300 amino acids were identified as SSPs and the remaining
271 secreted proteins were marked as OTHER. Secreted effectors were identified using EffectorP version 3.0
272 (Sperschneider and Dodds 2022) with the option of fungal mode. There was no intersection between each group.
273 Furthermore, we followed the grouping criteria in the study of Mesny et al. (2021), and further classified secreted
274 CAZymes into the plant cell wall-degrading enzymes (PCWDEs), fungal cell wall-degrading enzymes (FCWDEs),
275 Cellulose, Hemicellulose, Lignin, Pectin, Peptidoglycan, Mannan, Glucan and Sucrose.

276

277 **Analyses of numerical traits**

278

279 To explore which of the basic components of the genomes and the functional proteins determine the lifestyle, we
280 classified the numerical traits of genome assemblies into two categories and constructed two numerical matrices:
281 basic genomic features and functional protein features. The former includes 25 numerical features: genome size
282 with TEs, genome size without TEs, TE size, GC content of genomes, GC content of genome without TE, GC
283 content of TE, the number of genes, the number of tRNAs, the number of exons, the number of introns, the average
284 length of genes, the average length of tRNAs, the average length of exons, the average length of introns, the
285 average length of intergenic regions, the minimum length of genes, the minimum length of tRNAs, the minimum
286 length of exons, the minimum length of introns, the minimum length of intergenic regions, the maximum length
287 of genes, the maximum length of tRNAs, the maximum length of exons, the maximum of introns and the
288 maximum length of intergenic regions. The latter includes 24 numerical features: total secreted proteins, the
289 effectors, proteases, lipases, SSPs, CAZymes, GHs, GTs, PLs, CEs, AAs, CBMs, PCWDEs, FCWDEs, cellulose-,
290 hemicellulose-, lignin-, pectin-, peptidoglycan-, mannan-, glucan-, chitin-, sucrose-degrading enzymes and other
291 functional proteins. The numbers of these features were summarized using in-house python scripts.

292

293 Correlations were calculated for the two main categories, and details were characterized in the captions of the
294 corresponding figures. To make the comparative analysis more reliable, we excluded the groups with fewer than
295 10 genomes. Overall comparisons were conducted to detect changes in these numerical traits across taxonomic
296 ranks and lifestyles. Post hoc pairwise multiple comparisons were performed to discover how many pairwise
297 comparisons were significantly different based on different grouping criteria and to explore which features were
298 useful in differentiating taxonomic groups and lifestyles.

299

300 **Predicting lifestyles using machine learning algorithms**

301

302 Six commonly used machine learning algorithms for multi-class classification implemented in the python library
303 scikit-learn (<https://scikit-learn.org>): Random Forests (abbreviated as RF), Decision Tree (DT), Naive Bayes
304 (Bayes), Support Vector Machine (SVM), Logistic Regression (LR) and K-Nearest Neighbors (KNN). These
305 algorithms were used to predict fungal lifestyles, and the predictive accuracies of these algorithms were compared
306 to determine the best classifier. Three matrices including the basic genomic features (25 numerical traits),
307 functional protein groups (24 numerical traits), and combined dataset of them (49 numerical traits) were used
308 during the training and prediction stages for selecting the most suitable dataset. The genomes with undetermined
309 lifestyles were excluded from the datasets. First, we standardized the values of features using the function
310 StandardScaler. Next, features with low variances were detected and removed using the function
311 VarianceThreshold with default parameters. Then, the dataset was split into the train (70%) and test subsets (30%)
312 using the function train_test_split, and parameters of the best suitable estimator were determined using the
313 function GridSearchCV. The performance of the estimator was evaluated using the function cross_val_score with
314 5 duplicates based on the test subset. Finally, we used the best estimator to predict the lifestyles of unlabeled
315 genomes.

316

317 **Results**

318

319 **Genome information**

320

321 A total of 638 representative genomes from 5 subclasses, 17 orders, 50 families, 147 genera and 614 species, were
322 selected in this study. More detailed information is described in Supporting Information Table S1. The most

323 numerous subclass is Hypocreomycetidae, which occupies 73.20% (n=467) of the total genomes (Table S2: sheet
324 subclass-count). The 10 most-numerous orders are Hypocreales, Glomerellales, Microascales, Ophiostomatales,
325 Diaporthales, Xylariales, Sordariales, Amphisphaeriales, Magnaporthales, and Coniochaetales in descending
326 order, the number of which range from 3 to 363 (Fig. 1, Table S2: sheet order-count). The other six orders contain
327 only one genome except for three genomes that have not yet been classified in any of the established orders with
328 certainty. Through a comprehensive survey of scientific literature and related databases, we indirectly obtained
329 lifestyle descriptions of most strains (86.99%, n=555) and further classified these strains into eight groups by their
330 host and tropic mode (Table S2: sheet lifestyle-count). We marked the strains isolated from diseased plant tissues
331 as plant pathogens, from decaying woods as saprobes, from insects as entomopathogens, from fungi as
332 mycoparasite, from plant tissues without disease symptoms as endophytes and from diseased human tissues as
333 human pathogens. Moreover, four carnivorous fungi that feed on nematodes were marked as nematophagous fungi,
334 and other genomes that lacked descriptive information regarding lifestyle were marked as “Undetermined”. The
335 most common lifestyle is plant pathogen, which occupies 58.31% (n=372) of the total genomes, followed by
336 saprotrophs at 12.38% (n=79), entomopathogens at 6.27% (n=40), mycoparasites at 3.61% (n=23), endophytes of
337 3.29% (n=21), human pathogens of 2.51% (n=16) and nematophagous fungi of 0.63% (n=4). The remaining 83
338 genomes were temporarily marked as “Undetermined”. We also traced the sequencing technologies of these
339 genomes (Fig. 1, Table S2: sheet wgs-count), and summarized that 74.92% (n=478) of them were sequenced using
340 second-generation sequencing technologies, 24.92% (n=159) were sequenced using third-generation sequencing
341 technologies and only one genome was sequenced using Sanger sequencing technology.

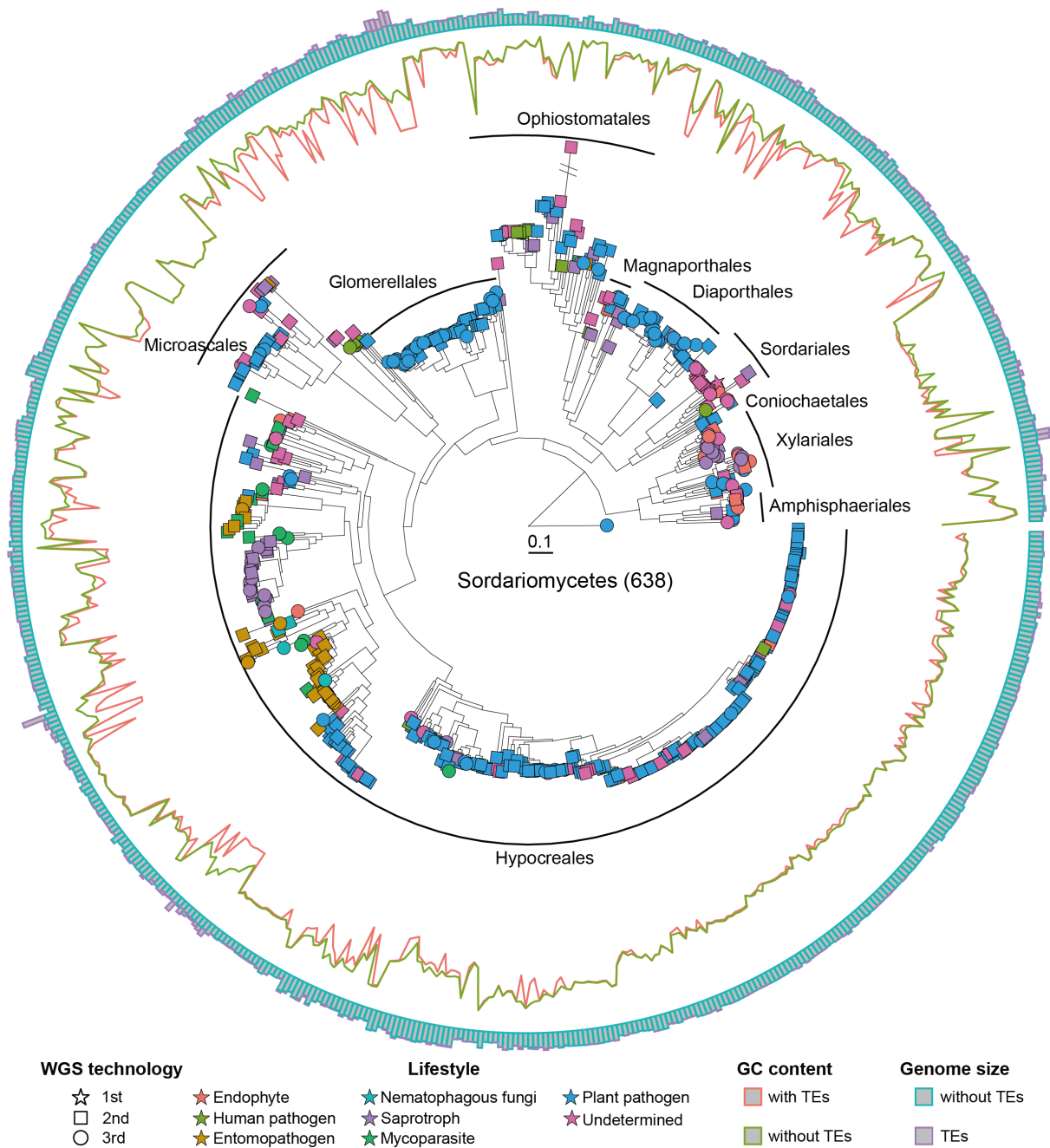
342

343 **Lifestyle occurrences in Sordariomycetes groups**

344

345 Based on the genome data in this study, seven kinds of lifestyles, *viz.* plant pathogens, saprotrophs,
346 entomopathogens, mycoparasites, endophytes, human pathogens and nematophagous fungi determined across 555
347 Sordariomycete genomes but with different occurrences at the subclass, order and family levels (Fig. 1, Table S2:
348 sheet subclass-lifestyle). More diverse lifestyle modes were observed in the more fully sampled groups. For
349 instance, the most-sampled subclasses Hypocreomycetidae and the subordinate order Hypocreales comprise all
350 seven lifestyles, whereas the subclass Sordariomycetidae and Xylariomycetidae only comprise four and three
351 kinds of lifestyles, respectively. At the order level (Table S2: sheet order-lifestyle), the order Ophiostomatales
352 comprise five kinds of lifestyles only inferior to the Hypocreales that includes seven lifestyles. We further compare
353 the occurrence of lifestyles in these two orders at the family level. The family Ophiostomataceae (Ophiostomatales)
354 features with plant pathogens; the family Nectriaceae (Hypocreales) features with plant pathogens; the family
355 Hypocreaceae features with saprotrophs; the family Ophiocordycipitaceae and the family Clavicipitaceae feature
356 with entomopathogens. We compared the distribution of lifestyles at different taxonomic levels (Table S2: sheets
357 lifestyle-subclass, lifestyle-order and lifestyle-family). Endophytes, saprotrophs and plant pathogens are present
358 in four subclasses, followed by human pathogens, present in three subclasses, and entomopathogens and
359 mycoparasites, present in two subclasses. The insufficient sampling lifestyle of nematophagous fungi is only
360 present in the subclass Hypocreomycetidae. At the order and family level, plant pathogen is the most common
361 lifestyle in 11 orders and 29 families, followed by saprotrophs in 9 orders and 19 families, endophytes and in 5
362 orders and 11 families, and human pathogens in 5 orders and 5 families.

363



364

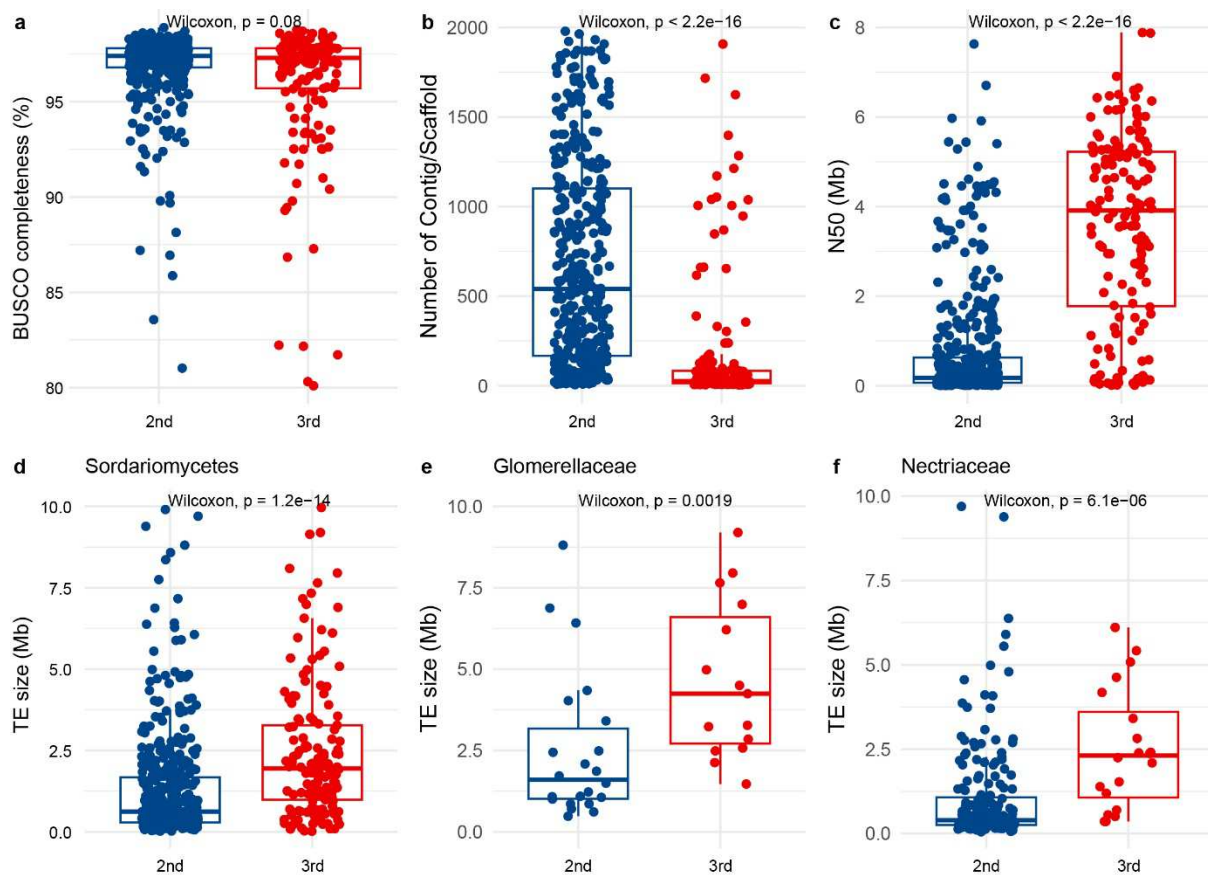
365 **Fig. 1. Maximum likelihood (ML) phylogeny of 638 taxa in the class Sordariomycetes.** The concatenation-
 366 based ML phylogeny ($\ln L = -134,234,602.321$) was reconstructed based on an amino acid dataset of 1,124
 367 BUSCO genes (total of 884,972 sites) under the LG + G4 evolution model. The sequencing strategies are shown
 368 in different shapes (when multiple sequencing strategies were conducted for generating the genomes, we just
 369 marked the sequencing strategy by the most advanced technology). Lifestyles are indicated using different fill
 370 colors. Guanine-cytosine (GC) content of the genome and genome without transposable elements (TEs) are
 371 indicated by a line chart. Genome size and TE sizes are indicated using stacked bar charts. This figure was plotted
 372 using the packages ggtree version 3.4.4 (Yu et al. 2017) and ggtreeExtra version 1.6.1 (Xu et al. 2021) in R (R
 373 Core Team 2022), with the dataset provided in Table S1.

374

375 **Influence of sequencing technologies on TE size**

376

377 The genomes were generated from first-generation, second-generation, and third-generation sequencing platforms,
 378 which account for 0.16% (n=1), 74.92% (n=478), and 24.92% (n=159) of the total number of genomes. To
 379 recognize the potential influences of sequencing technologies on subsequent numerical analysis, we compared the
 380 completeness, continuity, and TE sizes of genomes generated from second- and third-generation sequencing
 381 technologies (Table S3). There is no significant difference ($p = 0.08$) in BUSCO completeness (Fig. 2a), however,
 382 we observed significant differences in the number of contig/scaffold (Fig. 2b, $p < 2.2e-16$) and the N50 value (Fig.
 383 2c, $p < 2.2e-16$), which suggests that the genomes generated from third-generation technologies are better in
 384 genomic continuity than that generated second-generation sequencing technologies. We also investigated whether
 385 the sequencing technologies influence the TE size, and found that the genomes generated from third-generation
 386 sequencing technologies have a larger size of TEs than second-generation sequencing technologies (Fig. 2d). We
 387 compared TE size between the two well-sampled families and the significant differences were also observed in
 388 Glomerellaceae genomes (Fig. 2e, $p = 0.0019$) and Nectriaceae genomes (Fig. 2f, $p = 6.1e-06$). Due to the non-
 389 negligible impact of sequencing technology on TE size, we did not explore further the relationships between
 390 lifestyles and the abundance of TEs. The abundance of TEs is provided in Table S1 and visualized in Fig. S1.
 391



392
 393 **Fig. 2 Comparative analyses of genome completeness, continuities, and TE sizes of genomes generated by**
 394 **second- (2nd) and third-generation (3rd) sequencing strategies. a** Bar plot of BUSCO completeness to
 395 represent the genome completeness; **b, c** Bar plots of the number of contigs/scaffolds and the value of N50 to
 396 represent the continuities. N50 is the shortest contig length that needs to be included for covering 50% of the
 397 genome, which is a measure to indicate the quality of assembled genomes that are fragmented in contigs of
 398 different lengths. The larger number of contigs/scaffolds means a more fragmented genome. The larger N50 value
 399 means a more contiguous genome. **d-f** Bar plots of TE size at the class and family levels to present the influence
 400 of sequencing technologies on TE size. Shapiro-Wilk test was conducted using the function shapiro.test (the

401 package stats) to check whether the compared datasets follow a normal distribution, and the results suggested that
402 these datasets are not normally distributed. Thus, Wilcoxon Rank Sum and Signed Rank Tests were conducted
403 using the function `stat_compare_means` (the package `ggpubr`) to test whether the compared datasets are
404 significantly different ($p \leq 0.05$). All bar plots were plotted using the package `ggpubr`. For visualization, a small
405 number of data points above 2,000 in subfigure b, data points above 8 Mb in subfigure c, and data points above
406 10 Mb in subfigure d and e, are not displayed. The input dataset is given in Table S1, and all resulting tables are
407 given in Table S3. Statistical analyses and visualization were done in R (R Core Team 2022).

