

Lifestyles shape genome size and gene content in fungal pathogens

Anna Fijarczyk^{1,2,3,4,*}, Pauline Hessenauer^{1,2}, Richard C. Hamelin^{1,2,5}, Christian R. Landry^{1,2,3,4,6}

¹Département de Biologie, Université Laval, Québec, G1V 0A6 Québec, Canada;

²Institut de Biologie Intégrative et des Systèmes (IBIS), Université Laval, Québec, G1V 0A6 Québec, Canada;

³PROTEO, Le réseau québécois de recherche sur la fonction, la structure et l'ingénierie des protéines, Université Laval, Québec, G1V 0A6 Québec, Canada;

⁴Centre de Recherche en Données Massives (CRDM), Université Laval, Québec, G1V 0A6 Québec, Canada;

⁵Department of Forest and Conservation Sciences, The University of British Columbia, V6T 1Z4 Vancouver, Canada;

⁶Département de Biochimie, Microbiologie et Bioinformatique, Université Laval, Québec, G1V 0A6 Québec, Canada;

*corresponding author: anna.fijarczyk.1@ulaval.ca

Keywords:

fungal pathogens, genome evolution, fungal lifestyles, genome size, insect vectors, Sordariomycetes

Summary

Fungi have a wide range of lifestyles and hosts. We still know little about the impact of lifestyles on their genome architecture. Here, we combined and annotated 562 fungal genomes from the class Sordariomycetes and examined the coevolution between 12 genomic and two lifestyle traits: pathogenicity and insect association. We found that many pathogens tend to evolve a larger number of protein-coding genes, tRNA genes, and have larger non-repetitive genome sizes than non-pathogenic species. In contrast, species with a pathogenic or symbiotic relationship with insects have smaller genome sizes and genes with longer exons; they also have fewer genes if they are vectored by insects, compared to species not associated with insects. Our study demonstrates that pathogen genome size and complexity are the result of an interplay between drift, imposed by symbiosis and small effective population size, which leads to genome contraction, and the adaptive role of gene amplification, which leads to genome expansion.

Introduction

Variation in genome size has been explained through mutually non-exclusive hypotheses invoking neutral and adaptive processes (reviewed in (Blommaert, 2020)). For instance, neutral hypotheses are supported by the correlation of genome size with non-coding DNA content (Lynch & Conery, 2003; Petrov, 2002), whereas natural selection is invoked to explain the strong correlation between genome size and cell size (Gregory, 2002). So far, genome size has been shown to correlate with various other traits in eukaryotes, including body size in many species (e.g. salamanders (Decena-Segarra et al., 2020) and amphipods (Hultgren et al., 2018)), life-history traits (e.g. Pinus (Grotkopp et al., 2004), or birds (J. P. Yu et al., 2020)), flowering time across latitudinal cline in *Zea mays* (Bilinski et al., 2018) or higher fitness in algae (Malerba et al., 2020). Understanding the contribution of neutral and adaptive processes to genome size evolution may thus require analyzing multiple mechanisms underlying the structural complexity of genomes, such as the amount and distribution of coding and non-coding DNA (Petrov, 2001).

Fungi display a great variety of genome sizes (from 8 to 177.6 Mbp) (Mohanta & Bae, 2015). Observations support both neutral and adaptive mechanisms in shaping this diversity. The evolution of genome size in ascomycetes on a broad phylogenetic scale supports the mutation-hazard hypothesis (MHH) according to which stronger drift facilitates the accumulation of slightly deleterious non-coding DNA, including mobile genetic elements and of introns (Lynch & Conery, 2003). As a consequence, species with smaller effective population sizes (N_e) are more prone to genome expansion. As expected with MHH in eukaryotes, ascomycetes with larger genomes have more introns, lower gene density, and higher transposable elements (TE) activity than species with smaller genomes (Kelkar & Ochman, 2012). On the other hand, some host-specialized fungi and obligate pathogens, species that are also associated with small N_e and strong drift, experience genome degradation and gene loss, which is explained by the elimination of accumulated deleterious DNA in the long term due to bias towards deletions (Kelkar & Ochman, 2012; Mira et al., 2001).

The role of adaptive evolution in shaping genome size has also been supported in fungi. The most obvious observation is the strong correlation between the number of genes and genome size (Stajich, 2017). For instance, the diversification of specific gene families is often associated with pathogenicity in fungi (Baroncelli et al., 2016; Muszewska et al., 2011; Sipos et al., 2017). Coincidentally, fungal plant pathogens harbor some of the largest known genomes across fungi (Stajich, 2017). Some other genetic mechanisms that contribute to genome expansion have been attributed to rapid adaptation of pathogens to hosts and these include accessory chromosomes (Croll & McDonald, 2012) or TE activity (Mat Razali et al., 2019), often linking the emergence of

pathogenicity to genome size expansion (Raffaele & Kamoun, 2012). For instance, TE expansions are known to play a role in the emergence of new virulence genes in emerging pathogens (Bao et al., 2017; Wacker et al., 2021).

The ecological aspects of fungal lifestyles, for instance, the dependence on the host species, can have an effect both on the gene repertoire of the species as well as on their N_e , and as a consequence on the amount of non-coding DNA and genome size. Pathogens can undergo phases of high clonality which significantly reduces their N_e , however other mechanisms such as accessory chromosomes, or chromosome duplication can increase N_e in parts of the genomes (Stukenbrock & Croll, 2014). Other pathogens, such as the opportunistic human pathogen *Aspergillus fumigatus*, persist in the environment where they can maintain high levels of nucleotide polymorphism, consistent with large N_e (Rhodes et al., 2022). On the contrary, obligate endoparasites, endosymbionts, and pathogens requiring vectors have typically small N_e due to their association with a single host species and transmission bottlenecks. For instance, the endoparasites Microsporidia have one of the smallest known fungal genomes, showing signatures of strong genome contractions, including reduced gene repertoire (Katinka et al., 2001). Therefore studying the impact of pathogenic lifestyle on genome size in fungi requires a wider ecological context, including their symbiotic or mutualistic relationships with other species.

In this study, we investigate how variation in fungal genome size and its underlying potentially neutral or slightly deleterious genomic features (e.g. intron number, repeat content) and adaptation (number of genes) coevolve with lifestyles. We consider pathogenic and symbiotic relationships with plants and insects. We aim to answer the following questions: 1) whether the fungal genome size and complexity coevolve with lifestyle, and 2) if there are common genomic features related to the evolution of pathogenicity. To answer these questions, we focus on the ascomycete class of Sordariomycetes fungi, which have rich genomic resources and include species with a wide range of lifestyles, from plant and insect pathogens, through wood decomposers, and endophytes, to opportunistic human pathogens. We studied over 560 Sordariomycetes genomes and examined how various genomic features coevolved with the diversity of lifestyles in this taxonomic group.

Results

Quality of genome assemblies of 562 Sordariomycetes fungal species

We analyzed 562 genomes from the class Sordariomycetes (Ascomycota) and 11 outgroup species together comprising fungi mostly represented by plant pathogens, saprotrophs, and entomopathogens (Supplementary Table 1). All species were

classified as pathogenic (including pathogens of plants, animals, and fungi) or non-pathogenic based on a literature search. Initially, representative genome assemblies for known species of Sordariomycetes and an outgroup (n=584 in total) were downloaded from NCBI (n=580) or other resources (n=4, Supplementary Table 1) and complemented with 21 genomes sequenced and assembled in this study (Supplementary Table 2). In total, 605 assemblies were filtered for contaminants and short contigs and analyzed for quality. Gene models were inferred *ab initio* for all assemblies. Thirty-two assemblies were excluded due to lack of many single-copy conserved genes (more than 25% of conserved genes were missing) or an excessive number of inferred gene models (*Blumeria graminis*), arriving at a final number of 573 genome assemblies for further analyses. The number of genes was systematically underestimated (10% on average) compared to the gene numbers that had been submitted to NCBI for the corresponding species, but underestimation was not different between pathogenic or non-pathogenic fungi (Wilcoxon rank sum test, $w = 2994$, $p = 0.57$). Species genomes ranged from 20.7 Mbp in *Ceratocystiopsis brevicornis* (CBS 137839) to 110.9 Mbp in *Ophiocordyceps sinensis* (IOZ07) and the *ab initio* gene models ranged from 6,280 in *Ambrosiella xylebori* (CBS 110.61) to 17,878 in *Fusarium langsethiae* (Fe2391).

Genomic traits correlating with genome size

The maximum likelihood tree generated from 1000 concatenated single-copy conserved proteins has a 100% bootstrap support for all major nodes except for one (support of 82% for the split between two subclades of Ophiostomatales, Figure 1A). The topology placed the Microascales order as a sister clade to Hypocreales, which is consistent with a topology based on four nuclear genetic markers (Hongsanan et al., 2017) but not consistent with the topology presented on the JGI MycoCosm (accessed July 10th 2022, <https://mycocosm.jgi.doe.gov/mycocosm/home>), where Hypocreales and Glomerellales form sister clades, and Microascales is their outgroup. To check if the discrepancy is caused by the inference method, we reconstructed the topology with two other methods: i) from 1000 individual protein trees using a coalescent-based method, and ii) with a maximum likelihood approach from concatenated 250 single-copy conserved proteins with separate partition models assigned to each protein. Both methods gave nearly identical topology to the first one and showed maximal support for Microascales and Hypocreales as sister clades. Thus, the primary maximum likelihood phylogeny (Figure 1A) based on 1000 concatenated protein alignments was used for all downstream analyses.

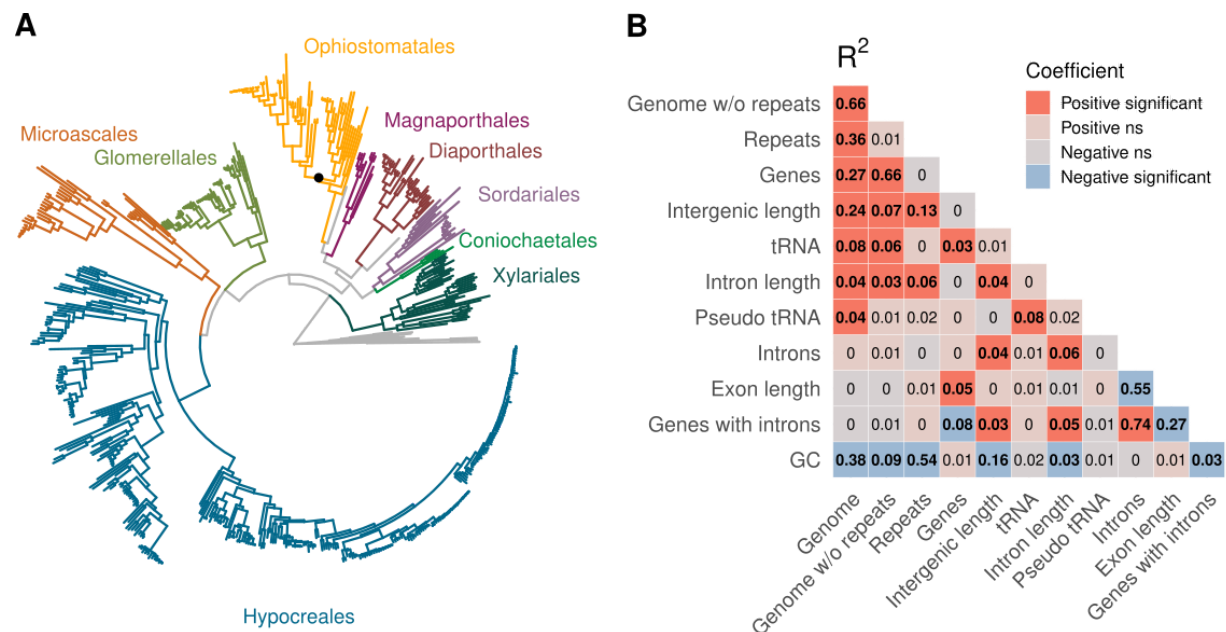


Figure 1. Several genomic traits are correlated with genome size in Sordariomycetes. A. Maximum likelihood tree based on 1000 concatenated protein sequence alignments calculated with IQ-TREE using ultrafast bootstrap approximation (n=573 species). The largest orders are indicated with different colors. Bootstrap support for all major clades except one within Ophiostomatales (82%, black dot) reaches 100%. **B.** R-squared values derived from generalized least squares models fitted to all pairwise combinations of genomic traits. Colors indicate positive (red) or negative (blue) coefficients obtained from fitting one trait to another. Coefficients significantly different from zero (t-test, adjusted $p < 0.05$) are highlighted with dark colors and bold R^2 values. Genomic traits include genome size (genome), size of the assembly excluding repeat content (genome w/o repeats), the fraction of repeat content (repeats), the number of genes (genes), mean intergenic length (intergenic length), the number of tRNA genes (tRNA), mean intron length (intron length), the number of pseudo tRNA genes (pseudo tRNA), the mean number of introns per gene (introns), the mean exon length (exon length), the fraction of genes with introns (genes with introns), and GC content (GC). Correlation plots are shown in figure supplement 1, and principal component analysis on all genomic traits in figure supplement 2. Raw data underlying figures are in Figure 1-Source Data 1-3.

We used phylogenetic generalized least squares to estimate the correlation between genome size (bp) and 11 genomic traits, including the number of genes, the fraction of repeat content, size of the assembly excluding repeat content (bp), GC content, the mean number of introns per gene, mean intron length (bp), mean exon length (bp), the fraction of genes with introns, mean intergenic length (bp) and the number of tRNA and pseudo tRNA genes. As expected from the MHH in eukaryotes, in Sordariomycetes, genome size is positively correlated with the fraction of repeat content, intergenic length, and mean intron length (Figure 1B, Figure 1-figure supplement 1). Genome size is also negatively correlated with the GC content, as a consequence of the negative correlation between repeat and GC content. Similar to other fungi, genome size is also

strongly positively correlated with the number of genes (including tRNA genes) and consequently to genome size without the repeat content. The numbers of introns and exon length are strongly negatively correlated with each other, but weakly with genome size, and they explain mostly the second principal component in the PCA analysis (Figure 1-figure supplement 2). Similarly, the fraction of genes with introns is also not correlated with genome size (Figure 1B, Figure 1-figure supplement 1). These results support that drift has shaped genome size of Sordariomycetes fungi through expansion of repeat content and intergenic DNA, and that adaptive evolution has shaped genome size through the proliferation of genes. However, some signatures of drift, such as frequency of introns and size of exons do not show a clear relationship with genome size, suggesting opposing consequences of drift on small and large genomes.

Pathogenicity coevolves with the number of genes

Our dataset comprises 357 pathogens, 202 non-pathogens, and 4 species with undetermined pathogenicity traits. Pathogens can be found in every order except for Sordariales (Figure 2A). Even though the trait could not be inferred with high accuracy at the root, the dispersion of pathogenic species across the phylogeny implies that the evolution and loss of pathogenicity have occurred multiple times in Sordariomycetes history (Figure 2A).

We used three methods (trait coevolution in BayesTraits, phylogenetic logistic regression, and machine learning-based classification of pathogens vs. non-pathogens) to test which genomic traits coevolve with pathogenicity. Using a discrete reversible jump mcmc (rjMCMC) method in BayesTraits, we tested dependent versus independent models of evolution of genomic traits with pathogenicity, by transforming all genomic traits into binary traits (high/low) with respect to their median. All genomic traits except for total genome size, fraction of genes with introns, and exon length show correlated evolution with pathogenicity (Figure 2B, Supplementary Table 3). Posterior probabilities of transition rates show that the most frequent gain of pathogenicity occurs when the number of genes, introns, and genome size without repeats are high (Figure 2B, Figure 2-figure supplement 1). We also observe gains of pathogenicity for low values of repeats (Figure 2B, Figure 2-figure supplement 1). We fitted a covarion model to test if dependent evolution with pathogenicity is present across the tree or is rather limited to some branches. There is a significantly better fit for varying rates of evolution than the strictly dependent model for all genomic traits, except for the number of genes, in which the support for the varying rates of evolution was the weakest (average log bayes factor across five runs is 4.5, Supplementary Table 3).

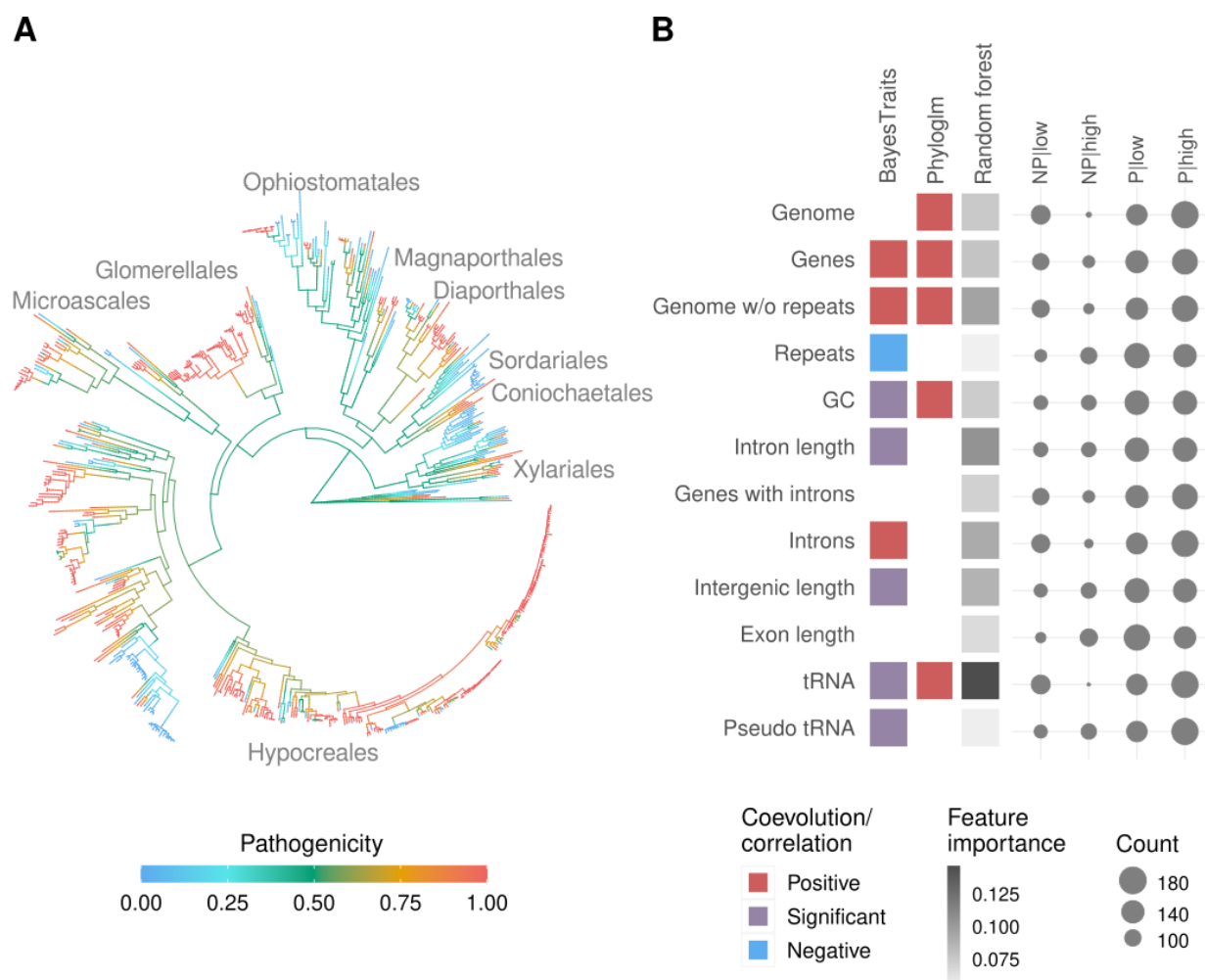


Figure 2. Genomic traits are evolving in concert with pathogenicity. **A.** Ancestral state reconstruction of pathogenicity. Colors correspond to estimated likelihoods of pathogenicity, calculated with ace function from the phytools package in R, using the “ER” model of discrete evolution with a single rate for the transition from pathogen to non-pathogen and vice versa. **B.** Results of the three methods for detecting correlation between pathogenicity and genomic traits: BayesTraits, phylogenetic logistic regression (phyloglm), and random forest classifier. Red and blue colors indicate positive and negative, respectively, coevolution/correlation of the trait with pathogenicity. Violet indicates those traits which coevolve with pathogenicity according to BayesTraits but a single direction of correlation could not be determined. Circles indicate frequencies of pathogenic (P) and nonpathogenic (NP) species with low (below median) or high (above median) values for 12 genomic traits. Posterior probabilities of transition rates are shown in figure supplement 1, details of the Random Forest classification in figure supplement 2, and results of analysis for ten subsets are in figure supplement 3. Raw data underlying figures are in Figure2-Source Data 1-4.