408

409 **Variations of basic genomic features**

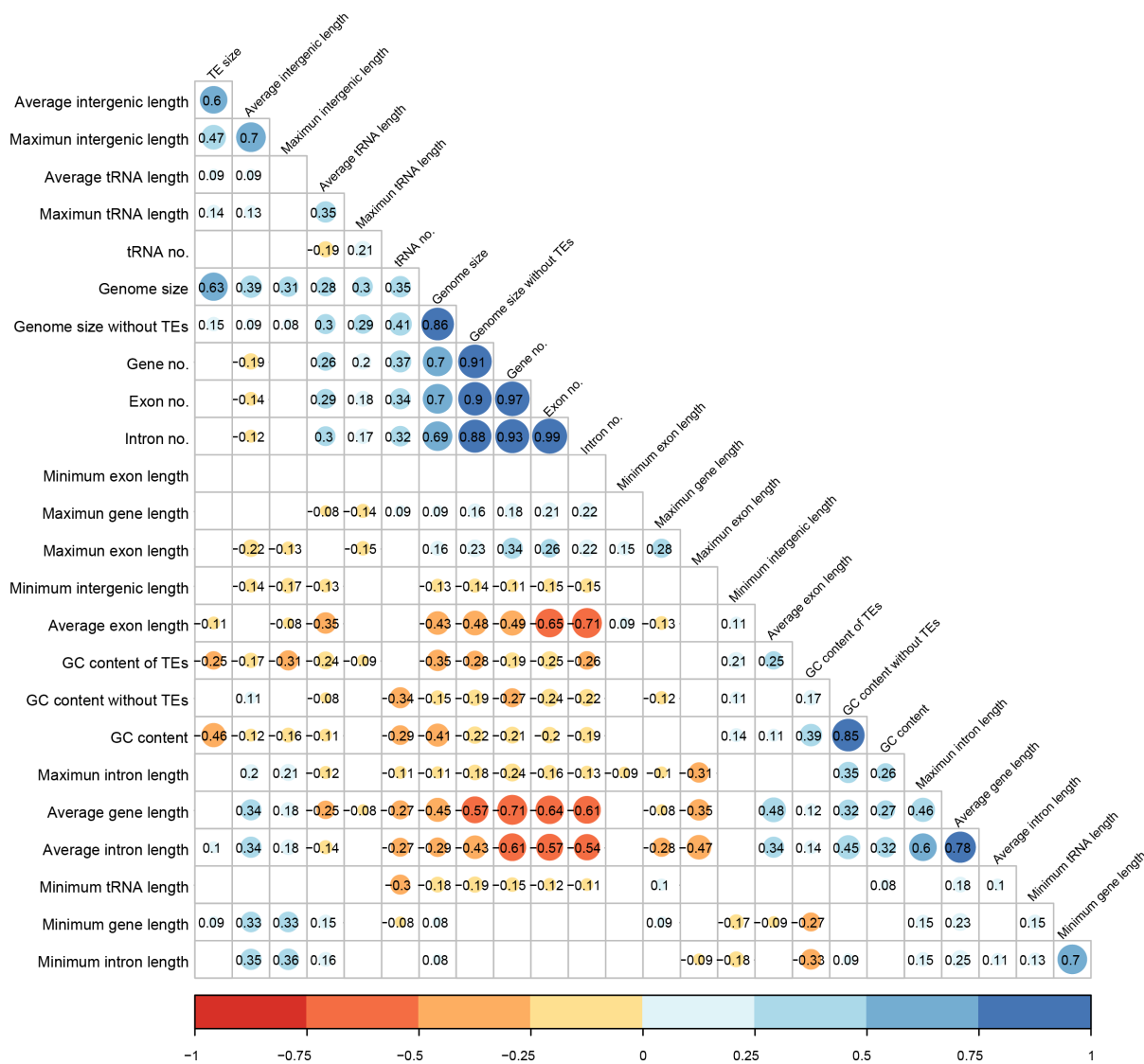
410

411 We counted a total of 25 basic genomic features, which are summarized in Table S1. Results of correlation
412 analyses among these features suggested that some features are highly correlated (Fig. 3, Table S4). Genome size
413 is positively correlated with TE size with a Pearson's correlation coefficient of 0.63, which is smaller than its
414 correlation coefficient with the genome size without TEs ($r = 0.86$), suggesting that the TEs can increase the
415 genome size but not the dominant factor. GC content is positively correlated to the GC content without TEs ($r =$
416 0.85) but negatively related to the TE size ($r = -0.46$). In addition, GC content with TEs or without TEs is
417 influenced by TE size, the larger TE size caused the larger difference between them, suggesting that TEs decrease
418 the GC content of genomes. Genome size without TE is positively correlated to the number of genes ($r = 0.91$),
419 the number of exons ($r = 0.90$), and the number of introns ($r = 0.88$). The latter two features, exons, and introns
420 are important structural components of genes, the numbers of which reasonably displayed high correlations with
421 the number of genes ($r = 0.97$; $r = 0.93$). The average length of genes is correlated to the average length of introns
422 ($r = 0.78$) and the exons ($r = 0.48$), indicating that changes in intron length are the main cause of the variation of
423 gene length compared to the exon. TE size is positively correlated to the average and maximum lengths of
424 intergenic regions ($r = 0.60$; $r = 0.47$), but not displays significant correlations with gene structures including gene
425 length, exon length, and intron length, suggesting that TEs are the main factor to change the distance between
426 genes without significant influence on the gene structures. The minimum and maximum length of multiple features
427 (genes, intergenic regions, introns, exons) exhibit relatively low correlation with other features, or correlations are
428 not significant, except for the maximum length and the average length of intergenic regions ($r = 0.70$), the
429 maximum length and the average length of introns ($r = 0.6$) and the minimum length of introns and genes ($r =$
430 0.7). Overall, most basic genomic features display a low correlation with each other, suggesting some of which
431 are stable and independent in evolution.

432

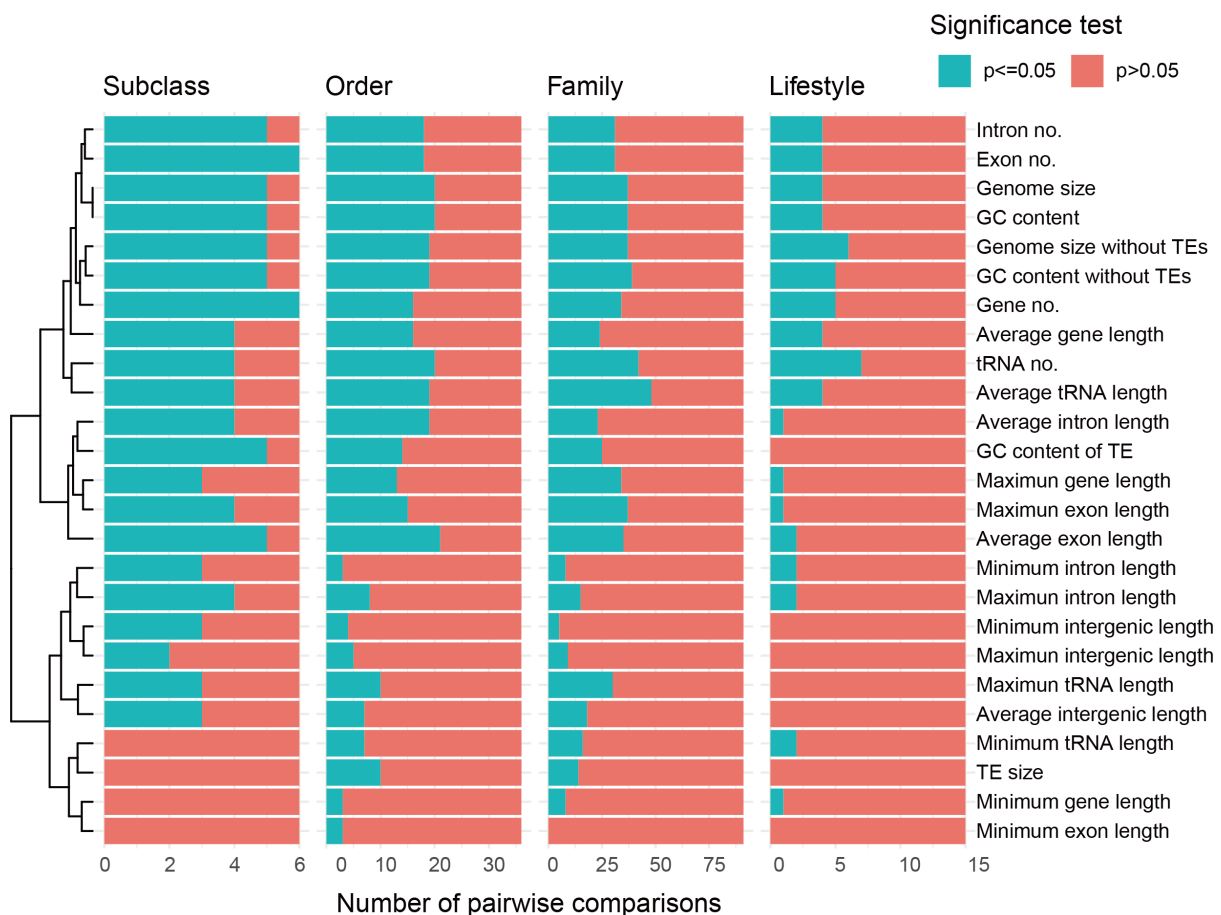
433 We also compared the group means of these 25 genomic features over all groups of different taxonomic ranks and
434 lifestyles (Table S5). We observed overall statistically significant differences in most genomic features (22/25) at
435 the subclass level, excluding the minimum length of exons, TE sizes, and the minimum length of tRNAs (Table
436 S5: sheet subclass). The minimum length of exons is the only feature that does not show a significant difference
437 at the order level (Table S5: sheet order). And at the family, all features display significant differences (Table S5:
438 sheet family). Considering the groups with different lifestyles, there are 6 genomic features without significant
439 difference (Table S5: sheet lifestyle), which are the minimum length of exons, the average length of intergenic
440 regions, the minimum length of intergenic regions, the size of TEs, the GC content of TEs and the maximum
441 length of tRNAs. In paired comparison analysis (Fig. 4), we included 4 subclasses Diaporthomycetidae,
442 Hypocreomycetidae, Sordariomycetidae and Xylariomycetidae, which formed 6 pairwise comparisons, 5 of which
443 are significantly different in most of the features (Table S5: sheet pairwise-subclass). Specially, the number of
444 genes and the number of exons display the most powerful resolution to differentiating the taxonomic groups at

445 the subclass level. At the order level (36 pairwise comparisons in total) and family level (91 pairwise comparisons
 446 in total), we observed a clear downward trend of significant differences, suggesting that all features lack
 447 resolutions at lower taxonomic levels (Table S5: sheets pairwise-order and pairwise-family). However, fairly low
 448 proportions of significantly different comparisons (15 pairwise comparisons in total) were observed across all
 449 features in the groups with different lifestyles (Table S5: sheet pairwise-lifestyle). Moreover, clustering analysis
 450 shows that several features (TE size, the minimum length of tRNAs, the minimum length of exon, and the
 451 minimum length of gene) display little usefulness in distinguishing different taxonomic groups, and most features
 452 are useless in differentiating different lifestyles.
 453



454
 455 **Fig. 3 Correlation analysis of 25 basic genomic features.** Ladder heatmap of Pearson correlation coefficients
 456 of all pairwise genomic features. The colors and values in small squares indicate the degree of positive correlation
 457 (red) or negative correlation (blue). No significant correlated comparisons ($p > 0.05$) were displayed in white and
 458 blank squares. Pearson correlation coefficients were calculated using the function cor (the package stats), and the
 459 significance test was conducted using the function cor.mtest (the package corrplot). The figure was plotted using
 460 the package corrplot with the resulting datasets in Table S4. Values of these 25 features are provided in Table S1.
 461

462 Although, not all features show strong discrimination in distinguishing one group from the other groups, a high
 463 proportion of significant differences of some genomic features was observed in specified comparisons. For
 464 instance, at the subclass level (Table S5: sheet class-class), there are 18, 17, 15, 15 and 15 significantly different
 465 features present in the pairwise comparisons of Hypocreomycetidae-Xylariomycetidae, Hypocreomycetidae-
 466 Sordariomycetidae, Diaporthomycetidae-Hypocreomycetidae, Diaporthomycetidae-Xylariomycetidae and
 467 Sordariomycetidae-Xylariomycetidae. Likewise, a high proportion of some features also were observed at the
 468 Order and Family levels (Table S5: sheets order-order and family-family). These results suggest that some features
 469 are useful in differentiating specified taxonomic groups, especially in phylogenetic distant comparisons. As for
 470 lifestyles, the largest difference in genomic features was only observed in the comparisons of entomopathogens-
 471 plant pathogens (14/25), followed by entomopathogens-endophytes (10/25), and the rest of comparisons display
 472 no difference or very small differences, especially in the comparisons of endophytes-saprotrophs (0/25),
 473 entomopathogens-mycoparasites (0/25), mycoparasites-saprotrophs (0/25), endophytes-mycoparasites (1/25),
 474 human pathogens-mycoparasites (1/25), endophytes-saprotrophs, human pathogens-plant pathogens (1/25), and
 475 human pathogens-saprotrophs (1/25) (Table S4: sheet lifestyle-lifestyle). It suggests that based on these basic
 476 genomic features it is difficult to differentiate compared lifestyles. In another word, we could not correctly assign
 477 a lifestyle label to a new taxon with very similar genomic features, to endophytes, saprotrophs, mycoparasites and
 478 entomopathogens.
 479



480
 481 **Fig. 4 Resolution powers of 25 basic genomic features in differentiating different taxonomic groups and**
 482 **lifestyles.** Stacked bar plots of the number of significantly (orange; $p \leq 0.05$) and non-significantly (green; $p >$
 483 0.05) different comparisons across all features based on their taxonomic ranks and lifestyles. The cluster analysis
 484 was performed using the function `dist` (the package `stats`) with the dataset in Table S4 sheet: clustering-matrix to

485 obtain a Euclidean distance matrix, then using the function hclust (the package stats) to cluster these features with
486 the “complete” agglomeration method. All datasets are given in corresponding sheets in Table S5.