In the second approach, the phylogenetic logistic regression was fitted to the pathogenicity trait. This analysis confirmed a positive correlation between pathogenicity and the number of genes and genome size without repeats, similarly to the analysis with

BayesTraits, and revealed a positive correlation with GC content, and the number of tRNA genes (Figure 2B, Supplementary Table 4). This method also positively correlated genome size with pathogenicity.

Finally, we used machine learning classification to rank the most important among 12 genomic traits for predicting pathogenicity, with the random forest classifier providing the highest accuracy predictions. We tested the performance of the classifier including phylogenetic relationships, which did not substantially improve the accuracy (AUC=0.7 vs AUC=0.71, see Materials and Methods, Figure 2-figure supplement 2). Top-ranked genomic features included the number of tRNAs, intron length, assembly size without repeats, and the number of introns (Figure 2B, Figure 2-figure supplement 2).

To account for the uneven contribution of pathogenic and non-pathogenic species in the analysis we run the BayesTraits and phyloglm analysis for subsets of species, with the same number of pathogenic and non-pathogenic species. Pathogenic species were randomly subsampled ten times giving a total of ten subsets. BayesTraits analysis confirmed coevolution of six genomic traits with pathogenicity in all ten subsets (genes, genome without repeats, repeats, intron length, intergenic length and tRNAs), and the coevolution of GC, intron number, and pseudo tRNAs in six to nine subsets (Figure 2-figure supplement 3). Advantage of gains over loss of genomic features was, however, not evident in BayesTraits analysis. Phyloglm analysis confirmed positive correlation between genome size, genes, genome without repeats, and tRNAs in most subsets, but not GC content (Figure 2-figure supplement 3).

All analyses were also repeated for a dataset that excluded repeated species (n=563 distinct species in total), giving nearly identical results, except that BayesTraits inferred weaker coevolution between pathogenicity and the number of introns, and the phyloglm model gave no support for the significant correlation with GC content (Supplementary Tables 3 and 4, Materials and Methods).

Overall, several lines of evidence indicate that most pathogens evolve through an increase in gene number, genome size excluding repeats, number and/or length of introns, and the number of tRNAs. This is consistent with the observation that most non-pathogenic species very rarely show high values for these traits (Figure 2B).

Genome reductions in species associated with insects

After reconstructing ancestral gene numbers in the Sordariomycetes phylogeny, as predicted with phylogenetic approaches (Figure 2B), several clades enriched for plant pathogens show an increased number of genes (Figure 3). These include Diaporthales (D), Magnaporthales (Ma), Glomerellales (G), and many *Fusarium* pathogens within the

Hypocreales (H1) clade. On the other hand, a couple of clades including insect pathogens (eg. H2.4) and a few plant pathogens (M1, O) carry some of the smallest gene numbers and genome sizes without repeats (Figure 3, Figure 3-figure supplement 1).

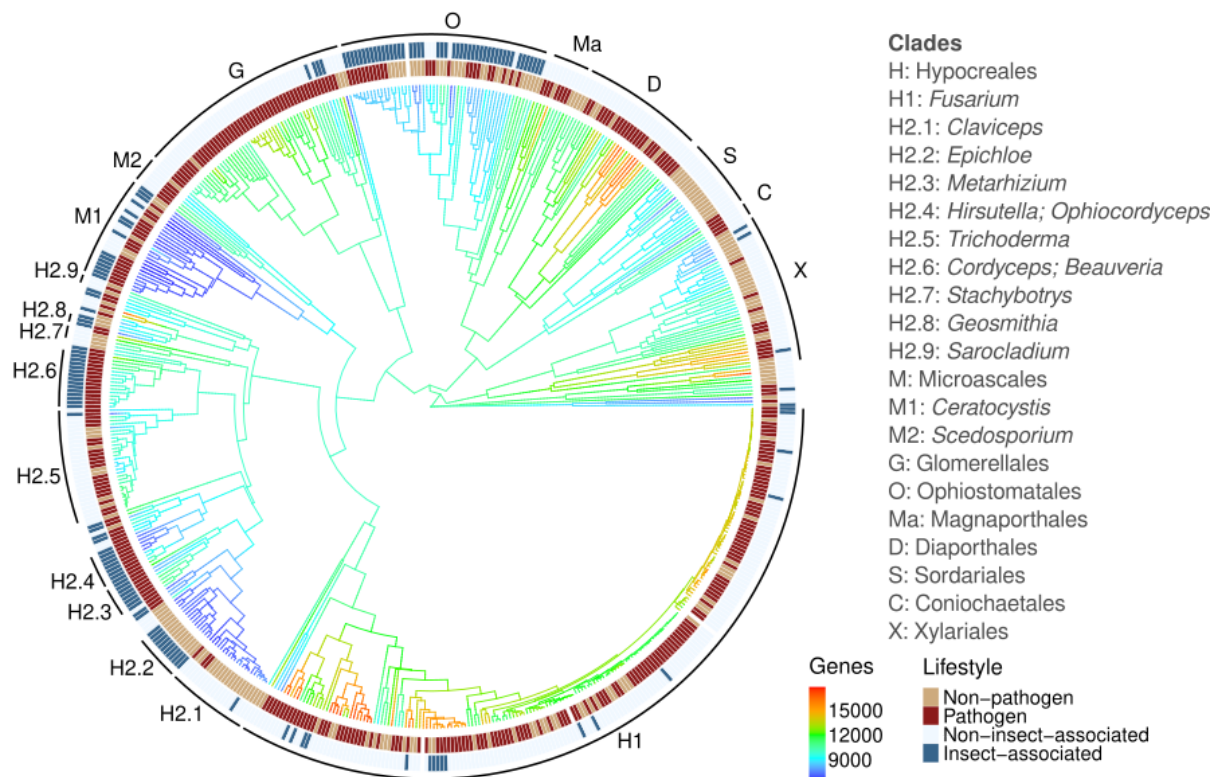


Figure 3. Substantial gene losses in some Sordariomycetes clades. Ancestral gene numbers were estimated with fastAnc function in R package phytools. A circular heatmap on the inside indicates species annotated as pathogens (red) or non-pathogenic (yellow), and the heatmap on the outside indicates species annotated as insect-associated (blue) or not (yellow). Black lines mark selected clades. The list on the right shows a corresponding order name or representative genera for highlighted clades. Ancestral states for other genomic features are shown in figure supplement 2. Raw data underlying figures are in Figure3-Source Data 1-2.

Symbiosis is one trait that often leads to streamlining of microbial genomes (Katinka et al., 2001), therefore, we checked if association with insects can explain gene and genome reductions in some clades. We consider species to be associated with insects if they are symbionts, mutualists, or pathogens of insects, and we annotated a subset of clades in the phylogeny that are either composed of mostly insect-associated species or not (Figures 3-4, Supplementary Table 5). First, we compared genomic traits between several pairs of insect- and non-insect-associated clades (groups a to e, Figure 4, Figure 4-figure supplement 1). Clades O (group a), M1 (group b), and H2.8 (group c) comprise insect vector-transmitted tree endophytes and pathogens. Clades H2.6 (group

d), H2.4, and H2.3 (group e) are specialized entomopathogens. Clade H2.2 comprises fly-transmitted grass symbionts (group e). In four groups (a to d), insect-associated species have smaller genomes, smaller genomes without repeats, and longer exons than their sister clades (Figure 4). Smaller genomes are accompanied in most cases by fewer introns and fewer genes (groups b, c, and a except for clade S) or/and smaller intergenic sizes (group d). All clades from group e show overall fewer differences between insect and non-insect-associated clades, which can be explained by the fact that, unlike other groups, they originated from an insect-associated ancestor (probability of insect-association in the ancestor = 0.97, 95% CI [0.96 - 0.99], Supplementary Table 5). Consequently, all clades in group e except for one (H2.3) have among the smallest genomes without repeats and the fewest genes out of all groups. By comparing current values of genomic traits with estimated ancestral states in the most recent common ancestor for each group, we found a consistent decrease in the total size of the genome, genome size without repeats, number of genes, and number of introns in all clades which are in symbiotic/mutualistic relationship with insects (clades O, M1, H2.8, and H2.2, Figure 4-figure supplement 2).

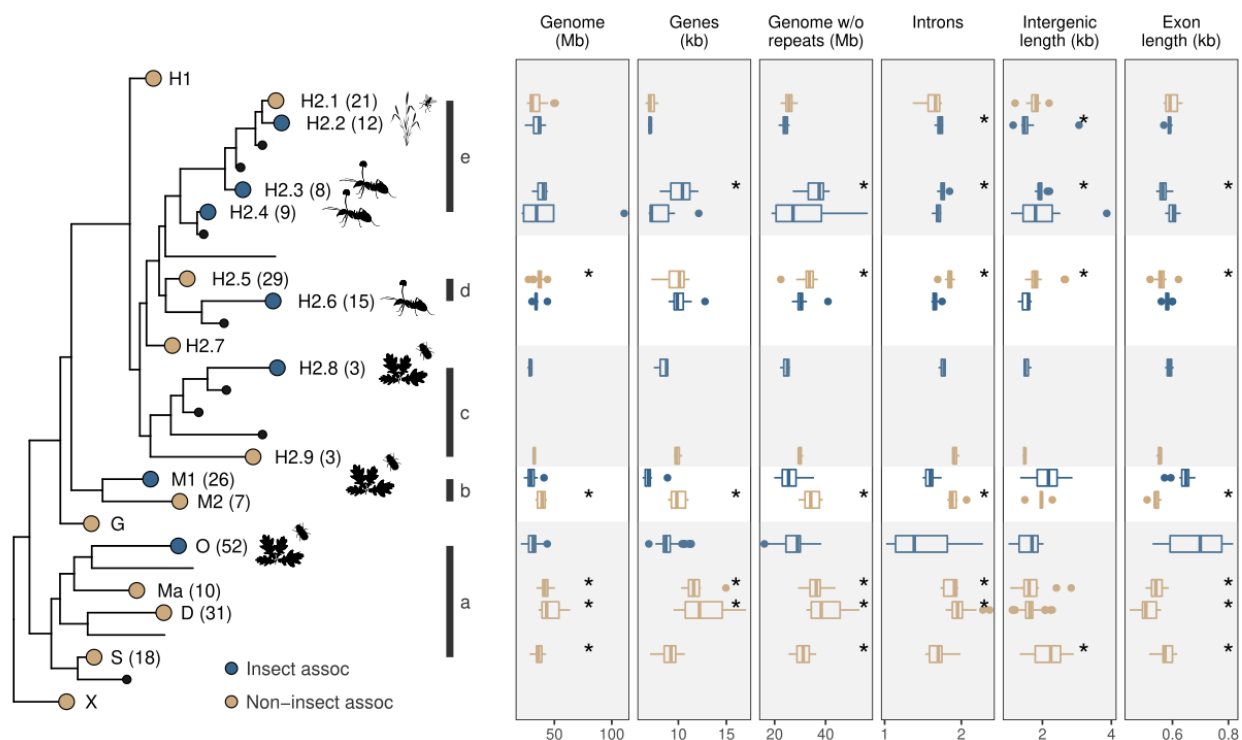


Figure 4. Genomic traits vary across insect and non-insect-associated clades. Five groups (a-e marked with vertical black lines) comparing insect-associated (blue nodes) and non-insect-associated (yellow) clades are shown. Black filled circles indicate clades composed of >1 species. Numbers shown in the parentheses indicate clade abundances. Asterisks indicate the comparisons with statistically significant pairwise differences, between blue and brown clades within each group (Mann-Whitney test, adjusted $p < 0.05$). Clades in group c have too few species for testing. Clades O (group a), M1 (group b), H2.8 (group c), and H2.2 (group e) are insect mutualists or symbionts, whereas clades H2.6 (group d) and clades H2.4 and H2.3 (group e) are insect pathogens. Comparison for the rest of the traits is shown in figure supplement 1, and the fold change of extant species relative to the ancestral node is shown in figure supplement 2. Raw data underlying figures are in Figure4-Source Data 1-2.

Next, we used three methods as described in the previous section, to test the coevolution of genomic traits with insect association. All three methods consistently suggest the presence of longer exons in insect-associated species, and two out of three methods confirm smaller genome size but not fewer genes in insect-associated species (Figure 5, Figure 5-figure supplements 1 and 2, Supplementary Tables 6-7). Analyses conducted on a dataset that excluded repeated species ($n=563$ in total), gave nearly identical results (Supplementary Tables 6-7, Materials and Methods), with the exception that BayesTraits inferred additional inverse coevolution between insect association and genome size without repeats, whereas phyloglm model gave no support for negative correlation of insect association with genome size without repeats and the number of introns.

Overall, these results show that insect-associated fungi are characterized by genes with longer exons and usually fewer introns, and have smaller genomes. In addition, clades in symbiotic/mutualistic relationships with insects experience a reduction in genome size excluding repeats and a reduction in the number of genes.

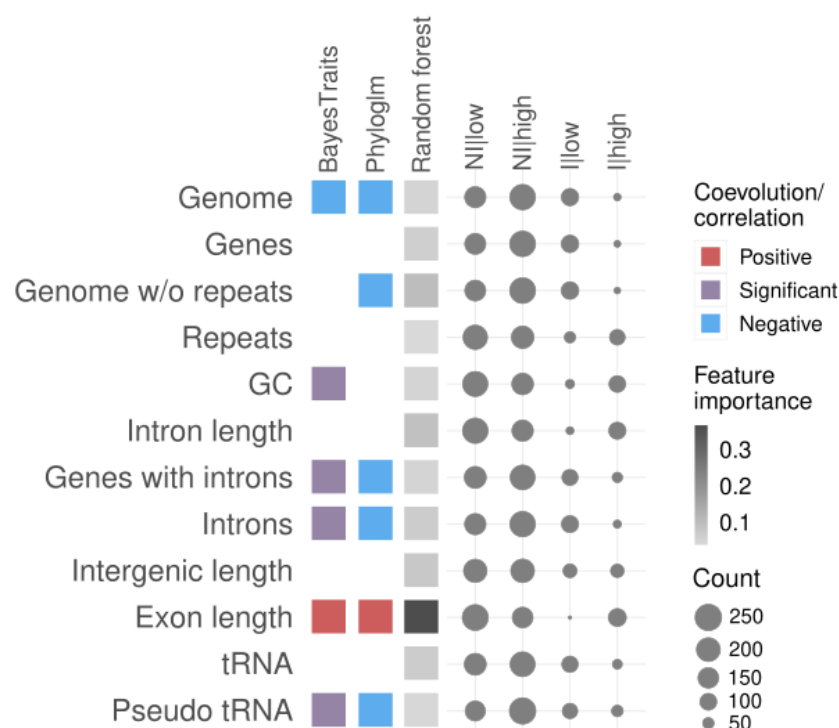


Figure 5. Genomic traits coevolve with insect-associated taxa. Results of the three methods used for detecting correlation between insect-associated species and genomic traits: BayesTraits, phylogenetic logistic regression (phyloglm), and random forest classifier. Circles indicate frequencies of insect-associated (I) and non-insect-associated (NI) species with low (below median) or high (above median) values for 12 genomic traits. I - insect association, NI - no insect association, low - genomic trait value below the median, high - genomic trait value above the median. Posterior probabilities of transition rates are shown in figure supplement 1, details of the Random Forest classification in figure supplement 2, and losses and gains of orthogroups are shown in figure supplement 3. Raw data underlying figures are in Figure5-Source Data 1-3.

Losses in gene families are more frequent than gains

To evaluate the dynamics of gene losses and gains in Sordariomycetes, we estimated rates of gains and losses of gene families (groups of orthologous genes with ≥ 1 member) from the root of the tree for a subset (112) of species. We found that gene family losses are overall more frequent than gene family gains, in particular at the deeper branches of the tree (Figure 5-figure supplement 3). The strongest gene family losses are located on deep branches leading to subclades of H2 clade, on the branch leading to clade M1, and on the branches leading to clades O, Ma, and D. Gene family losses dominate over gains in most species but are overcome by gains only in pathogenic species from clades H1, M2, G, D, and some X. The clades H1, M2, and G did not undergo major ancient losses in their genomes. This analysis shows that gene losses are prevalent across most Sordariomycete clades, and are not restricted to

insect-associated fungi. Net gains of genes since the Sordariomycete common ancestor are observed only in strictly pathogenic clades.

Insect-associated plant pathogens preferentially lose genes important for host colonization

To resolve whether gene loss in insect-associated fungi is uniform across all genes or constrained to specific functional classes, we annotated *ab initio* gene models with KOG functions (Tatusov et al., 2003), annotated genome assemblies with Secondary Metabolite Clusters (Blin et al., 2019) (SMC), and searched gene models against databases of CAZymes (Cantarel et al., 2009), peptidases (Rawlings et al., 2018) (MEROPS, M), and transcription factor pfam domains (Mistry et al., 2021) (TF). The CAZymes and SMC genes are critical for overcoming external host barriers and for entry into the host (Cantarel et al., 2009; Scharf et al., 2014). Therefore, selection pressures acting on them are expected to vary with lifestyles. Ancestral content of annotated gene classes was estimated for nodes at the split of each clade with its sister clade (most recent common ancestor with a sister clade) and compared with mean content observed in the clade to obtain a fold change.

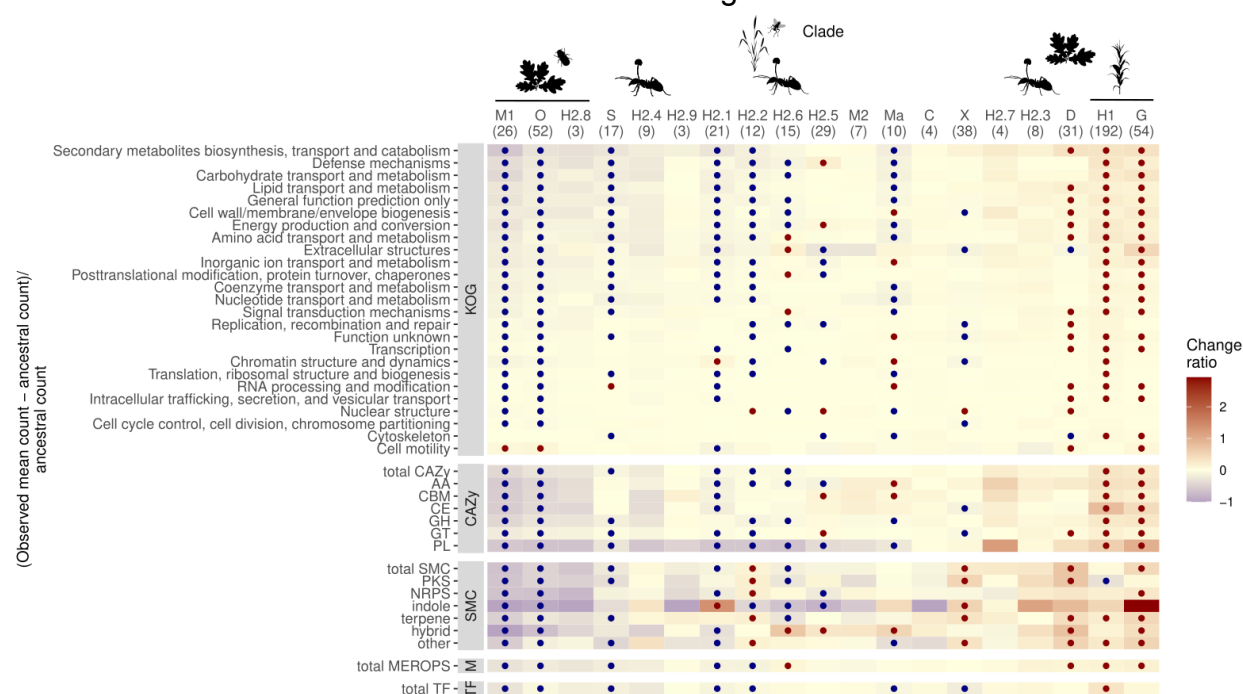


Figure 6. Insect-vectored clades lose genes involved in breaking plant host barriers.