487

488 **Overview of functional protein groups**

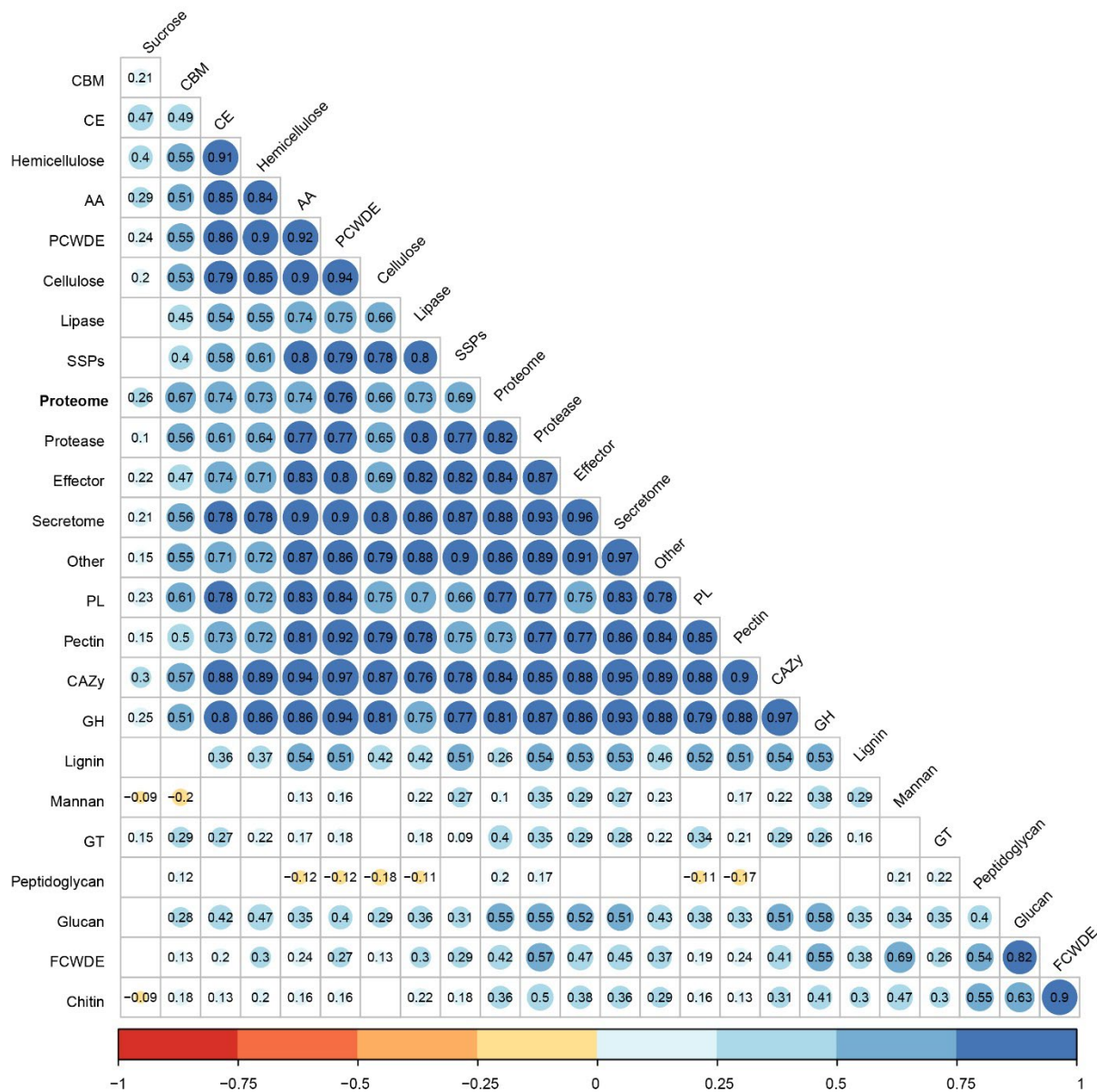
489

490 A total of 24 functional protein groups were summarized in Table S1 and visualized in Fig. S2. To explore the
491 correlation between the number of the proteome and the number of each functional protein group we include the
492 feature of proteomes, which is equivalent to the number of protein-coding genes in the last part during correlation
493 analysis (Fig. 5; Table S6). The result shows that 66.67% (16/24) of protein groups are highly positively correlated
494 ($r > 0.6$) with the total number of the proteome. The main subgroups of the secretome, the number of CAZymes,
495 protease, lipase, SSPs, secreted effectors and other functional proteins are highly positively correlated with the
496 total number of secretomes with the Pearson correlation coefficient of 0.95, 0.93, 0.86, 0.87, 0.96 and 0.97,
497 respectively. The six subgroups of CAZymes display varying degrees of correlation with the total number of
498 CAZymes. The numbers of AAs, GHs, CEs and PLs display high correlation with the Pearson correlation
499 coefficient of 0.97, 0.97, 0.88 and 0.88, respectively. The number of CBMs displays a relatively high correlation
500 ($r = 0.57$) with CAZymes, whereas the GTs display a low correlation ($r = 0.29$) with CAZymes. As for the more
501 specified functional subgroups of CAZymes, the numbers of PCWDEs, pectin-degrading enzymes, hemicellulose-
502 degrading enzymes, and cellulose-degrading enzymes, are highly correlated with the total number of CAZymes
503 with the Pearson correlation coefficients of 0.97, 0.90, 0.89 and 0.87, respectively, followed by lignin-degrading
504 enzymes and glucan-degrading enzymes with relatively high correlation coefficients of 0.54 and 0.51. The
505 numbers of FCWDEs, chitin-degrading enzymes and mannan-degrading enzymes display relatively low
506 correlation with CAZymes, the correlation coefficients of which are 0.41, 0.31 and 0.22 respectively, and no
507 significant correlation was observed between peptidoglycan-degrading enzymes and CAZymes. We also noticed
508 the high correlations between several specified functional subgroups of CAZymes, such as FCWDEs and chitin-
509 degrading enzymes with correlation coefficients of 0.9, FCWDEs and glucan-degrading enzymes with correlation
510 coefficients of 0.82, which are mainly due to the overlapping functional proteins (Table S6). Compared with the
511 correlation matrix of genomic features (Fig. 3), most functional proteins are more stable in number, showing a
512 trend of co-evolution except for mannan-degrading enzymes, GTs, and peptidoglycan-degrading enzymes.

513

514 The discrimination of these 24 functional protein groups was visualized by comparing the numbers of significantly
515 different pairwise comparisons and not significantly different pairwise comparisons (Fig. 6, Table S7). Compared
516 with the discrimination of 25 basic genomic features, clear increases in functional protein groups are observed at
517 the taxonomic levels and lifestyles. At the subclass level, more than half (15/24) of these protein groups are
518 powerful in differentiating subclasses ($n > 3$, Table S7: sheet cluster-matrix), especially the number of CBMs and
519 mannan-degrading enzymes with 100% resolution (Table S7: sheet pairwise-subclass). However, CEs,
520 hemicellulose-degrading enzymes and PCWDEs display very poor resolution, especially the latter two. At the
521 order and family levels (Table S7: sheets pairwise-order and pairwise-family), the numbers of significantly
522 different pairwise comparisons increase with the total number of pairwise comparisons, but the proportion of
523 significantly different pairwise comparisons for each protein group decreases, most notably in CBMs and mannan-
524 degrading enzymes. Although the numbers of PCWDEs and hemicellulose-degrading enzymes are useless in
525 differentiating subclasses, we notice that PCWDEs can distinguish more than half of the pairwise comparisons at
526 the order level (23/36) and the family level (48/91), and hemicellulose-degrading enzymes can distinguish more
527 than half of the pairwise comparisons at the order level (19/36) and nearly half at the family level (39/91). On the
528 subject of lifestyles (Table S7: sheet pairwise-lifestyle), we observed clear drops in the proportion of significantly

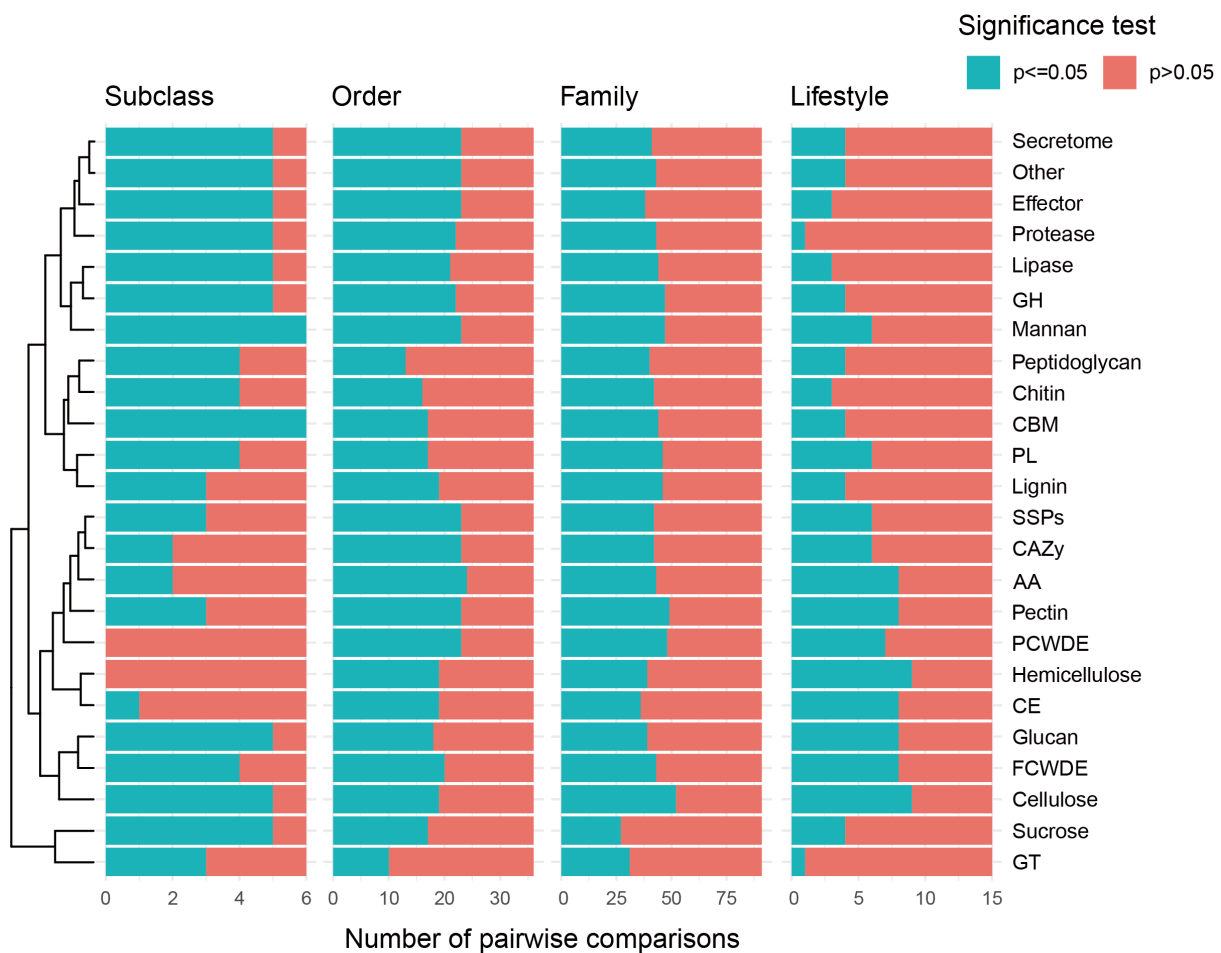
529 different pairwise comparisons for some protein groups, and also noticed some increased proportions, such as the
 530 glucan-, cellulose- and hemicellulose-degrading enzymes.
 531



532
 533 **Fig. 5 Correlation analysis of 24 functional protein groups and proteomes.** Ladder heatmap of Pearson
 534 correlation coefficients of all pairwise genomic features. The colors and values in small squares indicate the degree
 535 of positive correlation (red) or negative correlation (blue). No significant correlated comparisons ($p > 0.05$) were
 536 displayed in white and blank squares. Pearson correlation coefficients were calculated using the function cor (the
 537 package stats), and the significance test was conducted using the function cor.mtest (the package corrplot). The
 538 figure was plotted using the package corrplot with the resulting datasets in Table S6. Values of these 24 functional
 539 protein groups and the total number of proteomes are provided in Table S1.

540
 541 We also counted the number of significantly different protein groups in each pairwise comparison. At the class
 542 level (Table S7: sheets subclass-subclass), the most notable subclass is Xylariomycetidae, which has 17
 543 significantly different protein groups with Diaporthomycetidae, 16 with Hypocreomycetidae and
 544 Sordariomycetidae. The smallest difference is observed in the pairwise comparison of Diaporthomycetidae and

545 Sordariomycetidae with 12 significantly different protein groups. In other words, Xylariomycetidae is the easiest
 546 to distinguish from other subclasses. At the order level (Table S7: sheet order-order), the most notable order is
 547 Ophiostomatales, which has 22 significantly different protein groups with Glomerellales and Hypocreales, 21 with
 548 Amphisphaeriales, 20 with Diaporthales, 19 with Magnaporthales. The smallest differences are observed in the
 549 pairwise comparisons of Magnaporthales-Amphisphaeriales, and Magnaporthales-Diaporthales. Moreover,
 550 Magnaporthales has only 2 significantly different protein groups with Glomerellales, 4 with Hypocreales and
 551 Xylariales, indicating that it is not easy to distinguish Magnaporthales from the compared orders based on most
 552 functional protein groups. At the family level (Table S7: sheet family-family), the largest number of significantly
 553 different protein groups is 23, which is observed in three pairwise comparisons of Ceratocystidaceae-Nectriaceae,
 554 Glomerellaceae-Ophiostomataceae and Nectriaceae-Ophiostomataceae. Inversely, the smallest number is 1,
 555 which is observed in two pairwise comparisons of Bionectriaceae-Nectriaceae and Clavicipitaceae-
 556 Ophiocordycipitaceae. For lifestyles (Table S7: sheet lifestyle-lifestyle), plant pathogens are the easiest lifestyle
 557 to distinguish from saprotrophs, entomopathogens and mycoparasites, and they have 21, 20, and 17 significantly
 558 different protein groups respectively. At the same time, it is also the most difficult to distinguish from endophytes
 559 because that they only significant difference in the abundance of PCWDEs and peptidoglycan-degrading enzymes.
 560 No significantly different protein group is present in the comparison of endophytes-saprotrophs, indicating that
 561 we cannot differentiate them based on the number of functional protein groups.
 562



563
 564 **Fig. 6 Contributions of 24 functional protein groups in differentiating different taxonomic groups and**
 565 **lifestyles.** Stacked bar plots of the number of significantly (orange; p <= 0.05) and non-significantly (green; p >
 566 0.05) different comparisons across all features based on their taxonomic ranks and lifestyles. The cluster analysis

567 was performed using the function `dist` (the package `stats`) with the dataset in Table S7 sheet: `cluster-matrix`, to
568 obtain a Euclidean distance matrix, then using the function `hclust` (the package `stats`) to cluster these features with
569 the “complete” agglomeration method. All datasets are given in corresponding sheets in Table S7.

570

571 **Predicting lifestyles using machine learning approaches**

572

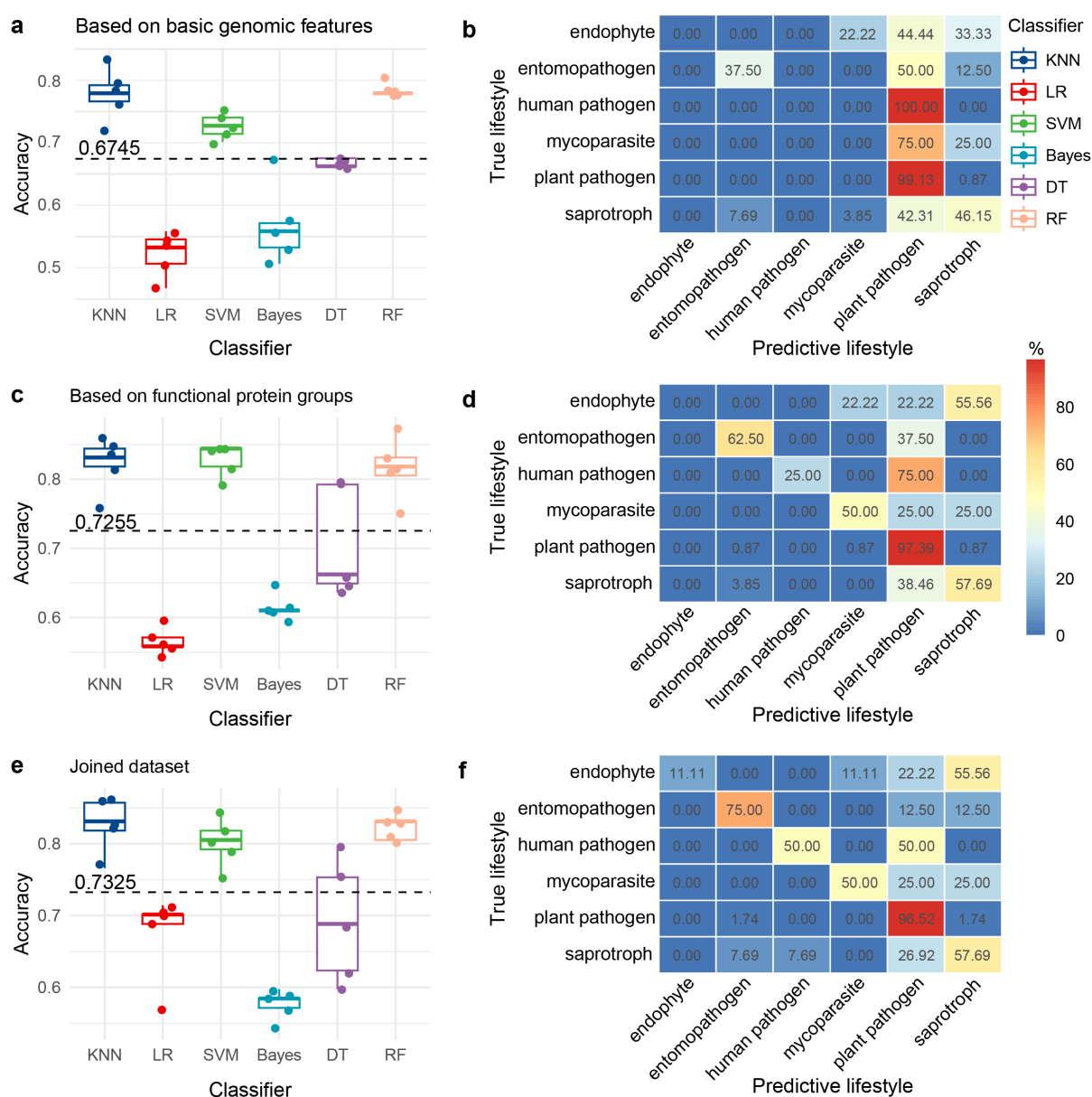
573 Predictive models of six commonly used machine learning algorithms were trained and optimized based on the
574 training subsets of three different datasets, and accuracies in predicting fungal lifestyles were compared and
575 visualized in Fig. 7 (Tables S8). For the dataset of basic genomic features, RF is the best classifier with an average
576 accuracy of 0.7844, followed by KNN (0.7766), SVM (0.7272), DT (0.6675) and Bayes (0.5688); LR is the worst
577 with an average accuracy of 0.5221. For the dataset of the functional protein groups, SVM is the best classifier
578 with an average accuracy of 0.8286, followed by KNN (0.8208), RF (0.8156), DT (0.7065) and Bayes (0.6156);
579 LR is the worst with an average accuracy of 0.5662. For the combined dataset including a total of 49 numerical
580 features, KNN is the best classifier with an average accuracy of 0.8260, followed by RF (0.8234), SVM (0.8026),
581 DT (0.6909) and LR (0.6753); Bayes is the worst with an average accuracy of 0.5766. In terms of machine learning
582 algorithms, KNN, SVM and RF perform better than LR, Bayes and DT in predictive accuracies across the three
583 datasets (Fig. 7 a, c, e). Bayes, DT, RF and SVM obtained the highest-average accuracies based on the functional
584 protein groups, and the other two methods, KNN and LR, were based on the combined datasets. We noticed that
585 all classifiers obtained the worst-average accuracies based on the basic genomic feature alone, and increased
586 accuracies were observed based solely on a functional protein dataset or combined dataset (Fig. S3), indicating
587 that numerical traits of functional protein groups are more useful than basic genomic features for predicting fungal
588 lifestyles.