Heatmap shows the fold change of genes/clusters relative to the ancestral state ((observed - ancestral)/ ancestral state). Clades are shown in columns (clade names correspond to the ones in Figure 3) with the number of clade members in parentheses, functional classes are shown in rows. The dots indicate significant gain (red) or loss (blue) of genes/clusters across clade members estimated from 100 rounds of bootstrapping of 10 species in clades with ≥ 10 members. SMC: secondary metabolite clusters, M: Merops, TF: transcription factors. Drawings indicate clades that are associated with insects (M1, O, H2.8, H2.4, H2.2, H2.6, H2.3), and three right-most clades, which are specialized tree (D) and crop pathogens (H1, G). Same heatmap but for pathogenic species only is shown in figure supplement 1. Raw data underlying figures are in Figure6-Source Data 1-2.

The insect-vectored clades Ophiostomatales (O), Microascales (M1), and *Geosmithia* (H2.8) exhibit the highest gene losses across nearly all KOG groups (Figure 6). To a lesser degree gene loss is also observed in non-pathogenic clades Sordariales (S), H2.1 as well as other insect-associated clades H2.2, H2.4, and H2.6. Only one insect-pathogenic clade H2.3 has an opposite pattern. On the opposite side of the spectrum, strictly pathogenic *Fusarium* from Hypocreales (H1), Glomerellales (G), and Diaporthales (D), exhibit gene gains across most KOG groups. Loss of genes is the strongest in fungal clades in symbiotic/mutualistic relationships with insects (M1, O, H2.8), with the most prominent loss in genes involved in secondary metabolites synthesis, transport, and catabolism, defense mechanisms, and carbohydrate transport and metabolism. Indeed, all analyzed types of CAZymes, SMC, as well as peptidases, and transcription factors have contracted in these clades (Figure 6).

Only genes related to the cytoskeleton and cell motility exhibit no gain or change in number in insect-vectored O, M1, and H2.8. These clades include some important emerging tree plant pathogens, therefore, to test whether losses are experienced equally by pathogenic and non-pathogenic members of these clades, gene gain/loss analyses were repeated for pathogenic species only. Results confirm extensive losses in similar clades and KOG groups (Supplementary Figure 7). Notably, KOG groups most commonly lost in insect-associated clades (and clades derived from them), include (apart from the ones mentioned above) lipid transport and metabolism, cell wall/membrane/envelope biogenesis, or energy production and conservation, and are the same KOG groups which undergo expansion in plant pathogenic clades (Figure 6, Figure 6-figure supplement 1). These results show that plant pathogens and those vectored by insects have distinct repertoires of genes and experience contrasting dynamics in genes important for host colonization.

Gene structure of lost genes

Genomes of insect-vectored species are characterized by fewer genes and genes with longer exons and fewer introns, the last two being strongly negatively correlated (Figure 1B). This pattern can be explained either by the more frequent loss of intron-rich genes or by a general trend towards a less complex gene structure. To test this, we first compared gene structure (number of exons, exon, and intron lengths) in single copy one-to-one orthologs between species belonging to two insect-vectored clades (O and M1) and their non-insect-vectored sister clades (D and M2). Longer exons among one-to-one orthologs of insect-vectored species would suggest a general trend toward less complex gene structure. Indeed, one-to-one orthologs of insect-vectored clades have longer exons (by 286 bp in O, Figure 7A, and by 140 bp in M1, Figure 7-figure supplement 1), and fewer introns (by 0.58 in O, Figure 7A, and by 0.33 in M1, Figure 7-figure supplement 1). Intron length is either longer (by 9 bp in M1, Figure 7-figure supplement 1) or shorter (by 10 bp in O, Figure 7A) compared to non-insect-vectored species.

As a second test, we looked at the gene structure of gene families relative to the occupancy of the gene family in the clade. If complex genes are more likely to get lost, we should observe rare genes (present in a few clade members) to have more introns (and exons) and shorter exons than common genes. Gene families that are frequently lost since the common ancestor of the clade, indeed have shorter exons, however, they also have fewer exons, contrasting our first hypothesis that more complex genes are lost more often (Figure 7B and Figure 7-figure supplement 2, Poisson regression with 2 predictors: exon length and count, and species occupancy as a response, $p < 3.61e-08$ for all clades). This trend is significant in both insect- and non-insect-vectored clades (Figure 7B and Figure 7-figure supplement 2).

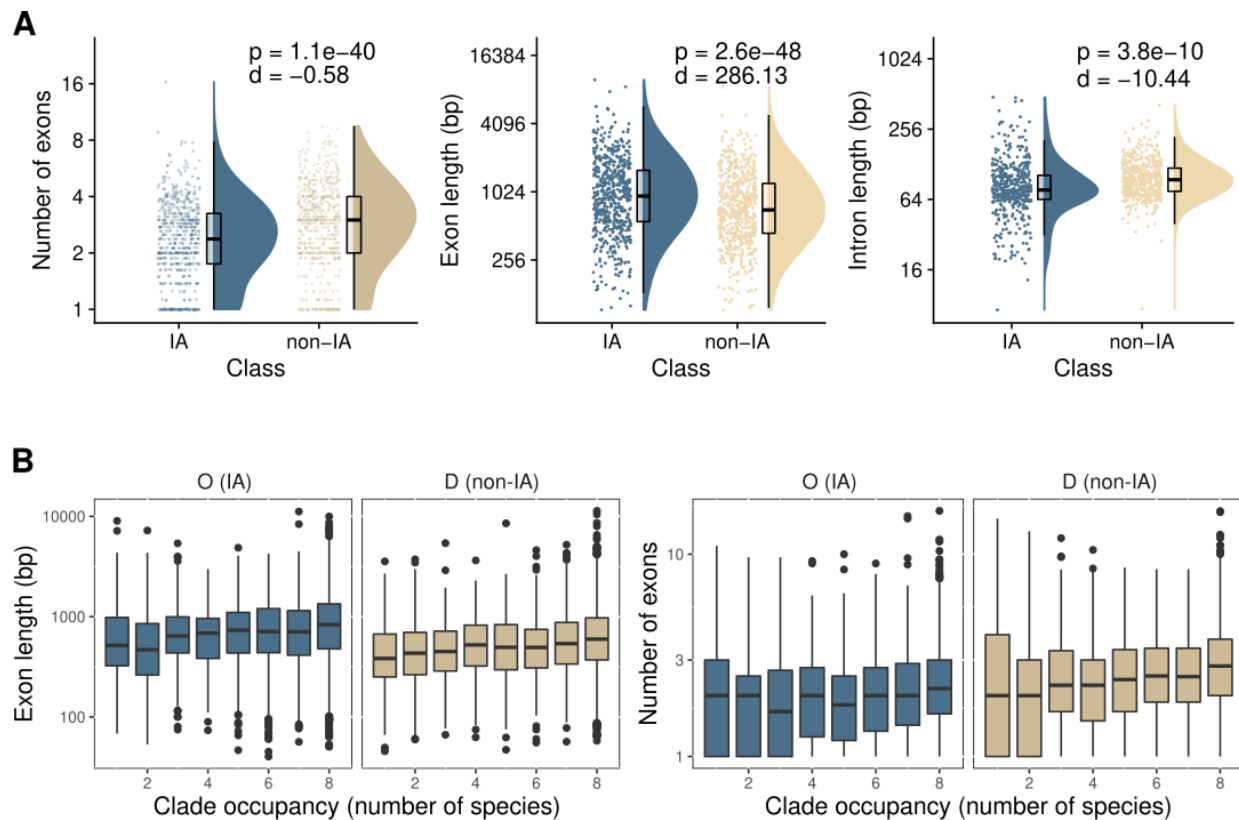


Figure 7. Gene structure changes in insect-associated clades affect all genes. A.

Comparison of gene structures between one-to-one single-copy orthologs ($n=583$) from insect-associated (IA) clade O (blue) and the corresponding non-insect-associated (non-IA) clade D (yellow). Intron length was compared only between orthologs with at least one intron ($n=476$). Ortholog features were averaged across eight species. Boxplots show medians, first and third quartiles, and lines span minimum to maximum values excluding outliers. p - p -values estimated with the paired two-sided Wilcoxon signed-rank test; d - mean differences between IA and non-IA orthologs. **B.** Average exon length and the number of exons of gene families (orthogroups) present in different frequencies across up to eight clade members. Corresponding figures for the clades M1 and M2 are shown in figure supplement 1. Raw data underlying figures are in Figure7-Source Data 1-2.

Discussion

We investigated the impact of fungal lifestyle on the evolution of genome size and complexity using 562 Sordariomycetes genomes. We analyzed 12 genomic traits and two lifestyle traits: pathogenicity and association with insects, and found several patterns of genome evolution in this group of fungi. First, genome size is strongly positively correlated with the repeat content and the number of genes in the genome. Second, fungi with the most streamlined genomes (smaller genomes, fewer introns, and genes) belong to groups that associate with insects that serve as hosts, vectors, or symbionts to fungi, or to groups directly derived from them. Third, pathogens exhibit distinct patterns in gene dynamics, with most plant and animal pathogens having increased gene numbers, tRNAs, and the number of introns per gene, and plant pathogens vectored by insects showing opposite patterns with frequent losses of genes and introns.

Impact of the lifestyle on genome size in Sordariomycetes

The strongest predictor of Sordariomycetes genome size is repeat content and to a lesser degree other non-coding DNA, such as intergenic length and intron length. Accumulation and expansion of non-coding DNA support the mutation-hazard hypothesis in explaining genome size expansions, according to which a small effective population size (N_e) leads to stronger drift and fixation of deleterious DNA. Eukaryotic species with smaller N_e accumulate more noncoding DNA (Lynch, 2006). In the case of fungi, drift can affect the amount of the repeat content and other non-coding DNA in some species or clades, but this trend is not associated with any particular lifestyle (note however that pathogens tend to have genes with more introns). Both pathogens (*Monosporascus*, *Ophiocordyceps*) and non-pathogens (*Claviceps*, *Epichloe*) experienced bouts of repeat expansions, which could be linked to some stochastic events of low N_e in these species, rather than a specific occupied ecological niche.

As one of the important contributors to repeat content, transposable elements (TEs) are often associated with pathogens. There are well-documented cases of virulence factors emerging due to transposable element activity. In spite of a few cases of pathogens with an expanded repeat content, suggesting an increased TE activity, in general, pathogenic Sordariomycetes do not carry many repeats. In contrast, an average pathogen experiences a decrease in the number of repeats. This pattern can be caused by an effective selection in removing deleterious noncoding DNA, and/or by a strong bias toward deletions in the genome. Pathogens with an expanded TE content may indicate rare cases of fungi that developed means to accommodate fast-evolving deleterious DNA, for instance in accessory chromosomes, chromosomal duplications,

separate genome compartments, or experienced events of low N_e , due to transmissions, or multiple rounds of population extermination (Stukenbrock & Croll, 2014).

The second major predictor of genome size is gene content. Expansion and diversification of genes are usually linked to adaptation to a specific niche, for instance in mycorrhizal (Miyauchi et al., 2020) or wood decomposing fungi (Franco et al., 2022). We find that an overall increase in gene repertoire is linked to a pathogenic lifestyle, and this pattern is visible in pathogens from different Sordariomycetes clades, except for vector-transmitted tree pathogens. The genomes of pathogens are located in the upper part of the genome size scale in our dataset, a pattern that is mostly driven by gene expansions since the repeat content is limited in most pathogens. Nevertheless, the largest genomes belong to species that exhibit both high repeat content and a large number of genes (for instance pathogenic *Colletotrichum* species).

Sordariomycetes with the smallest genomes are almost exclusively endosymbionts, insect endoparasites, are vectored by insects, or are directly derived from clades associated with insects. They evolved smaller genomes, genes with longer exons (and fewer introns), and in the case of insect-vectored species, they have also lost many genes. It has been suggested, that lack of diversification of carbohydrate enzymes in *Ophiostoma* pathogens, might have been caused by the fact that beetle vectors take responsibility for penetrating host tissues and entering their vascular system (Comeau et al., 2014). These patterns of genome streamlining are reminiscent of the reductive genome evolution in endosymbiotic and endoparasitic prokaryotes, as well as some endoparasitic fungi. Gene-rich genomes of prokaryotes show an inverse correlation between drift and genome size, explained mainly by the bias towards deletions (Kuo et al., 2009; Mira et al., 2001). Loss of genes in endosymbiotic and some free-living bacteria has been also explained by selection against non-essential genes (Giovannoni et al., 2014), a loss of redundant genes with drift (Moran et al., 2008), or increased mutation rate (Bourguignon et al., 2020). Our results show that in two major insect-vectored clades (Ophiostomatales and Microascales) the change in gene architecture towards longer exons and fewer introns is consistent across all the genes, and is not caused by the retention of genes with such structure. Moreover, these clades exhibit the longest branches in the phylogeny (measured by amino-acid substitutions), an observation that can be explained by the fixation of deleterious substitutions when N_e is small. Consistent gene structure and long branches support the impact of drift on the genome complexity in species with the reduced N_e .

Gene expansions in Sordariomycetes pathogens

Pathogenicity has evolved multiple times across Sordariomycetes in nearly every order. As we were unable to assign the probability of pathogenicity with high accuracy for deeper branches in the tree, our findings confirm that this is a fast-evolving trait. In spite of the emergence of pathogenicity across independent taxonomic groups, several genomic traits proved to be reappearing among most pathogens, and these are high overall numbers of genes, tRNA genes, larger genomes without repeats, genes with longer introns, and less repeat content. Expanded gene numbers are most evident in plant pathogenic species such as *Fusarium* (Hypocreales), *Colletotrichum*, and *Verticillium* (Glomerellales) or *Diaporthe* (Diaporthales), but also in some (not all) entomopathogens (*Metarhizium*), mycopathogens (*Trichoderma*) and opportunistic human pathogens (*Scedosporium*, *Sporothrix*). This implies that most pathogens (excluding those that are vectored by insects) benefit from expanded gene repertoires, for example in a more effective establishment in host tissues and/or defense against host immunological mechanisms. Indeed we see that the largest gains of genes in pathogens correspond to genes involved in secondary metabolite synthesis, carbohydrate metabolism, and defense mechanisms.

One group of genes that coevolves with pathogenicity are tRNAs. These observations are more difficult to interpret. However, in fungi and other microbes, post-translational modifications of tRNAs are known to play a role in triggering virulence and evading host defense (Chen et al., 2021; Hinsch et al., 2016; Morrison et al., 2017). Duplications of existing tRNAs may be an alternative way to maintain a variety of tRNAs to cope with the host immune system.

Another trait observed in pathogens is an increased number of introns per gene. Previous studies have shown that genes of fungal ancestors were intron-rich and several events of massive intron loss occurred across fungi, leading to intron-poor groups including Saccharomycotina as well as Pezizomycotina, which comprise Sordariomycetes (Lim et al., 2021). In spite of an overall low number of introns, we observe an increasing trend of introns in pathogens. Introns play an important role in alternative splicing, and this process has been found to occur more frequently in pathogenic species, affecting genes involved in dimorphism and stress response, essential functions when entering the host environment (Grützmann et al., 2014; Muzafar et al., 2021).

Different modes of genome evolution in Sordariomycetes pathogens

Our results show that the genome complexity of Sordariomycetes fungi is influenced by an interplay between drift and adaptation, both of which are affected by a specific lifestyle of the pathogen. The distinctions between lifestyles are visible in particular in the evolution of the gene content and gene structure. Based on these traits alone, three

groups of pathogens can be distinguished: i) specialized plant pathogens with expanded gene repertoires and intron-rich genes, ii) insect pathogens with low to high gene expansions counteracting strong ancient gene contractions, and with long-exon genes, and iii) insect-vectored plant pathogens with overall gene contractions including many genes responsible for host adaptation and pathogenicity, and long-exon genes. Considering these observations, one can imagine a hypothetical scenario of genome evolution in Sordariomycetes, in which species constantly reduce their genome size due to the overall bias towards deletions. This trend would be exacerbated in species that evolve in symbiosis, where drift is strong and many genes become redundant and thus dispensable due to the presence of host gene products, eventually leading to genome size contraction. But the trend of gene loss would be reversed if the species evolved as pathogens, or switched from insect-symbiosis to insect-pathogenicity, in which case their gene repertoires would increase through duplications. Finally, species that are neither pathogenic nor symbiotic, would constantly reduce their genomes due to an overall high deletion rate (but not small N_e), though not to the same degree as symbiotic species (as can be exemplified by non-pathogenic Sordariales).

We show that fungi follow distinct evolutionary trajectories to gain their pathogenic potential. It is worth noting, however, that our dataset contains more pathogenic than non-pathogenic species, and some genetic clades are dominated by pathogens. Therefore, to fully capture all evolutionary trends during transition to pathogenicity, more effort needs to be placed on sequencing non-pathogenic species.

Materials and Methods

Genome assemblies

14 strains (11 x *Ophiostoma* and 3 x *Leptographium*, Supplementary Table 2) were selected for high-depth short-read Illumina sequencing. Isolates were grown on Malt Extract Agar medium (15 g l⁻¹ agar, 30 g l⁻¹ malt extract, and 5 g l⁻¹ mycological peptone), and DNA was extracted with acetyl trimethylammonium bromide chloroform protocol. The fungal isolates were handled in a facility that has received a Plant Pest Containment Level 1 certification by the Canadian Food Inspection Agency. Library preparation and sequencing were conducted at the G  nome Qu  bec Innovation Center (Montr  al, Canada). One strain (*O. quercus*) was sequenced with Illumina NovaSeq (paired-end 150 bp), and the rest with Illumina HighSeq X (paired-end 150 bp). Data quality was inspected with Fastqc v0.11.2/8 (Andrews, 2010). Ten strains were used for *de novo* genome assembly (Supplementary Table 2). The depth of coverage of the generated data was between 32x and 555x (Supplementary Table 2). Reads were trimmed for adapters with Trimmomatic v0.33/0.36 (Bolger et al., 2014) using options 'ILLUMINACLIP:adapters.fa:6:20:10 MINLEN:21' and overlapping reads were merged with bbmerge from BBTools v36/v37 (Bushnell et al., 2017). *De novo* genome assemblies were generated with SPAdes v3.9.1 (Bankevich et al., 2012) with options '-k 21,33,55,77,99 --careful'. Mitochondrial DNA was searched in contigs with NOVOPlasty v3.8.3 (Dierckxsens et al., 2017) using the mitochondrial sequence of *O. novo-ulmi* as a bait (CM001753.1) and matching contigs were removed. Reads were remapped to the nuclear assembly with bwa mem v0.7.17 (Li & Durbin, 2009), and contigs with normalized mean coverage < 5% and shorter than 1000 bp were also removed.