589

590 Based on the test subsets, we tested the performance of the three best classifiers, RF for the dataset of basic
591 genomic features and combined dataset and SVM for the functional protein groups. For the dataset of basic
592 genomic features (Fig. 7b), we noticed that 99.13% of plant pathogens were assigned the correct lifestyles,
593 suggesting that RF is reliable for distinguishing plant pathogens from other lifestyles. However, it performed
594 worse in differentiating endophytes, human pathogens and mycoparasites from other lifestyles. Predictive results
595 of all endophytes, human pathogens and mycoparasites did not match the assigned lifestyles that we determined
596 by a literature survey or the genomic descriptions. About half of endophytes (44.44%) were incorrectly predicted
597 as plant pathogens, and some other genomes were incorrectly recognized as saprotrophs and mycoparasites. Of
598 mycoparasites, 75% were incorrectly predicted as human pathogens and 25% as saprotrophs. Of human pathogens,
599 all of them were incorrectly predicted as plant pathogens. As for the other three lifestyles, RF obtained an increased
600 accuracy. Of entomopathogens, 37.5% were correctly classified, and 50% were incorrectly predicted as plant
601 pathogens and 12.5% as saprotrophs. Of plant pathogens, 99.13% were correctly classified, and the rest were
602 incorrectly predicted as saprotrophs (0.87%). Of saprotrophs, 46.15% were correctly predicted as saprotrophs,
603 42.13 % were incorrectly predicted as plant pathogens, 7.69% as entomopathogens, 3.85% as mycoparasites.
604 Concerning the dataset of function protein groups (Fig. 7d), we used the SVM algorithm and observed a clear
605 improvement in differentiating entomopathogens (37.50% to 62.50%), human pathogens (0 to 25%),
606 mycoparasites (0 to 50%) and saprotrophs (46.15% to 57.69%) from other lifestyles. Compared with RF
607 predication based on the dataset of genomic features, SVM resulted in the same incorrect results in differentiating
608 human pathogens, with a similar result for endophytes and slightly decreased accuracy in predicting saprotrophs.
609 As for the combined dataset (Fig. 7f), the KNN algorithm was used to predict lifestyles, and we observed a clear
610 improvement in predictive accuracies for endophytes, entomopathogens and human pathogens.

611
 612
 613
 614
 615
 616
 617
 618
 619
 620
 621
 622

Based on the combined dataset, we obtained the highest average accuracy of 0.7325 for the six algorithms. It was seen that KNN was the best classifier with an average accuracy of 0.8260. Therefore, we used KNN to conduct the prediction of 83 Sordariomycetes genomes with undetermined lifestyles, and the predicted lifestyles with probabilities were listed in Table S9. KNN classified these 83 genomes into 5 lifestyles, 1 endophyte, 4 human pathogens, 4 entomopathogens, 6 saprotrophs, and 68 plant pathogens. We further checked the taxonomic positions of strains, and only one endophyte is distributed in the family Bionectriaceae; 4 human pathogens in Sordariaceae; 4 entomopathogens in Ophiocordycipitaceae, Ophiostomataceae, Clavicipitaceae and Cordycipitaceae; 6 saprotrophs in Bionectriaceae, Diatrypaceae, Sordariaceae and Hypoxylaceae; and 68 plant pathogens in 20 families. We traced the lifestyles of phylogenetically closed groups with predicted genomes, and most of the observed lifestyles were consistent with our predictions.



623
 624
 625

Fig. 7 Lifestyle predictions using machine learning methods. a Boxplots of predictive accuracies using six machine learning algorithms for predicting fungal lifestyles based on the train subset of the basi

626 c genomic features. **b** Confusion matrix, a performance matrix, to evaluate the performance of the best
627 classifier (RF, average accuracy = 0.7844) in predicting fungal lifestyles based on the test subset of th
628 e basic genomic features. **c** Predictive accuracies of the six commonly used machine learning algorithm
629 s based on the train subset of the functional protein features. **d** Confusion matrix of the best classifier
630 (SVM, average accuracy = 0.8286) in predicting fungal lifestyles based on the test subset of the functi
631 onal protein groups. **e** Predictive accuracies of the six commonly used machine learning algorithms bas
632 ed on the combined datasets. **f** Confusion matrix of the best classifier (KNN, average accuracy = 0.82
633 60) in predicting fungal lifestyles based on the test subset of the functional protein groups. For the co
634 nfusion matrix, the diagonal elements show the proportion of correctly classified genomes, while the of
635 f-diagonal elements show the number of misclassified genomes.

636

637 **Discussion**

638

639 **Diverse lifestyles but unbalanced whole genome sequencing**

640

641 Sordariomycetes has a large number of available genome sequences for an ascomycetes class in public databases;
642 however, many of these genomes are restricted to economically important groups such as plant pathogens
643 (*Fusarium*, *Diaporthe*, *Calonectria*, *Claviceps*, *Collectotrichum*), entomopathogens (*Cordyceps*, *Metarhizium*,
644 *Ophiocordyceps*, *Tolyposcladium*), mycoparasites (*Clonostachys*), human pathogens (*Sporothrix*, *Sarocladium*,
645 *Scedosporium*) model organism (*Neurospora*), and biocontrol and secondary metabolites producers (*Trichoderma*,
646 *Daldinia*, *Xylaria*). For instance, Hypocreomycetidae includes plant pathogens, entomopathogens, mycoparasites,
647 human pathogens and biocontrol agents and is responsible for 73.20% of the total Sordariomycete genome used
648 in this study. However, the Sordariomycetes include other ecologically important saprotrophs, epiphyllous,
649 hypophyllous, facultative lichenised, fungicolous and extreme inhibiting groups primarily overlooked due to their
650 economically insignificance. Therefore, the current genomic data are largely incomplete and cannot be used to
651 make reliable conclusions about the overall lifestyle of Sordariomycetes fungi. Saprobies are the most common
652 type of fungi, and Sordariomycetes now comprises 195 families, and 171 have a saprobic lifestyle. This is true as
653 many of these fungi can degrade polymers of varying complexity by releasing extracellular enzymes that break
654 down plant and animal debris. We suspect that saprobic Sordariomycete families will likely be more than this as
655 the remaining families are poorly sampled or monotypic. Plant pathogens are the second most abundant lifestyle
656 in Sordariomycetes, distributed over 93 families. The five largest Sordariomycetes orders, Diaporthales,
657 Glomerellales, Hypocreales, Microascales and Ophiostomatales, each contain a large number of highly destructive
658 plant pathogens. These include some of the most important diseases of the cereal (rice, wheat, barley and maize)
659 ornamental, fruit, vegetable and wild crops (Chang et al. 2018; Talhinhas and Baroncelli 2021; Liu et al. 2022;
660 Han et al. 2023). Endophytes are distributed over 40 families of Sordariomycetes. There is publishable evidence
661 that fungal endophytes can switch lifestyles to saprotrophs and pathogens and vice versa (Promputtha et al. 2007;
662 Promputtha et al. 2010). Human pathogens, entomopathogens, mycoparasites and nematophagous fungi are
663 distributed over 17, 11, 5 and 2 families of Sordariomycetes, respectively. The least distributed nematophagous
664 fungi are only present in Hypocreales families Clavicipitaceae and Ophiocordycipitaceae. Their diverse lifestyles
665 and ability to switch to other life modes and inhibit diverse ecological niches that include extreme environmental
666 constraints allow Sordariomycetes to adapt and distributed over all ecosystems on earth and to be the second
667 largest ascomycetes class.

668

669 **Influence of sequencing technologies on genome assemblies**

670

671 High-quality genome assemblies are fundamental for genomic studies. Therefore, when we used genomes from
672 public databases, we were meticulous in checking their quality that was inevitably affected by the methods of
673 DNA extraction (Nouws et al. 2020), sequencing technologies (Lang et al. 2020; Murigneux et al. 2020) and
674 assembly algorithms (Miller et al. 2010; Meng et al. 2022). As a user of public genomes, although we cannot
675 improve genome assemblies by optimizing these steps, recognizing the inaccuracies of genome assemblies
676 reduces the possibility of drawing incorrect conclusions. Repetitive DNA sequences present technical challenges
677 for assembly algorithms by bringing in ambiguous alignment during genome assemblies, leading to biases and
678 errors in final assembly results (Treangen and Salzberg 2012; Tørresen et al. 2019). For instance, fungal ribosomal
679 RNA genes (rDNA) as multiple-copy segments organized in tandem arrays exist in genomes (Cooper 2000). Each
680 repeat unit (18S rRNA-internal transcribed spacer 1-5.8S rRNA-internal transcribed spacer 2-28S rRNA-
681 intergenic spacer) is approximately 9kb in length (SONE et al. 2000; Salim et al. 2017), which far exceeds the
682 read length limit of second-generation sequencing, and the reads generated from second-generation sequencers
683 cannot span this kind of long repetitive sequence (Treangen and Salzberg 2012). Assembly algorithms, such as
684 the Greedy strategy, Overlap-Layout-Consensus strategy, and de Bruijn graph strategy, tend to assemble these
685 highly similar or identical sequences into single, collapsed contig (Treangen and Salzberg 2012). Although third-
686 generation sequencing technologies, also called long-read sequencing technologies, can overcome the read length
687 limit by producing 20–200 kb reads (Goodwin et al. 2016), the high cost per genome hinders its widespread
688 application, especially in some fungal species that lack of direct economic interest. Furthermore, our previous
689 study (Chen et al. 2022) found that second-generation sequencing technologies can provide reliable genome
690 assemblies for phylogenomic analyses, which focus on protein-coding genes rather than repetitive sequences. In
691 this study, we included 638 genomes, most of which were generated using second-generation sequencing
692 technologies (n=478, 74.92%). We set the completeness threshold at 80% to remove the unreliable genomes, and
693 confirmed that each group included at least 10 genomes during statistical analyses. Hence, we believe that
694 sequencing strategies did not influence the numerical traits meaningfully.

695

696 TEs are mobile genetic elements that are composed of diverse members, including short interspersed nuclear
697 elements (SINEs), Helitrons, Alus, endogenous retroviruses (ERVs), DNA transposons and retrotransposons
698 (Wicker et al. 2007). The ability to move and their repetitive nature make TEs key drivers of genome evolution
699 (Dhillon et al. 2019; Senft and Macfarlan 2021). Many studies have shown that the expansion of TEs resulted in
700 a significantly expanded genome in fungal species, such as in *Cenococcum geophilum* (Peter et al. 2016),
701 *Zymoseptoria tritici* (Oggenfuss et al. 2021) and *Lactarius* species (Lebreton et al. 2022). Large-scale genomic
702 location analysis of TEs has indicated that most TEs are evolutionarily neutral, but animal-related and pathogenic
703 fungi include more TEs inserted in genes compared to fungi with other lifestyles (Muszewska et al. 2019).
704 Kirkland et al. (2018) reported that *hAT* or *Gypsy* TEs located within 1kb of protein-coding genes can decrease
705 the expression of related genes. LTR retrotransposons, a class I transposable element, inserted in the *MFS1*
706 promoter region, resulted in *MFS1* overexpression and the presence of multidrug resistance phenotype in the
707 wheat pathogen *Zymoseptoria tritici* (Omrane et al. 2017). TEs are important and biologically functional repetitive
708 sequences, the abundance of which in genomes is inevitably affected by sequencing technologies, especially by
709 second-generation sequencing technologies. In this study, we recognized that TE sizes in the genomes generated
710 from second-generation sequencing technologies are significantly smaller than those from third-generation
711 sequencing technologies. We also discovered the GC content of TEs is significantly lower than other regions in
712 the genomes, and that TE sizes are negatively correlated with the overall GC content of fungal genomes. Hu et al.
713 (2022) showed that GC content is positively correlated with growth temperature in prokaryotes, and Šmarda et al.

714 (2014) reported that increased GC content helps plants adapt to seasonally cold and/or dry climates. Considering
715 the clear influence of sequencing technologies, the true abundance of TEs in most genomes has been
716 underestimated in previous studies and in this study. Therefore, instead of providing a more in-depth analysis, we
717 only compared the abundance of TEs in multiple groups and displayed their diversity in Table S1 and Fig. S2. We
718 did not observe a significant difference in TE sizes between lifestyles; thus, the underestimated abundance in this
719 study did not affect our statistical and predicted results. However, future studies related to TEs should take into
720 account the influence of sequencing technologies.

721

722 **Effectors are not a reliable indicator for disease-related fungi but are useful for differentiating specific** 723 **lifestyles**

724

725 Effectors, important virulence factors secreted by bacteria (Yu et al. 2020), fungi (Stergiopoulos and Wit 2009),
726 and Oomycetes (Birch et al. 2006), either function in the interaction space between hyphae and host cells or are
727 transferred into host cells to subvert host immunity. A successful fungal infection with significant disease
728 symptoms is a complicated process that depends on the result of the battle between the pathogen and its host (GS
729 1996). When pathogens start to invade a host, the innate immune system is activated by recognizing microbial
730 invariant molecular patterns (also known as pathogen-associated molecular patterns, PAMPs) (Akira et al. 2006).
731 In fungi, chitin, the important cell wall component, is one of the main PAMPs, which is recognized by pattern-
732 recognition receptors (PRRs) located in the host membrane (Boller 1995), and further activates important
733 chemical pathways and specific gene expressions to eliminate pathogens (Macho and Zipfel 2014). The PAMP-
734 triggered immunity (PTI) is the frontline of the plant host's immune system; if fungi seek to successfully colonize
735 the host, they must avoid inducing PTI or suppress it. Effectors can suppress PTI, but they also can be captured
736 by effector-triggered immunity (ETI). Therefore, linking the disease symptoms and effectors or elucidating their
737 relationships remains a significantly challenging task. We hypothesize that this is why we did not observe a
738 significantly different abundance in the average number of effectors between plant pathogens (the average number
739 = 216) and endophytes (the average number = 207) in our analysis. There is limited capacity to validate the
740 function of effectors in pathogen-host interactions experimentally; accordingly, only a small part of effectors are
741 well studied in model fungi and economically important fungi (Stergiopoulos and Wit 2009), and many effectors
742 have been identified in newly sequenced non-model fungal genomes or not economically important genomes
743 using bioinformatic approaches (Jones et al. 2018). *PgtSRI*, a novel fungal effector identified by Yin et al. (2019)
744 from the wheat rust pathogen *Puccinia graminis*, decreases the abundance of small RNAs by suppressing RNA
745 silencing in plant cells, and further obstructs small RNA-regulated host immune reactions. Czisowski et al. (2021)
746 showed that endophytic *Fusarium oxysporum* strains display different *SIX* gene profiles (a family of effector genes
747 secreted in xylem) with pathogenic strains. However, a larger-scale study of fusarioid fungi did not find a
748 significant difference in the number of effectors (Hill et al. 2022). In this study, we observed that plant pathogens
749 (the average number = 216) include more effectors ($p < 0.05$) than saprotrophs (the average number = 162) and
750 entomopathogens (the average number = 142). The abundance of effectors in endophytes (the average number =
751 207) is significantly higher ($p < 0.05$) than that found in entomopathogens. To explain these differences, we
752 speculate that the pathogenic *F. oxysporum* isolates and non-pathogenic isolates might have similar numbers of
753 effectors that differ in composition. Compared to results from Hill et al. (2022), we include more genomes with
754 lifestyle information ($n = 555$ VS $n = 61$), which provides more numerical information for conducting statistical
755 analysis. Moreover, we believe that our dataset includes more reliable lifestyle information. For instance,
756 entomopathogens are mainly from the families Cordycipitaceae, Clavicipitaceae and Ophiocordycipitaceae;
757 species in these families are more specifically parasitic on insects than species from the family Nectriaceae

758 (fusarioid fungi) (Simmons et al. 2015; Luangsa-ard et al. 2017; Araújo et al. 2018). We confirmed that the
759 abundance of effectors is significantly different between several compared lifestyles, and future extended studies
760 should focus on the composition to verify whether it is a possible indicator for differentiating lifestyles.

761

762 **Basic genomic features are generally consistent with higher taxonomic ranks rather than lifestyles**

763

764 In the genomic era, the rapid development of sequencing technologies and affordable cost of WGS have brought
765 new insights to taxonomy. Genome Taxonomy Database (GTDB) exemplifies the important contribution of
766 genomes in bacterial and archaeal taxonomy (Parks et al. 2018; Rinke et al. 2021). In fungal taxonomy, Gostinčar
767 (2020) first tried to use the genomic distance to delineate fungal species, and obtained a relatively high degree of
768 accuracy in delineating species according to the assumed threshold of genomic distances. However, the proposed
769 criteria have not been widely utilized. Compared with the multilocus phylogenetic taxonomy, huge computational
770 resource requirements, higher sequencing cost, more complicated analytic methods and lower accuracy at higher
771 taxonomic ranks render it useless. In this study, we initially planned to differentiate lifestyles based on the basic
772 numerical features of genomes and exclude the influence of phylogenetic signals. However, we unexpectedly
773 discovered that some basic numerical features, such as genome size, GC content, and gene number, easily accessed
774 from public databases, display powerful resolution for differentiating genomes at the higher levels, especially at
775 the subclass. Inversely, most of these basic genomic features are useless only using the two features tRNA number
776 and genome size without TEs displaying a certain degree of resolving power. To some extent, our discovery agrees
777 with the conclusion of Li et al. (2021), in which fungal genome divergence is broadly consistent with the current
778 taxonomic scheme at higher ranks, even using different genomic information. Fijarczyk et al. (2022) reported that
779 pathogenic fungi include a higher number of protein-coding genes, tRNA genes, and larger genome sizes without
780 repeats than non-pathogenic fungi. Compared with insect-unrelated fungi, they also found that insect-related fungi
781 have smaller genome sizes, gene numbers and exon numbers but increased exon length. In this study, we divided
782 638 genomes into more specific lifestyles instead of only marking them as pathogenic or non-pathogenic, and our
783 results are partially consistent with the previous discoveries by Fijarczyk et al. (2022). More specifically, we
784 observed that plant pathogens have the largest average gene number of 11858, which is significantly larger than
785 the average gene number of saprotrophs (the average number = 10581) and entomopathogens (the average number
786 = 8821). However, entomopathogens have the smallest average gene number, which is significantly smaller than
787 that of endophytes (the average number = 11577). As for genome size and tRNA number, we observed a similar
788 pattern when we compared both features across lifestyles. In aggregate, although several basic genomic features
789 display a certain degree of discrimination for differentiating lifestyles, we prefer to conclude that differences
790 across these basic genomic features reflect taxonomic ranks rather than lifestyles.