A total of 11 *Ophiostoma* strains were selected for long-read sequencing with PacBio (Supplementary Table 2). DNA extraction was done in the same way as for Illumina libraries, except that no vortexing and shaking were done to avoid DNA fragmentation. Library preparation and sequencing with the PacBio SMRTcell Sequel system were conducted at the G  nome Qu  bec Innovation Center (Montr  al, Canada). The average depth of coverage of the generated data was between 47x and 269x. *De novo* genome assemblies were generated with pb-falcon v2.2.0 (Chin et al., 2016). The range of parameters was tested and the final configuration files with the best-performing parameters for each species are on https://github.com/aniafijarczyk/Fijarczyk_et_al_2022. *O. quercus* was sequenced in two runs and read data from the two runs were combined together. *O. novo-ulmi* (H294) was sequenced in two runs but only read data from one run was used for an assembly due to sufficient coverage. Assemblies were ordered according to *O. novo-ulmi* H327 genome (GCA_000317715.1) using Mauve snapshot-2015-02-13 (Darling et al., 2004). Assemblies were polished between two to four times by remapping long reads with minimap2 (pbmm v1 (Li, 2018)) and correcting assemblies using arrow (pbgcpp v1). We

also performed one round of polishing with pilon v1.23 (Walker et al., 2014) after mapping short-read Illumina reads (bwa v0.7.17 (Li & Durbin, 2009)) from this study (three assemblies) or from the previous study (Hessenauer et al., 2020) (seven assemblies, Supplementary Table 2). The only exception is *O. quercus*, for which we had no corresponding short-read data, therefore we performed four rounds of correction using arrow. The effectiveness of polishing was assessed by analyzing the completeness of genes with Busco v3 (Waterhouse et al., 2018). Mitochondrial genomes were assembled using Illumina assemblies as baits (or those of related species). Long reads were mapped to bait assembly with pbmm2 (Li, 2018), and a subsample of mapped reads was used for mtDNA assembly with mecat v2 (Xiao et al., 2017). Consecutive rounds of mapping and assembly were conducted until circularized assemblies were obtained. Nuclear assembly contigs with more than 50% of low-quality bases, those mapping to mtDNA, or with no mapping of Illumina reads (standardized mean coverage < 5%) were filtered out. The genome assembly of *O. montium* was very fragmented and had a high percentage of missing conserved genes (74.3%) so instead, an Illumina *de novo* assembly for this species was considered in further analysis.

In both short-read and long-read assemblies, contigs were searched for viral or bacterial contaminants by BLAST searches against bacterial or viral UniProt accessions and separately against fungal UniProt accessions. Contigs having more bacterial/viral hits in length than fungal hits were marked as contaminants and removed. Finally, repeats were identified with RepeatModeller v2.0.1 (Flynn et al., 2020) using option -LTRStruct, and recovered repeat families together with fungal repeats from RepBase were used for masking assemblies with RepeatMasker v4.1.0 (Tarailo-Graovac & Chen, 2009).

Genes were annotated with Augustus v3.3.2 (Stanke & Morgenstern, 2005) and Breaker v2.1.2 (Hoff et al., 2019). To obtain *ab initio* gene models in *Ophiostoma novo-ulmi*, *ulmi*, *himal-ulmi*, *quercus* and *triangulosporum*, a training file was generated using RNA-seq reads from *O. novo-ulmi* H327 (SRR1574322, SRR1574324, SRR2140676) mapped onto a *O. novo-ulmi* H294 genome with STAR v2.7.2b (Dobin et al., 2013). A training set of genes was generated with GeneMark-ES-ET v4.33 (Lomsadze et al., 2005). Final gene models were merged from *ab initio* models, models inferred from the alignment of RNA-seq reads (except *O. triangulosporum*) and *O. novo-ulmi* H327 proteins. The rest of the genome assemblies were annotated only *ab initio* using Magnaporthe grisea species training file.

Sordariomycetes genomes

580 Sordariomycetes genome assemblies (one reference genome per species) were downloaded from NCBI (21/05/2021), the two assemblies were downloaded from JGI Mycocosm (*S. kochii*, *T. guianense*), and two from other resources (*O. ulmi* W9 (Christendat, 2013), *L. longiclavatum* (Wong et al., 2020)). All assemblies were filtered

for short contigs (< 1000 bp) and contaminants as described above. Gene completion was assessed with Busco v3 (Waterhouse et al., 2018) using an orthologous gene set from Sordariomycetes. *Ab initio* gene models were obtained with Augustus v3.3.2 (Stanke & Morgenstern, 2005) with different species-specific training files: *Magnaporthe grisea*, *Fusarium*, *Neurospora*, *Verticillium longisporum*1 or *Botrytis cinerea* (Supplementary Table 1). 31 assemblies were removed based on the low percentage of complete conserved genes (Busco score < 85%) and one due to an excessive number of gene models (*Botrytis cinerea*).

Phylogeny

The maximum likelihood tree was built from 1000 concatenated genes (retrieved with Busco v3 (Waterhouse et al., 2018)) with the highest species representation. Protein alignments were generated with mafft v7.453 (Katoh & Standley, 2013) with E-INS-i method (option '--genafpair --ep 0 --maxiterate 1000'), trimmed with trimal v1.4.rev22 (Capella-Gutiérrez et al., 2009) with option '-automated1' and converted into a matrix. Best protein evolution model JTT+I+G4 was chosen as the most frequent model across all protein alignments, based on the BIC score in IQ-TREE v1.6.12 (Kalyaanamoorthy et al., 2017; Nguyen et al., 2015). The maximum likelihood tree was inferred with ultrafast bootstrap (Hoang et al., 2018), seed number 17629, and 1000 replicates implemented in IQ-TREE. Time-scaled phylogeny was inferred with program r8s v1.81 (Sanderson, 2003), setting the calibration point of 201 My at the split of *Neurospora crassa* and *Diaporthe ampelina*, retrieved from TimeTree (Hedges & Kumar, 2005).

To confirm the branching order of Hypocreales, Microascales, and Glomerellales, another maximum likelihood tree was built, based on a subset of 250 protein alignments (250 longest sequences out of 1000), with a separate model partition assigned to each protein (Chernomor et al., 2016). The maximum likelihood tree was inferred with ultrafast bootstrap, seed number 17629, and 1000 replicates implemented in IQ-TREE v1.6.12 (Hoang et al., 2018; Nguyen et al., 2015).

In the third method for phylogeny inference, we used the multispecies coalescent-based method to obtain a consensus topology based on 1000 separate protein-based trees. The maximum likelihood tree was generated for each of 1000 protein alignments, with their corresponding protein evolution model, ultrafast bootstrap, and 1000 replicates using IQ-TREE v1.6.12 (Hoang et al., 2018; Nguyen et al., 2015). Tree nodes with support of less than 30 were contracted using newick-utils v1.6 (Junier & Zdobnov, 2010). Consensus topology was obtained with Astral v5.7.8 (Zhang et al., 2018).

Genomic and ecological traits

Pathogenicity and insect association were assigned based on a literature search. Pathogenicity was assigned if the species caused a well-recognized disease, or

pathogenicity was experimentally documented at least on one host species. We considered pathogens of plants, animals, and fungi, both obligatory and opportunistic. Insect association included all types of relationships with insects, including pathogenic, symbiotic, and mutualistic. Insect-vectored species were limited to documented cases of insect transmission. Analyzed genomic traits included genome size (bp), number of genes, a fraction of repeat content, size of the assembly excluding repeat content (bp), GC content, the mean number of introns per gene, mean intron size (bp), mean exon size (bp), a fraction of genes with introns, mean intergenic length (bp) and number of tRNA and pseudo tRNA genes. Genome size is equivalent to assembly size after filtering. Genome size excluding repeat content is assembly size excluding regions masked by RepeatMasker (in this study or downloaded from NCBI), and repeat content is the proportion of masked bases compared to total assembly size. Gene number corresponds to the number of all *ab initio* gene models obtained with Augustus. GC content is the proportion of GC bases in the total assembly. The mean number of introns per gene, intron and exon size, a fraction of genes with introns and intergenic length were estimated from gff files with annotated gene models. tRNA and pseudo tRNA genes were obtained with tRNAscan-SE v2.0.9 (Chan et al., 2021).

Trait correlations

Correlations between genomic traits were calculated using phylogenetic generalized least squares (gls) function from the R package nlme and independent contrasts using pic function from the R package ape (Paradis et al., 2004). P-values in phylogenetic gls models were obtained using anova() function, and adjusted for multiple tests with Benjamini-Hochberg method at false discovery rate of 0.05. Some genomic traits were rescaled: genome (bp $\times 10^{-7}$), genes ($\times 10^{-4}$), assembly w/o repeats (bp $\times 10^{-7}$), intron length (bp $\times 10^{-4}$), exon length (bp $\times 10^{-3}$), intergenic length (bp $\times 10^{-4}$), tRNAs ($\times 10^{-3}$), and pseudo tRNAs ($\times 10^{-3}$).

Coevolution of ecological traits (pathogenicity and insect association) with genomic traits was estimated using three approaches, i) a reversible jump MCMC discrete model of evolution implemented in BayesTraits v3 (Pagel & Meade, 2006), ii) a phylogenetic logistic regression with phyloglm function in R package phylolm (Ho & Ané, 2014), iii) an investigation of feature's importances from machine learning classification using scikit-learn python (v3) package. All analyses were additionally run for the subset of species (n=563), in which any repeating species were excluded.