791

792 **Functional proteins are useful for differentiating lifestyles**

793

794 Compared with basic genomic features, numerous studies have demonstrated that functional proteins, responsible
795 for degrading substrates, invading host cells and obtaining nutrition are biologically more convincing in
796 differentiating lifestyles (Feldman et al. 2017; Muszewska et al. 2017b; Seong and Krasileva 2023). In the present
797 study, we divided the functional proteins into multiple groups and discovered that these functional proteins
798 generally display relatively high discrimination for differentiating taxonomic groups at different ranks and slightly
799 reduced for distinguishing lifestyles.

800

801 Secretome, a collective term representing all secreted proteins of an organism, is assumed to be related to fungal
802 lifestyles. Krijger et al. (2014) reported that plant pathogens and saprotrophs include larger secretomes than animal
803 pathogens, also indicated that differences in fungal secretome size reflects more on the phylogenetic relationships
804 and less on lifestyle differences. Alfaro et al. (2014) believed that lifestyle is correlated to the composition of the
805 secretome rather than its size. Recently, Chang et al. (2022) reported that the secretome size is mainly determined
806 by phylogeny and lifestyle plays an important auxiliary role. Our results (Table S5: sheet pairwise-lifestyle) reveal
807 that plant pathogens have the largest secretomes (the average number = 847), whereas entomopathogens have the
808 smallest secretomes (the average number = 513). Based on the average number, we can clearly differentiate ($p <$
809 0.05) plant pathogens from entomopathogens, mycoparasites (the average number = 671), saprotrophs (the
810 average number = 667) and entomopathogens, as well as differentiate entophytes (the average number = 828)
811 from entomopathogens. With respect to the main protein groups, including CAZymes, lipases and SSPs, they
812 display similar or higher discrimination than secretome, but lipases display lower discrimination.

813
814 PCWDEs play key roles in obtaining nutrients and degrading the main structural components of the plant cell
815 wall, i.e., cellulose, hemicellulose and pectin. Lichenized fungi live as symbionts of green algae or cyanobacteria,
816 obtaining diverse nutrients from their partners; therefore, they have fewer PCWDEs than non-lichenized fungi
817 (Song et al. 2022). The reduction of PCWDEs is a prevailing trend in ectomycorrhizal Russulaceae (Looney et al.
818 2022), but they retain a certain degree of diversity in components (Kohler et al. 2015). The reduced abundance of
819 PCWDEs in fungi might help in facilitating symbiosis by decreasing the expression of PCWDEs to reduce plant
820 immune responses (Plett and Martin 2011). As for other kinds of lifestyles, the compositions of PCWDEs are
821 different between saprophytic and plant-pathogenic fungi (Zhao et al. 2013; Kubicek et al. 2014). To the best of
822 our knowledge, the present study is the first to conduct a comprehensively comparative analysis on the abundance
823 of PCWDEs across multiple lifestyles. Plant-related fungi including endophytes (the average number = 75), plant
824 pathogens (the average number = 81) and saprotrophs (the average number = 62) have a significantly larger
825 repository of PCWDEs compared with entomopathogens (the average number = 12). For the plant-unrelated fungi,
826 entomopathogens feature the smallest repository of PCWDEs. However, interestingly, human pathogens feature
827 relatively high abundance of PCWDEs (the average number = 72). We investigated the lifestyles of these human
828 pathogens, which belong to *Scedosporium* (Kaur et al. 2019), *Phialemoniopsis* (Alvarez Martinez et al. 2021),
829 *Lomentospora* (Ramirez-Garcia et al. 2018), *Fusarium* (Zhang et al. 2020), *Sporothrix* (Rodrigues et al. 2016),
830 and *Madurella* (Ahmed et al. 2004), also confirmed that these groups are indeed associated with human diseases.
831 However, we did not receive any clues to help explain the high abundance of PCWDEs in human pathogens. We
832 speculate that these species mainly exist as non-human pathogens, but they rarely infect humans as opportunistic
833 pathogens. Therefore, the contraction of PCWDEs has not yet occurred or is in an early evolutionary stage, while
834 still featuring a large number of PCWDEs. More in-depth studies should be carried out to trace changes of
835 PCWDEs in human pathogens.

836
837 FCWDEs are critical for degrading the cell wall of fungal hosts during mycoparasitism. Mycoparasitic species
838 tend to have an expanded repository of FCWDEs (Gruber and Seidl-Seiboth 2012). Our results show that
839 mycoparasites have the largest repository of FCWDEs (the average number = 41), which is significantly larger
840 than entomopathogens (the average number = 28), human pathogens (the average number = 24), and plant
841 pathogens (the average number = 27). To date, there are few studies that investigate the relationship between
842 FCWDEs and fungal lifestyles. Results in the present study represent an important addition to this field.

843
844 **The promising but limited potential of machine learning for lifestyle prediction**

845

846 Machine learning algorithms heavily rely on massive amounts of data, the accuracy of which dramatically depends
847 on not only the correctness of the training data and test data but also the quantity of input data (Raudys and Jain
848 1991; Sordo and Zeng 2005; Read et al. 2011). In classification tasks, inaccurately labeled datasets and inadequate
849 sampling can lead to incorrect predictions. In the present study, there were two main challenges: inadequate
850 sampling in several lifestyles and inaccurate lifestyle labels for some genomes. Unbalanced lifestyle distribution
851 of genomes from public databases is common and unavoidable. Distribution largely depends on economic and
852 medical importance, as well as the availability of samples. In our dataset, we include enough genomes of plant
853 pathogens (n = 372), but fewer genomes of mycoparasites (n = 23), human pathogens (n = 16), and nematophagous
854 fungi (n = 4). We excluded nematophagous fungi during analysis, but the relatively small sample sizes for multiple
855 lifestyles affected the predictive accuracies to some extent, as shown in Fig. 7. Another challenge is assigning
856 lifestyle labels to each genome. We attempted to determine the lifestyle of each genome, but for most genomes,
857 the lifestyle is determined based on published literature or the submitter's description. Moreover, most studies
858 directly characterize fungi isolated from diseased plants as plant pathogens, which does not follow Koch's
859 postulates (van Wyk et al. 2012; Oberti et al. 2020; Telenko et al. 2020). Our predictive models display a high
860 degree of accuracy in differentiating plant pathogens from other lifestyles, and adequate sampling reduced the
861 error caused by inaccurate labeling. In predicting the lifestyle of unlabeled genomes, we further compared the
862 predicted lifestyles and observed lifestyles in phylogenetically closed groups, and most of our predicted lifestyles
863 are consistent with the observed lifestyles. Taken together, we suggest that using machine learning algorithms to
864 predict fungal lifestyles is promising and can be improved with more sequenced genomes in the future.

865

866 **Predicting potentially adverse fungal lifestyle**

867

868 Fungi provide food and important medical and industrial secondary metabolites, as well as promote the global
869 carbon cycle (Hyde et al. 2019; Lücking et al. 2021; Maharachchikumbura et al. 2021). However, the past two
870 decades have witnessed the occurrence of new and emerging disease-causing fungi that infect plants, animals and
871 humans (Fisher et al. 2012). Human activities have largely expanded fungal distribution and brought pathogenic
872 fungal species accidentally to new ecosystems (Santini et al. 2013). *Pseudogymnoascus destructans*, an emerging
873 fungal pathogen causing white-nose syndrome in bats, was initially detected in a commercial tourist cave, and it
874 was speculated that the species was brought to external environments by tourist movements and further spread
875 across North America, resulting in widespread mortality of hibernating bats (Blehert et al. 2009; Frick et al. 2015;
876 Langwig et al. 2016). During the long-term interaction between fungal pathogens and hosts, both the fungi and
877 the host have developed mechanisms to counteract each other's actions. Therefore, the hosts do not develop disease
878 symptoms even if the fungi express abundant virulent factors. However, the fungi are introduced to new habitats
879 and colonize new hosts, disease-causing interactions do develop (Parker and Gilbert 2004). *Phytophthora*
880 *ramorum*, an alien plant pathogen to California and Oregon, causes a disease known as sudden oak death, that led
881 to the death of a large number of trees, seriously threatening the local forest ecosystem (Rizzo and Garbelotto
882 2003). In addition, some fungal species or strains have multiple lifestyles, including non-pathogenic and
883 pathogenic. Cannon et al. (2012) and Liu et al. (2022) demonstrated that endophytic fungi can switch to pathogenic
884 lifestyle and cause disease symptoms. Due to the lack of effective analytical methods, some potential fungal
885 pathogens were neglected until they caused devastating impacts on human health, food security and ecosystem
886 stability (Anderson et al. 2004; Fisher et al. 2012; McDonald and Stukenbrock 2016). In scientific investigations
887 and daily practices, we only observe one specific lifestyle of a certain fungal isolate under the current condition.
888 Therefore, the experimentally exploring the potential lifestyles is impractical. In the study, our machine learning

889 model determine the fungal lifestyles according the corresponding probabilities, the highest probability represents
890 the final predictive results, and the secondary high but non-zero probabilities imply that the strain might have
891 other kind of lifestyles. For instance, *Arthrimum puccinioides* CBS 549.89 was predicted as a plant pathogen with
892 a probability of 0.6689, but it may also be an endophyte or saprotroph with a probability of 0.1683 and 0.1628
893 respectively. Through a literature survey, we did observe endophytic and saprotrophic lifestyles in other species
894 within the genus *Arthrimum* (Wang et al. 2018). With more fungal genomes sequenced and added to the dataset,
895 the accuracy of our predictive model for determining fungal lifestyles using machine learning algorithms will
896 become more reliable. The relatively high probability of harmful lifestyles can be used as an early warning of
897 some devastating fungi. By identifying these harmful fungi early on, appropriate measures can be taken to prevent
898 their spread and minimize their impact.

899

900 **Supplementary information**

901

902 Fig. S1. Distribution and proportion (%) of TE families in 638 genome assemblies. The bubble size represents the
903 proportion of the TE in the genome. The bar represents the proportion of total TE size to the genome size.

904 Fig. S2. Composition and abundance of functional protein groups in 638 genome assemblies. The bubble size
905 represents the number of the protein group. The bar represents the proportion of the secretome size to the total
906 number of proteins per genome.

907 Fig. S3. Predictive accuracies of six machine learning algorithms based on three data matrices.

908

909 Table S1. A summary table containing genome information of 638 genome assemblies, lineage information, and
910 statistics of TE categories, basic genomic features and functional protein groups.

911 Table S2. Taxonomic and lifestyle coverage of 638 Sordariomycete genomes.

912 Table S3. Number and proportion of different sequencing technologies.

913 Table S4. Results of Pearson Correlation of 25 basic genomic features.

914 Table S5. Comparative analysis results of 25 basic genomic features.

915 Table S6. Results of Pearson Correlation of 24 functional protein features.

916 Table S7. Comparative analysis results of 24 functional protein features.

917 Table S8. Predictive accuracies of six machine learning algorithms

918 Table S9. Prediction results of 83 undetermined genomes and the observed lifestyles of phylogenetically closed
919 groups.

920

921 **Code availability**

922 All the scripts used for statistics, visualization and machine learning are written in R or Python. Scripts are
923 available at GitHub (<https://github.com/ypchan/Predict-fungal-lifestyles>).

924

925 **Acknowledgments**

926

927 This research was funded by the Talent Introduction and Cultivation Project, University of Electronic Science and
928 Technology of China, grant number A1098531023601245. Several genomes were produced by the US Department
929 of Energy Joint Genome Institute (JGI) (<https://ror.org/04xm1d337>; operated under Contract No. DE-AC02-
930 05CH11231) in collaboration with the user community and we acknowledge Professor J.W. Spatafora for allowing
931 us to use these genomes submitted in JGI. K.D. Hyde acknowledges the National Research Council of Thailand
932 (NRCT) grant “Total fungal diversity in a given forest area with implications towards species numbers, chemical

933 diversity and biotechnology” (grant no. N42A650547).

934

935 **Author contributions**

936

937 SSNM and YPC designed the study. YPC collected genome data and performed all bioinformatic analyses. SSNM,
938 PWS, WHT and YPC checked the tables and performed lifestyle assignments. YPC and SSNM wrote the first
939 draft of the manuscript. RX checked all R codes and Python codes, and provided a portion of computing resources.
940 SSNM, HKD and MS help in revision. All authors provided valuable comments on the manuscript. All authors
941 read and approved the final manuscript.

942

943 **Funding**

944

945 This research was funded by Talent Introduction and Cultivation Project, University of Electronic Science and
946 Technology of China, grant number A1098531023601245.

947

948 **Declarations**

949

950 The authors declare that there is no conflict of interest related to this study.

951

952 **Reference**

953

954 Ahmed AOA, van Leeuwen W, Fahal A, van de Sande W, Verbrugh H, van Belkum A (2004) Mycetoma caused
 955 by *Madurella mycetomatis*: a neglected infectious burden. *Lancet Infect Dis* 4 (9):566-574.
 956 [https://doi.org/10.1016/S1473-3099\(04\)01131-4](https://doi.org/10.1016/S1473-3099(04)01131-4)

957 Akira S, Uematsu S, Takeuchi O (2006) Pathogen recognition and innate immunity. *Cell* 124 (4):783-801.
 958 <https://doi.org/10.1016/j.cell.2006.02.015>

959 Alfaro M, Oguiza JA, Ramírez L, Pisabarro AG (2014) Comparative analysis of secretomes in basidiomycete
 960 fungi. *J Proteomics* 102:28-43. <https://doi.org/10.1016/j.jprot.2014.03.001>

961 Alvarez Martinez D, Alberto C, Riat A, Schuhler C, Valladares P, Ninet B, Kraak B, Crous PW, Hou LW, Toutous
 962 Trelu L (2021) *Phialemoniopsis limonesiae* sp. nov. causing cutaneous phaeohyphomycosis in an
 963 immunosuppressed woman. *Emerging Microbes Infect* 10 (1):400-406.
 964 <https://doi.org/10.1080/22221751.2021.1892458>

965 Anderson PK, Cunningham AA, Patel NG, Morales FJ, Epstein PR, Daszak P (2004) Emerging infectious diseases
 966 of plants: pathogen pollution, climate change and agrotechnology drivers. *Trends Ecol Evol* 19 (10):535-
 967 544. <https://doi.org/10.1016/j.tree.2004.07.021>

968 Araújo JPM, Evans HC, Kepler R, Hughes DP (2018) Zombie-ant fungi across continents: 15 new species and
 969 new combinations within *Ophiocordyceps*. I. Myrmecophilous hirsutelloid species. *Stud Mycol*
 970 90:119-160. <https://doi.org/10.1016/j.simyco.2017.12.002>

971 Bao W, Kojima KK, Kohany O (2015) Repbase update, a database of repetitive elements in eukaryotic genomes.
 972 *Mobile DNA* 6 (1):11. <https://doi.org/10.1186/s13100-015-0041-9>

973 Barros MBdL, Paes RdA, Schubach AO (2011) *Sporothrix schenckii* and Sporotrichosis. *Clin Microbiol Rev* 24
 974 (4):633-654. <https://doi.org/10.1128/CMR.00007-11>

975 Bartlett P, Eberhardt U, Schütz N, Beker HJ (2022) Species determination using AI machine-learning algorithms:
 976 *Hebeloma* as a case study. *IMA Fungus* 13 (1):13. <https://doi.org/10.1186/s43008-022-00099-x>

977 Birch PRJ, Rehmany AP, Pritchard L, Kamoun S, Beynon JL (2006) Trafficking arms: oomycete effectors enter
 978 host plant cells. *Trends Microbiol* 14 (1):8-11. <https://doi.org/10.1016/j.tim.2005.11.007>

979 Blehert DS, Hicks AC, Behr M, Meteyer CU, Berlowski-Zier BM, Buckles EL, Coleman JTH, Darling SR, Gargas
 980 A, Niver R, Okoniewski JC, Rudd RJ, Stone WB (2009) Bat white-nose syndrome: an emerging
 981 fungal pathogen? *Science* 323 (5911), 227-227. <https://doi.org/10.1126/science.1163874>

982 Boddy L (2016) Chapter 9 - Interactions with humans and other animals. In: Watkinson SC, Boddy L, Money NP
 983 (eds) *The Fungi* (Third Edition). Academic Press, Boston, pp 293-336. <https://doi.org/10.1016/B978-0-12-382034-1.00009-8>

985 Boller T (1995) Chemoperception of microbial signals in plant cells. *Annu Rev Plant Phys* 46 (1):189-214.
 986 <https://doi.org/10.1146/annurev.pp.46.060195.001201>

987 Brûna T, Hoff KJ, Lomsadze A, Stanke M, Borodovsky M (2021) BRAKER2: automatic eukaryotic genome
 988 annotation with GeneMark-EP+ and AUGUSTUS supported by a protein database. *NAR Genomics*
 989 *Bioinf* 3 (1). <https://doi.org/10.1093/nargab/lqaa108>

990 Brûna T, Lomsadze A, Borodovsky M (2020) GeneMark-EP+: eukaryotic gene prediction with self-training in the
 991 space of genes and proteins. *NAR Genomics Bioinf* 2 (2). <https://doi.org/10.1093/nargab/lqaa026>

992 Bzdok D, Krzywinski M, Altman N (2018) Machine learning: supervised methods. *Nat Methods* 15 (1):5-6.
 993 <https://doi.org/10.1038/nmeth.4551>

994 Camacho DM, Collins KM, Powers RK, Costello JC, Collins JJ (2018) Next-generation machine learning for
 995 biological networks. *Cell* 173 (7):1581-1592. <https://doi.org/10.1016/j.cell.2018.05.015>

- 996 Cannon PF, Damm U, Johnston PR, Weir BS (2012) *Colletotrichum*: current status and future directions. *Stud.*
997 *Mycol* 73 (1):181-213. <https://doi.org/10.3114/sim0014>
- 998 Cantalapiedra CP, Hernández-Plaza A, Letunic I, Bork P, Huerta-Cepas J (2021) eggNOG-mapper v2: functional
999 annotation, orthology assignments, and domain prediction at the metagenomic scale. *Mol Biol Evol* 38
1000 (12):5825-5829. <https://doi.org/10.1093/molbev/msab293>
- 1001 Chan Patricia P, Lin Brian Y, Mak Allysia J, Lowe Todd M (2021) tRNAscan-SE 2.0: improved detection and
1002 functional classification of transfer RNA genes. *Nucleic Acids Res* 49 (16):9077-9096.
1003 <https://doi.org/10.1093/nar/gkab688>
- 1004 Chang TH, Hassan O, Lee YS (2018) First Report of Anthracnose of Japanese Plum (*Prunus salicina*) Caused by
1005 *Colletotrichum nymphaeae* in Korea. *Plant Disease* 102 (7):1461-1461. <https://doi.org/10.1094/pdis-01-18-0018-pdn>
- 1007 Chang Y, Wang Y, Mondo S, Ahrendt S et al (2022) Evolution of zygomycete secretomes and the origins of
1008 terrestrial fungal ecologies. *iScience* 25 (8). <https://doi.org/10.1016/j.isci.2022.104840>
- 1009 Chen YP, Wu T, Tian WH, Ilyukhin F, Hyde KD, Maharachchikumbura SSN (2022) Comparative genomics
1010 provides new insights into the evolution of *Colletotrichum*. *Mycosphere* 13 (2):56.
1011 <https://doi.org/10.5943/mycosphere/si/1f/5>
- 1012 Consortium TU (2020) UniProt: the universal protein knowledgebase in 2021. *Nucleic Acids Res* 49 (D1):D480-
1013 D489. <https://doi.org/10.1093/nar/gkaa1100>
- 1014 Cooper, G M. (2000). *The Cell: A Molecular Approach*. 2nd edition. Sinauer Associates.
- 1015 Crawford K, Heatley NG, Boyd PF, Hale CW, Kelly BK, Miller GA, Smith N (1952) Antibiotic production by a
1016 species of *Cephalosporium*. *J Gen Microbiol* 6 (1-2):47-59. <https://doi.org/10.1099/00221287-6-1-2-47>
- 1017 Crous PW, Lombard L, Sandoval-Denis M, Seifert KA et al (2021) *Fusarium*: more than a node or a foot-shaped
1018 basal cell. *Stud Mycol* 98:100116. <https://doi.org/10.1016/j.simyco.2021.100116>
- 1019 Cziślowski E, Zeil-Rolfe I, Aitken EAB (2021) Effector profiles of endophytic *Fusarium* associated with
1020 asymptomatic banana (*Musa* sp.) hosts. *Int J Mol Sci* 22 (5):2508. <https://doi.org/10.3390/ijms22052508>
- 1021 Dasari P, Shopova IA, Stroe M, Wartenberg D, Martin-Dahse H, Beyersdorf N, Hortschansky P, Dietrich S,
1022 Cseresnyés Z, Figge MT, Westermann M, Skerka C, Brakhage AA, Zipfel PF (2018) AspF2 from
1023 *Aspergillus fumigatus* recruits human immune regulators for immune evasion and cell damage. *Front*
1024 *Immunol* 9. <https://doi.org/10.3389/fimmu.2018.01635>
- 1025 de Jonge R, Bolton MD, Thomma BPHJ (2011) How filamentous pathogens co-opt plants: the ins and outs of
1026 fungal effectors. *Curr Opin Plant Biol* 14 (4):400-406. <https://doi.org/10.1016/j.pbi.2011.03.005>
- 1027 Dean R, Van Kan JL, Pretorius ZA, Hammond-kosack KE, Dipietro A, Spanu PD, Rudd JJ, Dickman M, Kahmann
1028 R, Ellis J, Foster GD (2012) The Top 10 fungal pathogens in molecular plant pathology. *Mol Plant Pathol*
1029 13 (4):414-430. <https://doi.org/10.1111/j.1364-3703.2011.00783.x>
- 1030 Delaye L, García-Guzmán G, Heil M (2013) Endophytes versus biotrophic and necrotrophic pathogens—are
1031 fungal lifestyles evolutionarily stable traits? *Fungal Divers* 60 (1):125-135.
1032 <https://doi.org/10.1007/s13225-013-0240-y>
- 1033 Deo RC (2015) Machine learning in medicine. *Circulation* 132 (20):1920-1930.
1034 <https://doi.org/10.1161/CIRCULATIONAHA.115.001593>
- 1035 Derkarabetian S, Castillo S, Koo PK, Ovchinnikov S, Hedin M (2019) A demonstration of unsupervised machine
1036 learning in species delimitation. *Mol Phylogenet Evol* 139:106562.
1037 <https://doi.org/10.1016/j.ympev.2019.106562>
- 1038 Eastwood DC, Floudas D, Binder M, Majcherczyk A et al (2011) The plant cell wall-decomposing machinery
1039 underlies the functional diversity of forest fungi. *Science* 333 (6043):762-765.

1040 <https://doi.org/10.1126/science.1205411>

1041 Emanuelsson O, Brunak S, von Heijne G, Nielsen H (2007) Locating proteins in the cell using TargetP, SignalP
1042 and related tools. *Nat Protoc* 2 (4):953-971. <https://doi.org/10.1038/nprot.2007.131>

1043 Eriksson OE WK (1997) Supraordinal taxa of Ascomycota. *Myconet* 1 (1):1-16

1044 Feldman D, Kowbel DJ, Glass NL, Yarden O, Hadar Y (2017) A role for small secreted proteins (SSPs) in a
1045 saprophytic fungal lifestyle: ligninolytic enzyme regulation in *Pleurotus ostreatus*. *Sci Rep* 7 (1):14553.
1046 <https://doi.org/10.1038/s41598-017-15112-2>

1047 Fijarczyk A, Hesseuauer P, Hamelin RC, Landry CR (2022) Lifestyles shape genome size and gene content in
1048 fungal pathogens. *bioRxiv:2022.2008.2024.505148*. <https://doi.org/10.1101/2022.08.24.505148>

1049 Fisher MC, Henk DA, Briggs CJ, Brownstein JS, Madoff LC, McCraw SL, Gurr SJ (2012) Emerging fungal
1050 threats to animal, plant and ecosystem health. *Nature* 484 (7393), 186-194.
1051 <https://doi.org/10.1038/nature10947>

1052 Fontana DC, de Paula S, Torres AG, de Souza VHM, Pascholati SF, Schmidt D, Dourado Neto D (2021)
1053 Endophytic fungi: biological control and induced resistance to phytopathogens and abiotic stresses.
1054 *Pathogens* 10 (5):570. <https://doi.org/10.3390/pathogens10050570>

1055 Fouché S, Badet T, Oggenfuss U, Plissonneau C, Francisco CS, Croll D (2019) Stress-driven transposable element
1056 de-repression dynamics and virulence evolution in a fungal pathogen. *Mol Biol Evol* 37 (1):221-239.
1057 <https://doi.org/10.1093/molbev/msz216>

1058 Frey-Klett P, Burlinson P, Deveau A, Barret M, Tarkka M, Sarniguet A (2011) Bacterial-fungal interactions:
1059 hyphens between agricultural, clinical, environmental, and food microbiologists. *Microbiol Mol Biol*
1060 *Rev* 75 (4):583-609. <https://doi.org/10.1128/MMBR.00020-11>

1061 Frick WF, Puechmaille SJ, Hoyt JR, Nickel BA, Langwig KE, Foster JT, Barlow KE, Bartonička T, Feller D,
1062 Haarsma A-J, Herzog C, Horáček I, van der Kooij J, Mulkens B, Petrov B, Reynolds R, Rodrigues L,
1063 Stihler CW, Turner GG, Kilpatrick AM (2015) Disease alters macroecological patterns of North
1064 American bats. *Glob Ecol Biogeogr* 24 (7), 741-749. <https://doi.org/10.1111/geb.12290>

1065 Friesen TL, Stukenbrock EH, Liu Z, Meinhardt S, Ling H, Faris JD, Rasmussen JB, Solomon PS, McDonald BA,
1066 Oliver RP (2006) Emergence of a new disease as a result of interspecific virulence gene transfer. *Nat*
1067 *Genet* 38 (8):953-956. <https://doi.org/10.1038/ng1839>

1068 Fu L, Niu B, Zhu Z, Wu S, Li W (2012) CD-HIT: accelerated for clustering the next-generation sequencing data.
1069 *Bioinformatics* 28 (23):3150-3152. <https://doi.org/10.1093/bioinformatics/bts565>

1070 Gíslason MH, Nielsen H, Almagro Armenteros JJ, Johansen AR (2021) Prediction of GPI-anchored proteins with
1071 pointer neural networks. *Curr Res Biotechnol* 3:6-13. <https://doi.org/10.1016/j.crbiot.2021.01.001>

1072 Glazebrook J (2005) Contrasting mechanisms of defense against biotrophic and necrotrophic pathogens. *Annu*
1073 *Rev Phytopathol* 43 (1):205-227. <https://doi.org/10.1146/annurev.phyto.43.040204.135923>

1074 Goodwin S, McPherson JD, McCombie WR (2016) Coming of age: ten years of next-generation sequencing
1075 technologies. *Nat Rev Genet* 17 (6):333-351. <https://doi.org/10.1038/nrg.2016.49>

1076 Gostinčar C (2020) Towards genomic criteria for delineating fungal species. *J Fungi* 6 (4):246.
1077 <https://doi.org/10.3390/jof6040246>

1078 Grigoriev IV, Nikitin R, Haridas S, Kuo A, Ohm R, Otilar R, Riley R, Salamov A, Zhao X, Korzeniewski F,
1079 Smirnova T, Nordberg H, Dubchak I, Shabalov I (2013) MycoCosm portal: gearing up for 1000 fungal
1080 genomes. *Nucleic Acids Res* 42 (D1):D699-D704. <https://doi.org/10.1093/nar/gkt1183>

1081 Gruber S, Seidl-Seiboth V (2012) Self versus non-self: fungal cell wall degradation in *Trichoderma*. *Microbiology*
1082 158 (1):26-34. <https://doi.org/10.1099/mic.0.052613-0>

1083 GS K (1996) Disease mechanisms of fungi. In: Baron S (ed) *Medical Microbiology*. vol 4th edition. University of

1084 Texas Medical Branch at Galveston, Galveston (TX).

1085 Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O (2010) New algorithms and methods to
1086 estimate maximum-likelihood phylogenies: assessing the performance of PhyML 3.0. *Syst Biol* 59
1087 (3):307-321. <https://doi.org/10.1093/sysbio/syq010>

1088 Han S, Wang M, Ma Z, Raza M, Zhao P, Liang J, Gao M, Li Y, Wang J, Hu D (2023) *Fusarium* diversity associated
1089 with diseased cereals in China, with an updated phylogenomic assessment of the genus. *Stud Mycol*
1090 104:87-148. <https://doi.org/10.3114/sim.2022.104.02>

1091 Haridas S, Albert R, Binder M, Bloem J et al (2020) 101 Dothideomycetes genomes: a test case for predicting
1092 lifestyles and emergence of pathogens. *Stud Mycol* 96:141-153.
1093 <https://doi.org/10.1016/j.simyco.2020.01.003>

1094 Hill R, Buggs RJA, Vu DT, Gaya E (2022) Lifestyle transitions in fusarioid fungi are frequent and lack clear
1095 genomic signatures. *Mol Biol Evol* 39 (4). <https://doi.org/10.1093/molbev/msac085>

1096 Hoang DT, Chernomor O, von Haeseler A, Minh BQ, Vinh LS (2017) UFBoot2: improving the ultrafast bootstrap
1097 approximation. *Mol Biol Evol* 35 (2):518-522. <https://doi.org/10.1093/molbev/msx281>

1098 Hobbie EA, Horton TR (2007) Evidence that saprotrophic fungi mobilise carbon and mycorrhizal fungi mobilise
1099 nitrogen during litter decomposition. *New Phytol* 173 (3):447-449. <https://doi.org/10.1111/j.1469-8137.2007.01984.x>

1100

1101 Horton P, Park K-J, Obayashi T, Fujita N, Harada H, Adams-Collier CJ, Nakai K (2007) WoLF PSORT: protein
1102 localization predictor. *Nucleic Acids Res* 35 (suppl_2):W585-W587. <https://doi.org/10.1093/nar/gkm259>

1103 Hu E-Z, Lan X-R, Liu Z-L, Gao J, Niu D-K (2022) A positive correlation between GC content and growth
1104 temperature in prokaryotes. *BMC Genomics* 23 (1):110. <https://doi.org/10.1186/s12864-022-08353-7>

1105 Hubley R, Finn RD, Clements J, Eddy SR, Jones TA, Bao W, Smit AFA, Wheeler TJ (2015) The Dfam database
1106 of repetitive DNA families. *Nucleic Acids Res* 44 (D1):D81-D89. <https://doi.org/10.1093/nar/gkv1272>

1107 Hyde KD, Norphanphoun C, Maharachchikumbura SSN, Bhat DJ et al (2020) Refined families of
1108 Sordariomycetes. *Mycosphere* 11 (1):305-1059. <https://doi.org/10.5943/mycosphere/11/1/7>

1109 Hyde KD, Xu J, Rapior S, Jeewon R et al (2019) The amazing potential of fungi: 50 ways we can exploit fungi
1110 industrially. *Fungal Divers* 97 (1), 1-136. <https://doi.org/10.1007/s13225-019-00430-9>

1111 Ismaiel AA, Papenbrock J (2015) Mycotoxins: producing fungi and mechanisms of phytotoxicity. *Agriculture* 5
1112 (3):492-537. <https://doi.org/10.3390/agriculture5030492>

1113 Jenks JD, Reed SL, Seidel D, Koehler P, Cornely OA, Mehta SR, Hoenigl M (2018) Rare mould infections caused
1114 by Mucorales, *Lomentospora prolificans* and *Fusarium*, in San Diego, CA: the role of antifungal
1115 combination therapy. *Int J Antimicrob Agents* 52 (5):706-712.
1116 <https://doi.org/10.1016/j.ijantimicag.2018.08.005>

1117 Jia M, Chen L, Xin H-L, Zheng C-J, Rahman K, Han T, Qin L-P (2016) A friendly relationship between endophytic
1118 fungi and medicinal plants: a systematic review. *Front Microbiol* 7.
1119 <https://doi.org/10.3389/fmicb.2016.00906>

1120 Jones DAB, Bertazzoni S, Turo CJ, Syme RA, Hane JK (2018) Bioinformatic prediction of plant-pathogenicity
1121 effector proteins of fungi. *Curr Opin Microbiol* 46:43-49. <https://doi.org/10.1016/j.mib.2018.01.017>

1122 Kaewchai S, Soyong K, Hyde KD (2009) Mycofungicides and fungal biofertilizers. *Fungal Divers* 38:25-50

1123 Kale SD, Tyler BM (2011) Entry of oomycete and fungal effectors into plant and animal host cells. *Cell Microbiol*
1124 13 (12):1839-1848. <https://doi.org/10.1111/j.1462-5822.2011.01659.x>

1125 Kalyanamoorthy S, Minh BQ, Wong TKF, von Haeseler A, Jermiin LS (2017) ModelFinder: fast model selection
1126 for accurate phylogenetic estimates. *Nat Methods* 14 (6):587-589. <https://doi.org/10.1038/nmeth.4285>

1127 Katoh K, Misawa K, Kuma Ki, Miyata T (2002) MAFFT: a novel method for rapid multiple sequence alignment

1128 based on fast Fourier transform. *Nucleic Acids Res* 30 (14):3059-3066.
1129 <https://doi.org/10.1093/nar/gkf436>

1130 Kaur J, Kautto L, Penesyana A, Meyer W, Elbourne LDH, Paulsen IT, Nevalainen H (2019) Interactions of an
1131 emerging fungal pathogen *Scedosporium aurantiacum* with human lung epithelial cells. *Sci Rep* 9
1132 (1):5035. <https://doi.org/10.1038/s41598-019-41435-3>

1133 Kim K-T, Jeon J, Choi J, Cheong K, Song H, Choi G, Kang S, Lee Y-H (2016) Kingdom-wide analysis of fungal
1134 small secreted proteins (SSPs) reveals their potential role in host association. *Front Plant Sci* 7.
1135 <https://doi.org/10.3389/fpls.2016.00186>

1136 Kirkland TN, Muszewska A, Stajich JE (2018) Analysis of transposable elements in *Coccidioides* species. *J Fungi*
1137 4 (1):13. <https://doi.org/10.3390/jof4010013>

1138 Knapp DG, Németh JB, Barry K, Hainaut M, Henrissat B, Johnson J, Kuo A, Lim JHP, Lipzen A, Nolan M, Ohm
1139 RA, Tamás L, Grigoriev IV, Spatafora JW, Nagy LG, Kovács GM (2018) Comparative genomics
1140 provides insights into the lifestyle and reveals functional heterogeneity of dark septate endophytic fungi.
1141 *Sci Rep* 8 (1):6321. <https://doi.org/10.1038/s41598-018-24686-4>

1142 Kohler A, Kuo A, Nagy LG, Morin E et al (2015) Convergent losses of decay mechanisms and rapid turnover of
1143 symbiosis genes in mycorrhizal mutualists. *Nat Genet* 47 (4):410-415. <https://doi.org/10.1038/ng.3223>

1144 Krijger J-J, Thon MR, Deising HB, Wiersma SGR (2014) Compositions of fungal secretomes indicate a greater
1145 impact of phylogenetic history than lifestyle adaptation. *BMC Genomics* 15 (1):722.
1146 <https://doi.org/10.1186/1471-2164-15-722>

1147 Kubicek CP, Starr TL, Glass NL (2014) Plant cell wall-degrading enzymes and their secretion in plant-pathogenic
1148 fungi. *Annu Rev Phytopathol* 52 (1):427-451. <https://doi.org/10.1146/annurev-phyto-102313-045831>

1149 Kwon SL, Park MS, Jang S, Lee YM, Heo YM, Hong J-H, Lee H, Jang Y, Park J-H, Kim C, Kim G-H, Lim YW,
1150 Kim J-J (2021) The genus *Arthrimum* (Ascomycota, Sordariomycetes, Apiosporaceae) from marine
1151 habitats from Korea, with eight new species. *IMA Fungus* 12 (1):13. <https://doi.org/10.1186/s43008-021-00065-z>

1152

1153 Lang D, Zhang S, Ren P, Liang F, Sun Z, Meng G, Tan Y, Li X, Lai Q, Han L, Wang D, Hu F, Wang W, Liu S
1154 (2020) Comparison of the two up-to-date sequencing technologies for genome assembly: HiFi reads of
1155 Pacific Biosciences Sequel II system and ultralong reads of Oxford Nanopore. *GigaScience* 9 (12).
1156 <https://doi.org/10.1093/gigascience/giaa123>

1157 Langwig KE, Frick WF, Hoyt JR, Parise KL, Drees KP, Kunz TH, Foster JT, Kilpatrick AM (2016) Drivers of
1158 variation in species impacts for a multi-host fungal disease of bats. *Philos Trans R Soc B* 371 (1709):
1159 20150456. <https://doi.org/10.1098/rstb.2015.0456>

1160 Lebreton A, Tang N, Kuo A, LaButti K, Andreopoulos W, Drula E, Miyauchi S, Barry K, Clum A, Lipzen A,
1161 Mousain D, Ng V, Wang R, Dai Y, Henrissat B, Grigoriev IV, Guerin-Laguette A, Yu F, Martin FM (2022)
1162 Comparative genomics reveals a dynamic genome evolution in the ectomycorrhizal milk-cap (*Lactarius*)
1163 mushrooms. *New Phytol* 235 (1):306-319. <https://doi.org/10.1111/nph.18143>

1164 Lee S-J, Rose JKC (2010) Mediation of the transition from biotrophy to necrotrophy in hemibiotrophic plant
1165 pathogens by secreted effector proteins. *Plant Signaling Behav* 5 (6):769-772.
1166 <https://doi.org/10.4161/psb.5.6.11778>

1167 Li J, Cornelissen B, Rep M (2020) Host-specificity factors in plant pathogenic fungi. *Fungal Genet Biol*
1168 144:103447. <https://doi.org/10.1016/j.fgb.2020.103447>

1169 Li Y, Steenwyk JL, Chang Y, Wang Y, James TY, Stajich JE, Spatafora JW, Groenewald M, Dunn CW, Hittinger
1170 CT, Shen X-X, Rokas A (2021) A genome-scale phylogeny of the kingdom fungi. *Curr Biol* 31 (8):1653-
1171 1665.e1655. <https://doi.org/10.1016/j.cub.2021.01.074>

- 1172 Liu F, Ma ZY, Hou LW, Diao YZ, Wu WP, Damm U, Song S, Cai L (2022) Updating species diversity of
 1173 *Colletotrichum*, with a phylogenomic overview. *Stud Mycol* 101 (1):1-56.
 1174 <https://doi.org/10.3114/sim.2022.101.01>
- 1175 Looney B, Miyauchi S, Morin E, Drula E, Courty PE, Kohler A, Kuo A, LaButti K, Pangilinan J, Lipzen A, Riley
 1176 R, Andreopoulos W, He G, Johnson J, Nolan M, Tritt A, Barry KW, Grigoriev IV, Nagy LG, Hibbett D,
 1177 Henrissat B, Matheny PB, Labbé J, Martin FM (2022) Evolutionary transition to the ectomycorrhizal
 1178 habit in the genomes of a hyperdiverse lineage of mushroom-forming fungi. *New Phytol* 233 (5):2294-
 1179 2309. <https://doi.org/10.1111/nph.17892>
- 1180 Lorrain C, Feurtey A, Möller M, Haueisen J, Stukenbrock E (2021) Dynamics of transposable elements in recently
 1181 diverged fungal pathogens: lineage-specific transposable element content and efficiency of genome
 1182 defenses. *G3-Genes Genom Genet* 11 (4). <https://doi.org/10.1093/g3journal/jkab068>
- 1183 Lu S, Edwards MC (2016) Genome-wide analysis of small secreted cysteine-rich proteins identifies candidate
 1184 effector proteins potentially involved in *Fusarium graminearum*–wheat interactions. *Phytopathology*
 1185 106 (2):166-176. <https://doi.org/10.1094/phyto-09-15-0215-r>
- 1186 Luangsa-ard JJ, Mongkolsamrit S, Thanakitpipattana D, Khonsanit A, Tasanathai K, Noisripoom W, Humber RA
 1187 (2017) Clavicipitaceous entomopathogens: new species in *Metarhizium* and a new genus *Nigelia*. *Mycol*
 1188 *Prog* 16 (4):369-391. <https://doi.org/10.1007/s11557-017-1277-1>
- 1189 Lücking R, Aime MC, Robbertse B, Miller AN, Aoki T, Ariyawansa HA, Cardinali G, Crous PW, Druzhinina IS,
 1190 Geiser DM, Hawksworth DL, Hyde KD, Irinyi L, Jeewon R, Johnston PR, Kirk PM, Malosso E, May
 1191 TW, Meyer W, Nilsson HR, Öpik M, Robert V, Stadler M, Thines M, Vu D, Yurkov AM, Zhang N,
 1192 Schoch CL (2021) Fungal taxonomy and sequence-based nomenclature. *Nat Microbiol* 6 (5), 540-548.
 1193 <https://doi.org/10.1038/s41564-021-00888-x>
- 1194 Luo Z-L, Hyde KD, Liu J-K, Maharachchikumbura SSN, Jeewon R, Bao D-F, Bhat DJ, Lin C-G, Li W-L, Yang
 1195 J, Liu N-G, Lu Y-Z, Jayawardena RS, Li J-F, Su H-Y (2019) Freshwater Sordariomycetes. *Fungal Divers*
 1196 99 (1):451-660. <https://doi.org/10.1007/s13225-019-00438-1>
- 1197 Ma C, Zhang HH, Wang X (2014) Machine learning for Big Data analytics in plants. *Trends Plant Sci* 19 (12):798-
 1198 808. <https://doi.org/10.1016/j.tplants.2014.08.004>
- 1199 Macho Alberto P, Zipfel C (2014) Plant PRRs and the activation of innate immune signaling. *Mol Cell* 54 (2):263-
 1200 272. <https://doi.org/10.1016/j.molcel.2014.03.028>
- 1201 Magyar D, Tartally A, Merényi Z (2022) *Hagnosa longicapillata*, gen. nov., sp. nov., a new sordariaceous
 1202 Ascomycete in the indoor environment, and the proposal of Hagnosaceae fam. nov. *Pathogens* 11 (5):593.
 1203 <https://doi.org/10.3390/pathogens11050593>
- 1204 Maharachchikumbura SSN, Hyde KD, Jones EBG, McKenzie EHC et al (2015) Towards a natural classification
 1205 and backbone tree for Sordariomycetes. *Fungal Divers* 72 (1):199-301. <https://doi.org/10.1007/s13225-015-0331-z>
- 1206
- 1207 Maharachchikumbura SSN, Wanasinghe DN, Cheewangkoon R, Al-Sadi AM (2021) Uncovering the hidden
 1208 taxonomic diversity of fungi in Oman. *Fungal Divers* 106 (1):229-268. <https://doi.org/10.1007/s13225-020-00467-1>
- 1209
- 1210 Maharachchikumbura SSN, Chen Y, Ariyawansa HA, Hyde KD, Haelewaters D, Perera RH, Samarakoon MC,
 1211 Wanasinghe DN, Bustamante DE, Liu J-K, Lawrence DP, Cheewangkoon R, Stadler M (2021)
 1212 Integrative approaches for species delimitation in Ascomycota. *Fungal Divers* 109 (1):155-179.
 1213 <https://doi.org/10.1007/s13225-021-00486-6>
- 1214 Mäkelä MR, Donofrio N, de Vries RP (2014) Plant biomass degradation by fungi. *Fungal Genet Biol* 72:2-9.
 1215 <https://doi.org/10.1016/j.fgb.2014.08.010>

1216 Manni M, Berkeley MR, Seppely M, Simão FA, Zdobnov EM (2021) BUSCO update: novel and streamlined
1217 workflows along with broader and deeper phylogenetic coverage for scoring of eukaryotic, prokaryotic,
1218 and viral genomes. *Mol Biol Evol* 38 (10):4647-4654. <https://doi.org/10.1093/molbev/msab199>
1219 Mapuranga J, Zhang N, Zhang L, Chang J, Yang W (2022) Infection strategies and pathogenicity of biotrophic
1220 plant fungal pathogens. *Front Microbiol* 13. <https://doi.org/10.3389/fmicb.2022.799396>
1221 McCotter SW, Horianopoulos LC, Kronstad JW (2016) Regulation of the fungal secretome. *Curr Genet* 62
1222 (3):533-545. <https://doi.org/10.1007/s00294-016-0578-2>
1223 McDonald BA, Stukenbrock EH (2016) Rapid emergence of pathogens in agro-ecosystems: global threats to
1224 agricultural sustainability and food security. *Philos Trans R Soc B* 371 (1709):9.
1225 <https://doi.org/10.1098/rstb.2016.0026>
1226 Melén K, Krogh A, von Heijne G (2003) Reliability measures for membrane protein topology prediction
1227 algorithms. *J Mol Biol* 327 (3):735-744. [https://doi.org/10.1016/S0022-2836\(03\)00182-7](https://doi.org/10.1016/S0022-2836(03)00182-7)
1228 Mendgen K, Hahn M (2002) Plant infection and the establishment of fungal biotrophy. *Trends Plant Sci* 7 (8):352-
1229 356. [https://doi.org/10.1016/S1360-1385\(02\)02297-5](https://doi.org/10.1016/S1360-1385(02)02297-5)
1230 Meng Y, Lei Y, Gao J, Liu Y, Ma E, Ding Y, Bian Y, Zu H, Dong Y, Zhu X (2022) Genome sequence assembly
1231 algorithms and misassembly identification methods. *Mol Biol Rep* 49 (11):11133-11148.
1232 <https://doi.org/10.1007/s11033-022-07919-8>
1233 Mengiste T (2012) Plant immunity to necrotrophs. *Annu Rev Phytopathol* 50 (1):267-294.
1234 <https://doi.org/10.1146/annurev-phyto-081211-172955>
1235 Mesny F, Miyauchi S, Thiergart T, Pickel B et al (2021) Genetic determinants of endophytism in the *Arabidopsis*
1236 root mycobiome. *Nat Commun* 12 (1):7227. <https://doi.org/10.1038/s41467-021-27479-y>
1237 Miller JR, Koren S, Sutton G (2010) Assembly algorithms for next-generation sequencing data. *Genomics* 95
1238 (6):315-327. <https://doi.org/10.1016/j.ygeno.2010.03.001>
1239 Minh BQ, Schmidt HA, Chernomor O, Schrempf D, Woodhams MD, von Haeseler A, Lanfear R (2020) IQ-TREE
1240 2: new models and efficient methods for phylogenetic inference in the genomic era. *Mol Biol Evol* 37
1241 (5):1530-1534. <https://doi.org/10.1093/molbev/msaa015>
1242 Miyauchi S, Kiss E, Kuo A, Drula E et al (2020) Large-scale genome sequencing of mycorrhizal fungi provides
1243 insights into the early evolution of symbiotic traits. *Nat Commun* 11 (1):5125.
1244 <https://doi.org/10.1038/s41467-020-18795-w>
1245 Murigneux V, Rai SK, Furtado A, Bruxner TJC, Tian W, Harliwong I, Wei H, Yang B, Ye Q, Anderson E, Mao Q,
1246 Drmanac R, Wang O, Peters BA, Xu M, Wu P, Topp B, Coin LJM, Henry RJ (2020) Comparison of long-
1247 read methods for sequencing and assembly of a plant genome. *GigaScience* 9 (12).
1248 <https://doi.org/10.1093/gigascience/giaa146>
1249 Muszewska A, Steczkiewicz K, Stepniewska-Dziubinska M, Ginalski K (2017a) Cut-and-paste transposons in
1250 fungi with diverse lifestyles. *Genome Biol Evol* 9 (12):3463-3477. <https://doi.org/10.1093/gbe/evx261>
1251 Muszewska A, Steczkiewicz K, Stepniewska-Dziubinska M, Ginalski K (2019) Transposable elements contribute
1252 to fungal genes and impact fungal lifestyle. *Sci Rep* 9 (1):4307. <https://doi.org/10.1038/s41598-019-40965-0>
1253
1254 Muszewska A, Stepniewska-Dziubinska MM, Steczkiewicz K, Pawlowska J, Dziedzic A, Ginalski K (2017b)
1255 Fungal lifestyle reflected in serine protease repertoire. *Sci Rep* 7 (1):9147.
1256 <https://doi.org/10.1038/s41598-017-09644-w>
1257 Newman TE, Derbyshire MC (2020) The evolutionary and molecular features of broad host-range necrotrophy in
1258 plant pathogenic fungi. *Front Plant Sci* 11. <https://doi.org/10.3389/fpls.2020.591733>
1259 Nielsen H, Engelbrecht J, Brunak S, von Heijne G (1997) Identification of prokaryotic and eukaryotic signal

1260 peptides and prediction of their cleavage sites. *Protein Eng Des Sel* 10 (1):1-6.
1261 <https://doi.org/10.1093/protein/10.1.1>

1262 Nouws S, Bogaerts B, Verhaegen B, Denayer S, Piérard D, Marchal K, Roosens NHC, Vanneste K, De
1263 Keersmaecker SCJ (2020) Impact of DNA extraction on whole genome sequencing analysis for
1264 characterization and relatedness of Shiga toxin-producing *Escherichia coli* isolates. *Sci Rep* 10 (1):14649.
1265 <https://doi.org/10.1038/s41598-020-71207-3>

1266 O'Connell RJ, Thon MR, Hacquard S, Amyotte SG et al (2012) Lifestyle transitions in plant pathogenic
1267 *Colletotrichum* fungi deciphered by genome and transcriptome analyses. *Nat. Genet.* 44 (9):1060-1065.
1268 <https://doi.org/10.1038/ng.2372>

1269 Oberti H, Dalla Rizza M, Reyno R, Murchio S, Altier N, Abreo E (2020) Diversity of *Claviceps paspali* reveals
1270 unknown lineages and unique alkaloid genotypes. *Mycologia* 112 (2):230-243.
1271 <https://doi.org/10.1080/00275514.2019.1694827>

1272 Oggenfuss U, Badet T, Wicker T, Hartmann FE, Singh NK, Abraham L, Karisto P, Vonlanthen T, Mundt C,
1273 McDonald BA, Croll D (2021) A population-level invasion by transposable elements triggers genome
1274 expansion in a fungal pathogen. *eLife* 10:e69249. <https://doi.org/10.7554/eLife.69249>

1275 Omrane S, Audéon C, Ignace A, Duplaix C, Aouini L, Kema G, Walker A-S, Fillinger S (2017) Plasticity of the
1276 *MFS1* promoter leads to multidrug resistance in the wheat pathogen *Zymoseptoria tritici*. *mSphere* 2
1277 (5):e00393-00317. <https://doi.org/10.1128/mSphere.00393-17>

1278 Parker IM, Gilbert GS (2004) The evolutionary ecology of novel plant-pathogen interactions. *Annu Rev Ecol Evol*
1279 *Syst* 35 (1): 675-700. <https://doi.org/10.1146/annurev.ecolsys.34.011802.132339>

1280 Parks DH, Chuvochina M, Waite DW, Rinke C, Skarszewski A, Chaumeil P-A, Hugenholtz P (2018) A
1281 standardized bacterial taxonomy based on genome phylogeny substantially revises the tree of life. *Nat*
1282 *Biotechnol* 36 (10):996-1004. <https://doi.org/10.1038/nbt.4229>

1283 Pellegrin C, Morin E, Martin FM, Veneault-Fourrey C (2015) Comparative analysis of secretomes from
1284 ectomycorrhizal fungi with an emphasis on small-secreted proteins. *Front Microbiol* 6.
1285 <https://doi.org/10.3389/fmicb.2015.01278>

1286 Peter M, Kohler A, Ohm RA, Kuo A, Krützmann J, Morin E, Arend M, Barry KW, Binder M, Choi C, Clum A,
1287 Copeland A, Grisel N, Haridas S, Kipfer T, LaButti K, Lindquist E, Lipzen A, Maire R, Meier B,
1288 Mihaltcheva S, Molinier V, Murat C, Pöggeler S, Quandt CA, Sperisen C, Tritt A, Tisserant E, Crous PW,
1289 Henrissat B, Nehls U, Egli S, Spatafora JW, Grigoriev IV, Martin FM (2016) Ectomycorrhizal ecology
1290 is imprinted in the genome of the dominant symbiotic fungus *Cenococcum geophilum*. *Nat Commun* 7
1291 (1):12662. <https://doi.org/10.1038/ncomms12662>

1292 Petersen TN, Brunak S, von Heijne G, Nielsen H (2011) SignalP 4.0: discriminating signal peptides from
1293 transmembrane regions. *Nat Methods* 8 (10):785-786. <https://doi.org/10.1038/nmeth.1701>

1294 Phurailatpam L, Mishra S (2020) Role of plant endophytes in conferring abiotic stress tolerance. In:
1295 Hasanuzzaman M (ed) *Plant Ecophysiology and Adaptation under Climate Change: Mechanisms and*
1296 *Perspectives II: Mechanisms of Adaptation and Stress Amelioration*. Springer Singapore, Singapore, pp
1297 603-628. https://doi.org/10.1007/978-981-15-2172-0_22

1298 Plett JM, Martin F (2011) Blurred boundaries: lifestyle lessons from ectomycorrhizal fungal genomes. *Trends*
1299 *Genet* 27 (1):14-22. <https://doi.org/10.1016/j.tig.2010.10.005>

1300 Presti LL, Lanver D, Schweizer G, Tanaka S, Liang L, Tollot M, Zuccaro A, Reissmann S, Kahmann R (2015)
1301 Fungal effectors and plant susceptibility. *Annu Rev Plant Biol* 66 (1):513-545.
1302 <https://doi.org/10.1146/annurev-arplant-043014-114623>

1303 Promputtha I, Hyde KD, McKenzie EHC, Peberdy JF, Lumyong S (2010) Can leaf degrading enzymes provide

1304 evidence that endophytic fungi becoming saprobes? *Fungal Divers* 41 (1):89-99.
1305 <https://doi.org/10.1007/s13225-010-0024-6>

1306 Promputtha I, Lumyong S, Dhanasekaran V, McKenzie EHC, Hyde KD, Jeewon R (2007) A phylogenetic
1307 evaluation of whether endophytes become saprotrophs at host senescence. *Microb Ecol* 53 (4):579-590.
1308 <https://doi.org/10.1007/s00248-006-9117-x>

1309 R Core Team (2022) R: A language and environment for statistical computing. in R Foundation for Statistical
1310 Computing. (2020).

1311 Rai M, Agarkar G (2016) Plant–fungal interactions: What triggers the fungi to switch among lifestyles? *Crit Rev*
1312 *Microbiol* 42 (3):428-438. <https://doi.org/10.3109/1040841X.2014.958052>

1313 Ramirez-Garcia A, Pellon A, Rementeria A, Buldain I et al (2018) *Scedosporium* and *Lomentospora*: an updated
1314 overview of underrated opportunists. *Med Mycol* 56 (suppl_1):S102-S125.
1315 <https://doi.org/10.1093/mmy/myx113>

1316 Raudys SJ, Jain AK (1991) Small sample size effects in statistical pattern recognition: recommendations for
1317 practitioners. *IEEE Trans Pattern Anal Mach Intell* 13 (3):252-264. <https://doi.org/10.1109/34.75512>

1318 Rawlings ND, Barrett AJ, Thomas PD, Huang X, Bateman A, Finn RD (2017) The MEROPS database of
1319 proteolytic enzymes, their substrates and inhibitors in 2017 and a comparison with peptidases in the
1320 PANTHER database. *Nucleic Acids Res* 46 (D1):D624-D632. <https://doi.org/10.1093/nar/gkx1134>

1321 Read J, Pfahringer B, Holmes G, Frank E (2011) Classifier chains for multi-label classification. *Mach Learn* 85
1322 (3):333-359. <https://doi.org/10.1007/s10994-011-5256-5>

1323 Réblová M, Miller AN, Rossman AY, Seifert KA, Crous PW et al (2016) Recommendations for competing sexual-
1324 asexually typified generic names in Sordariomycetes (except Diaporthales, Hypocreales, and
1325 Magnaporthales). *IMA Fungus* 7 (1):131-153. <https://doi.org/10.5598/imafungus.2016.07.01.08>

1326 Řehulka J, Kubátová A, Hubka V (2016) *Cephalotheca sulfurea* (Ascomycota, Sordariomycetes), a new fungal
1327 pathogen of the farmed rainbow trout *Oncorhynchus mykiss*. *J Fish Dis* 39 (12):1413-1419.
1328 <https://doi.org/10.1111/jfd.12477>

1329 Rinke C, Chuvochina M, Mussig AJ, Chaumeil P-A, Davín AA, Waite DW, Whitman WB, Parks DH, Hugenholtz
1330 P (2021) A standardized archaeal taxonomy for the genome taxonomy database. *Nat Microbiol* 6 (7):946-
1331 959. <https://doi.org/10.1038/s41564-021-00918-8>

1332 Rizzo DM, Garbelotto M (2003) Sudden oak death: endangering California and Oregon forest ecosystems. *Front*
1333 *Ecol Environ* 1 (4):197-204. [https://doi.org/10.1890/1540-9295\(2003\)001\[0197:SODECA\]2.0.CO;2](https://doi.org/10.1890/1540-9295(2003)001[0197:SODECA]2.0.CO;2)

1334 Rodrigues AM, de Hoog GS, de Camargo ZP (2016) *Sporothrix* species causing outbreaks in animals and humans
1335 driven by animal–animal transmission. *PLoS Pathog* 12 (7):e1005638.
1336 <https://doi.org/10.1371/journal.ppat.1005638>

1337 Salim D, Bradford WD, Freeland A, Cady G, Wang J, Pruitt SC, Gerton JL (2017) DNA replication stress restricts
1338 ribosomal DNA copy number. *PLoS Genet* 13 (9):e1007006.
1339 <https://doi.org/10.1371/journal.pgen.1007006>

1340 Santini A, Ghelardini L, De Pace C, Desprez-Loustau ML et al (2013) Biogeographical patterns and determinants
1341 of invasion by forest pathogens in Europe. *New Phytol* 197 (1), 238-250. <https://doi.org/10.1111/j.1469-8137.2012.04364.x>

1343 Senft AD, Macfarlan TS (2021) Transposable elements shape the evolution of mammalian development. *Nat Rev*
1344 *Genet* 22 (11):691-711. <https://doi.org/10.1038/s41576-021-00385-1>

1345 Seong K, Krasileva KV (2023) Prediction of effector protein structures from fungal phytopathogens enables
1346 evolutionary analyses. *Nat Microbiol* 8 (1):174-187. <https://doi.org/10.1038/s41564-022-01287-6>

1347 Shang Y, Feng P, Wang C (2015) Fungi that infect insects: altering host behavior and beyond. *PLoS Pathog* 11

1348 (8):e1005037. <https://doi.org/10.1371/journal.ppat.1005037>

1349 Shen X, Opulente DAx, Kominek J, Zhou X et al (2018) Tempo and mode of genome evolution in the budding
1350 yeast subphylum. *Cell* 175 (6):1533-1545.e1520. <https://doi.org/10.1016/j.cell.2018.10.023>

1351 Shen X, Steenwyk JL, LaBella AL, Opulente DA, Zhou X, Kominek J, Li Y, Groenewald M, Hittinger CT, Rokas
1352 A (2020) Genome-scale phylogeny and contrasting modes of genome evolution in the fungal phylum
1353 Ascomycota. *Sci Adv* 6 (45). <https://doi.org/10.1126/sciadv.abd0079>

1354 Simão FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM (2015) BUSCO: assessing genome
1355 assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 31 (19):3210-3212.
1356 <https://doi.org/10.1093/bioinformatics/btv351>

1357 Simmons DR, Kepler RM, Renner SA, Groden E (2015) Phylogeny of *Hirsutella* species (Ophiocordycipitaceae)
1358 from the USA: remedying the paucity of *Hirsutella* sequence data. *IMA Fungus* 6 (2):345-356.
1359 <https://doi.org/10.5598/imafungus.2015.06.02.06>

1360 Singh VK, Meena M, Zehra A, Tiwari A, Dubey MK, Upadhyay RS (2014) Fungal toxins and their impact on
1361 living systems. In: Kharwar RN, Upadhyay RS, Dubey NK, Raghuwanshi R (eds) *Microbial Diversity
1362 and Biotechnology in Food Security*. Springer India, New Delhi, pp 513-530.
1363 https://doi.org/10.1007/978-81-322-1801-2_47

1364 Šmarda P, Bureš P, Horová L, Leitch IJ, Mucina L, Pacini E, Tichý L, Grulich V, Rotreklová O (2014) Ecological
1365 and evolutionary significance of genomic GC content diversity in monocots. *PNAS* 111 (39):E4096-
1366 E4102. <https://doi.org/https://doi.org/10.1073/pnas.1321152111>

1367 Smits THM (2019) The importance of genome sequence quality to microbial comparative genomics. *BMC
1368 Genomics* 20 (1):662. <https://doi.org/10.1186/s12864-019-6014-5>

1369 Solla A, Bohnens J, Collin E, Diamandis S, Franke A, Gil L, Burón M, Santini A, Mitterpergher L, Pinon J,
1370 Broeck AV (2005) Screening european elms for resistance to *Ophiostoma novo-ulmi*. *For Sci* 51 (2):134-
1371 141.

1372 Sone T, Fukiya S, Kodama M, Tomita F (2000) Molecular structure of rDNA repeat unit in *Magnaporthe grisea*.
1373 *Biosci Biotechnol Biochem* 64 (8):1733-1736. <https://doi.org/10.1271/bbb.64.1733>

1374 Song H, Kim K-T, Park S-Y, Lee G-W, Choi J, Jeon J, Cheong K, Choi G, Hur J-S, Lee Y-H (2022) A comparative
1375 genomic analysis of lichen-forming fungi reveals new insights into fungal lifestyles. *Sci Rep* 12
1376 (1):10724. <https://doi.org/10.1038/s41598-022-14340-5>

1377 Sordo M, Zeng Q On sample size and classification accuracy: a performance comparison. In: Berlin, Heidelberg,
1378 2005. *Biological and Medical Data Analysis*. Springer Berlin Heidelberg, pp 193-201

1379 Spanu PD, Abbott JC, Amselem J, Burgis TA et al (2010) Genome expansion and gene loss in powdery mildew
1380 fungi reveal tradeoffs in extreme parasitism. *Sciences* 330 (6010):1543-1546.
1381 <https://doi.org/10.1126/science.1194573>

1382 Sperschneider J, Dodds PN (2022) EffectorP 3.0: prediction of apoplasmic and cytoplasmic effectors in fungi and
1383 oomycetes. *Mol Plant-Microbe Interact* 35 (2):146-156. <https://doi.org/10.1094/mpmi-08-21-0201-r>

1384 Stanke M, Diekhans M, Baertsch R, Haussler D (2008) Using native and syntenically mapped cDNA alignments
1385 to improve de novo gene finding. *Bioinformatics* 24 (5):637-644.
1386 <https://doi.org/10.1093/bioinformatics/btn013>

1387 Stergiopoulos I, Wit PJGMd (2009) Fungal effector proteins. *Annu Rev Phytopathol* 47 (1):233-263.
1388 <https://doi.org/10.1146/annurev.phyto.112408.132637>

1389 Sugita R, Tanaka K (2022) *Thyridium* revised: Synonymisation of *Phialemoniopsis* under *Thyridium* and
1390 establishment of a new order, Thyridiales. *MycoKeys* 86:147-176. <https://10.3897/mycokeys.86.78989>

1391 Sun Y, Liu N, Samarakoon MC, Jayawardena RS, Hyde KD, Wang Y (2021) Morphology and phylogeny reveal

1392 Vamsapriyaceae fam. nov. (Xylariales, Sordariomycetes) with two novel *Vamsapriya* species. *J Fungi* 7
1393 (11):891. <https://doi.org/10.3390/jof7110891>

1394 Tongcham P, Supa P, Pornwongthong P, Prasitmeeboon P (2020) Mushroom spawn quality classification with
1395 machine learning. *Comput Electron Agric* 179:105865. <https://doi.org/10.1016/j.compag.2020.105865>

1396 Tørresen OK, Star B, Mier P, Andrade-Navarro MA, Bateman A, Jarnot P, Gruca A, Grynberg M, Kajava AV,
1397 Promponas VJ, Anisimova M, Jakobsen KS, Linke D (2019) Tandem repeats lead to sequence assembly
1398 errors and impose multi-level challenges for genome and protein databases. *Nucleic Acids Res* 47
1399 (21):10994-11006. <https://doi.org/10.1093/nar/gkz841>

1400 Tortorano AM, Prigitano A, Esposto MC, Arsic Arsenijevic V et al (2014) European Confederation of Medical
1401 Mycology (ECMM) epidemiological survey on invasive infections due to *Fusarium* species in europe.
1402 *Eur J Clin Microbiol Infect Dis* 33 (9):1623-1630. <https://doi.org/10.1007/s10096-014-2111-1>

1403 Treangen TJ, Salzberg SL (2012) Repetitive DNA and next-generation sequencing: computational challenges and
1404 solutions. *Nat Rev Genet* 13 (1):36-46. <https://doi.org/10.1038/nrg3117>

1405 Troy GC, Panciera DL, Pickett JP, Sutton DA, Gene J, Cano JF, Guarro J, Thompson EH, Wickes BL (2013)
1406 Mixed infection caused by *Lecythophora canina* sp. nov. and *Plectosphaerella cucumerina* in a German
1407 shepherd dog. *Med Mycol* 51 (5):455-460. <https://doi.org/10.3109/13693786.2012.754998>

1408 van Kan JAL (2006) Licensed to kill: the lifestyle of a necrotrophic plant pathogen. *Trends Plant Sci* 11 (5):247-
1409 253. <https://doi.org/10.1016/j.tplants.2006.03.005>

1410 Wang D, Tian L, Zhang DD, Song J, Song SS, Yin CM, Zhou L, Liu Y, Wang B-L, Kong Z-Q, Klosterman SJ, Li
1411 J-J, Wang J, Li T-G, Adamu S, Subbarao KV, Chen J-Y, Dai X-F (2020) Functional analyses of small
1412 secreted cysteine-rich proteins identified candidate effectors in *Verticillium dahliae*. *Mol Plant Pathol* 21
1413 (5):667-685. <https://doi.org/10.1111/mpp.12921>

1414 Wang M, Tan X-M, Liu F, Cai L (2018) Eight new *Arthrinium* species from China. *Mycosphere* 34.
1415 <https://doi.org/10.3897/mycokeys.34.24221>

1416 Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhou B, Flavell A, Leroy P, Morgante M, Panaud O,
1417 Paux E, SanMiguel P, Schulman AH (2007) A unified classification system for eukaryotic transposable
1418 elements. *Nat Rev Genet* 8 (12):973-982. <https://doi.org/10.1038/nrg2165>

1419 Wijayawardene NN, Hyde KD, Dai DQ, Sanchez-Garcia M, et al (2022) Outline of fungi and fungus-like taxa-
1420 2021. *Mycosphere* 13 (1):53-453. <https://doi.org/10.5943/mycosphere/13/1/2>

1421 Wingfield MJ, De Beer ZW, Slippers B, Wingfield BD, Groenewald JZ, Lombard L, Crous PW (2012) One fungus,
1422 one name promotes progressive plant pathology. *Mol Plant Pathol* 13 (6):604-613.
1423 <https://doi.org/10.1111/j.1364-3703.2011.00768.x>

1424 Wu W, Chen W, Liu S, Wu J, Zhu Y, Qin L, Zhu B (2021) Beneficial relationships between endophytic bacteria
1425 and medicinal plants. *Front. Plant Sci* 12. <https://doi.org/10.3389/fpls.2021.646146>

1426 Xu C, Jackson SA (2019) Machine learning and complex biological data. *Genome Biol* 20 (1):76.
1427 <https://doi.org/10.1186/s13059-019-1689-0>

1428 Xu J, Yang X, Lin Q (2014) Chemistry and biology of *Pestalotiopsis*-derived natural products. *Fungal Divers* 66
1429 (1):37-68. <https://doi.org/10.1007/s13225-014-0288-3>

1430 Xu S, Dai Z, Guo P, Fu X, Liu S, Zhou L, Tang W, Feng T, Chen M, Zhan L, Wu T, Hu E, Jiang Y, Bo X, Yu G
1431 (2021) ggtreeExtra: compact visualization of richly annotated phylogenetic data. *Mol Biol Evol* 38
1432 (9):4039-4042. <https://doi.org/10.1093/molbev/msab166>

1433 Yin C, Ramachandran SR, Zhai Y, Bu C, Pappu HR, Hulbert SH (2019) A novel fungal effector from *Puccinia*
1434 *graminis* suppressing RNA silencing and plant defense responses. *New Phytol* 222 (3):1561-1572.
1435 <https://doi.org/10.1111/nph.15676>

1436 Yu G, Smith DK, Zhu H, Guan Y, Lam TT-Y (2017) GGTREE: an R package for visualization and annotation of
1437 phylogenetic trees with their covariates and other associated data. *Methods Ecol Evol* 8 (1):28-36.
1438 <https://doi.org/10.1111/2041-210X.12628>

1439 Yu G, Xian L, Xue H, Yu W, Rufian JS, Sang Y, Morcillo RJL, Wang Y, Macho AP (2020) A bacterial effector
1440 protein prevents MAPK-mediated phosphorylation of SGT1 to suppress plant immunity. *PLoS Pathog*
1441 16 (9):e1008933. <https://doi.org/10.1371/journal.ppat.1008933>

1442 Zeng T, Holmer R, Hontelez J, te Lintel-Hekkert B, Marufu L, de Zeeuw T, Wu F, Schijlen E, Bisseling T, Limpens
1443 E (2018) Host- and stage-dependent secretome of the arbuscular mycorrhizal fungus *Rhizophagus*
1444 *irregularis*. *Plant J* 94 (3):411-425. <https://doi.org/10.1111/tpj.13908>

1445 Zhang H, Yohe T, Huang L, Entwistle S, Wu P, Yang Z, Busk PK, Xu Y, Yin Y (2018) dbCAN2: a meta server for
1446 automated carbohydrate-active enzyme annotation. *Nucleic Acids Res* 46 (W1):W95-W101.
1447 <https://doi.org/10.1093/nar/gky418>

1448 Zhang Y, Yang H, Turra D, Zhou S, Ayhan DH, DeIulio GA, Guo L, Broz K, Wiederhold N, Coleman JJ, Donnell
1449 KO, Youngster I, McAdam AJ, Savinov S, Shea T, Young S, Zeng Q, Rep M, Pearlman E, Schwartz DC,
1450 Di Pietro A, Kistler HC, Ma L-J (2020) The genome of opportunistic fungal pathogen *Fusarium*
1451 *oxysporum* carries a unique set of lineage-specific chromosomes. *Commun Biol* 3 (1):50.
1452 <https://doi.org/10.1038/s42003-020-0770-2>

1453 Zhao Z, Liu H, Wang C, Xu JR (2013) Comparative analysis of fungal genomes reveals different plant cell wall
1454 degrading capacity in fungi. *BMC Genomics* 14 (1):274. <https://doi.org/10.1186/1471-2164-14-274>

1455 Zieliński B, Sroka-Oleksiak A, Rymarczyk D, Piekarczyk A, Brzychczy-Włoch M (2020) Deep learning approach
1456 to describe and classify fungi microscopic images. *PLoS One* 15 (6):e0234806.
1457 <https://doi.org/10.1371/journal.pone.0234806>

1458

1459