For analysis in BayesTraits, genomic traits were converted into binary traits based on their median (0 if below median and 1 if above median). Four different models were investigated: independent evolution, dependent evolution, covarion model where dependent evolution varies across the tree (Venditti et al., 2011), and for a subset of significant genomic traits additional model testing equal transition rates. To determine

the direction of coevolution with ecological traits (positive or negative), two models were compared, one in which all transition rates were allowed to differ (dependent model), and the second one in which the two most frequent transition rates were set to be equal. If the model with differing transition rates performed better, the direction of coevolution was determined by the most frequent transition rate, otherwise, no specific direction was concluded. Models were tested by comparing complex to simpler models with log bayes factor, i.e. dependent vs. independent, covarion vs. dependent, dependent vs. dependent with selected equal transition rates. A complex model was chosen if the log bayes factor surpassed 5. Each model was run three to five times to check for consistency with 21 mln iterations, 1 mln burn-in, and thinning of 1000. Analyses for the subset of species (n=563), in which repeating species were excluded, were run in three independent runs for the dependent, independent model, and for the dependent model with constrained rates.

In the second approach, a phylogenetic logistic regression was used to fit each genomic trait to ecological trait using phyloglm function in R with the “logistic_MPLE” method, btol option (searching space limit) set to 30, and 1000 independent bootstrap replicates. Benjamini-Hochberg correction was applied to p-values. Species with missing information on the ecological trait were filtered out.

In the last approach, genomic traits (features) were used to train several machine learning classifiers for the prediction of ecological traits. Species with missing data were filtered out, and data were rescaled (between 0 and 1) for classification with SVC. Twenty percent of data was selected as held-out data. First, the performance of several classifiers was investigated (KNeighbors, SVC, DecisionTree, RandomForest, and GradientBoosting) with a balanced accuracy score, and 5-fold cross-validation, in which train set (80%) was split into 5 parts (with two classes of a predicted trait in same proportions), and accuracy measured 5 times, each time a different part being set as a test set and the remaining four parts as the training set. Hyperparameters for the best-performing classifier were chosen using grid search. Performance on held-out data was evaluated with a balanced accuracy score, precision, recall, and ROC curve. To account for the role of phylogeny, a matrix of internode distances was used in combination with other genomic features. The classification models were evaluated in the same way and compared with those not integrating phylogenetic distances. Feature importances were obtained with a mean decrease impurity method from the RandomForest classifier.

For the trait of pathogenicity, the RandomForest classifier performed best with an accuracy of 0.7. Phylogeny had a negligible effect on model performance, improving it only to 0.71. Top features included all genomic features, and two phylogenetic nodes

corresponding to two overlapping clades with *Claviceps* species excluding *C. paspali*, and *C. citrina*. The same model run on a filtered dataset (n=563), excluding repeating species, gave an accuracy below 0.7, therefore we did not report resulting feature importances.

For the trait of insect association, the RandomForest classifier also gave the best performance, giving an accuracy of 0.83, and 0.87 when the phylogenetic matrix was included. Both approaches revealed exon length and genome size without repeats as the two most important features, and intron length or number of introns as a third feature, respectively. The model ran on a filtered dataset (n=563), excluding repeating species, gave an accuracy of 0.82% and 91% for the approaches without and with the phylogenetic matrix, respectively. The most important features of the first approach were: exon length, assembly w/o repeats, intergenic length, and intron length, and of the second approach: exon length, assembly without repeats, number of genes, and number of introns.

To compare the ancestral with the current states of genomic traits at focal nodes in the tree, ancestral states were estimated with a continuous model for genomic traits implemented in BayesTraits v3 (Organ et al., 2007) using the same MCMC parameters as in the case of discrete traits models. Similarly, ecological traits were reconstructed with a discrete model in BayesTraits v3 (Pagel & Meade, 2006). To visualize evolution of traits on the tree, ancestral states of genomic traits were reconstructed with “ER” model implemented in ace function from R package phytools (Revell, 2011), and discrete ecological traits were reconstructed with fastAnc function from R package phytools (Revell, 2011). Traits were visualized on the tree with ggtree R package (G. Yu, 2020) using the continuous option.

Functional gene classes

Functional gene annotations were determined using several databases. Protein sequences were searched against KOG (release 2003 (Tatusov et al., 2003)) and MEROPS Scan Sequences (release 12.1 (Rawlings et al., 2018)) using diamond v2.0.9 (Buchfink et al., 2015) with options ‘--more-sensitive -e 1e-10’ and against CAZymes (download 24/09/2021 (Cantarel et al., 2009)) with options ‘--more-sensitive -e 1e-102’. Pfam domains were searched with pfam_scan.pl (Mistry et al., 2021) script with hmmer v3.2.1 (Mistry et al., 2013) and matched with Sordariomycetes transcription factors obtained from JGI Mycocosm database (download 9/12/2021). Secondary Metabolite Clusters were inferred from genome assemblies using antiSMASH v4.0.2 (Blin et al., 2019). Ancestral states of the gene and gene cluster numbers were inferred using the fastAnc function from R package phytools (Revell, 2011). Per clade change ratio was calculated as a mean difference between the observed number of genes in each clade member and ancestral state divided by ancestral state. Mean gains or losses of genes

per clade were tested with a bootstrap of 10 species with 100 replicates for clades with ≥ 10 members.

Orthologs and orthogroups

A subset of 112 species was selected, by randomly selecting 1-8 species per clade. Orthologs were searched using Orthofinder v2.5.2 (Emms & Kelly, 2019) with no *a priori* defined tree. Orthogroup expansions and contractions were detected with CAFE v4.2.1 (Han et al., 2013) with a p-value threshold of 0.05. Genes belonging to one-to-one orthologs and orthogroups were determined from the Orthofinder output, and the presence of orthogroups in the ancestral clades was determined from CAFE output.

Limitations of the study

The major challenge of this study is the assignment of lifestyle traits. Two reasons account for this. First, information on the given lifestyle is limited, and second, species can potentially exert a given phenotype only in specific conditions, but that has not been reported or tested. These reasons may account for some inevitable uncertainty in both traits: pathogenicity and insect association. The trait of insect association is typically characteristic of individual clades, therefore most members of the clade are expected to carry this trait or to be closely related to species with that trait. Therefore we can alleviate the misannotation by comparing insect-associated clades with their corresponding non-insect-associated sister clades, as was done in the fourth part of the Results. Pathogenicity on the other hand is a fast-evolving trait, exerting a large variety of effects on the host. Nevertheless, the inferences of coevolution between pathogenicity and genomic traits are generally consistent with observations in clades enriched for well-studied pathogenic species. Because of the lack of ecological information for many sequenced fungal genomes, many details on the fungal lifestyles could not be included. One trait that is known to impact genome evolution is the host range of pathogens (Badet et al., 2017). The availability of the host, as well as the ability to switch from pathogenic to saprotrophic lifestyle, can potentially have a great impact on Ne, and therefore genome evolution.

Gene annotations and genome processing was performed in a uniform way for all the species, in order to remove any software-related biases. Repeat content was obtained from NCBI repeat annotations together with downloaded assemblies, therefore the quality of annotations is expected to vary mostly with the assembly quality. Short-read assemblies which dominate in the dataset could cause an underestimation of the true amount of repeats, however, the cases of repeat expansions in our dataset can be detected in both contiguous and fragmented assemblies. Although the annotations are mostly underestimated, our dataset allows us to detect general trends in both genes and repeat dynamics.

Several more detailed aspects of genome architecture could be informative about pathogen evolution. For instance, we did not consider the level of compartmentalization of the genome, which plays an important role in the evolution of virulence genes (Möller & Stukenbrock, 2017). Accessory chromosomes or chromosome duplication are also important factors in generating variation and aiding in adaptation to the host. In fact, we know little about the frequency of these processes in non-pathogenic species. Similarly, the distribution of specific genes, such as tRNAs or repeats in the genome can inform us about the architecture of pathogenic genomes. Finally, whole genome duplications are processes that can lead to major lifestyle transitions and at the same time can influence the level of gene duplication and loss leading to major genome rearrangements.

Data Availability

Raw short and long reads from sequenced genomes are available at NCBI under project number PRJNA841745. Genome assemblies have been deposited at GenBank under the accessions JANS LN000000000-JANS MH000000000. The versions described in this paper are versions JANS LN010000000-JANS MH010000000. Code used in this study is available on github (https://github.com/aniafijarczyk/Fijarczyk_et_al_2022). Cleaned genome assemblies, sequences and coordinates of gene models were uploaded to Dryad.

Acknowledgments

We thank Erika Dort for an access to lifestyle database to confirm pathogenicity traits of 90 species from our dataset. We thank Rohan Dandage, Ilga Porth and Louis Bernier for comments and discussions about the manuscript. This project was funded by Genome Canada and Genome Québec BioSAFE project and a NSERC Discovery grant to CRL. CRL holds the Canada Research Chair in Cellular Systems and Synthetic Biology.

Figure Supplements

Figure 1–figure supplement 1. Correlations among 12 genomic traits. Squares within the matrix show scatterplots between pairs of genomic traits. K on the scale stands for x1000.

Figure 1–figure supplement 2. Principal component analysis on 12 genomic traits. The inset shows a biplot from the principal component analysis performed on contrasts of 12 genomic traits. The color of the variables indicates the contribution of each variable to the variance explained by principal components.

Figure 2–figure supplement 1. Posterior probabilities of transition rates estimated with BayesTraits from one run of a dependent model of the evolution of pathogenicity with each genomic trait. P - pathogen, NP - non-pathogen, low - genomic trait value below median, high - genomic trait value above the median. Asterisks indicate traits for which two dominant rates are not equal.

Figure 2–figure supplement 2. Results of the pathogen classification with Random Forest. A. ROC curves (where the positive label is the pathogen), and performance metrics of the best classifier (random forest) trained to distinguish pathogenic from non-pathogenic species. C. feature importances (and credible intervals) from most to least important for all genomic traits. D and E. Same as B and C, except that genomic traits were combined with internode phylogenetic distances.

Figure 2–figure supplement 3. Coevolution of pathogenicity and genomic traits in 10 subsets, with the pathogenic species randomly selected to match the number of non-pathogenic species. A. Cells show average log bayes factors across 3 runs equal to 4 or more. Log bayes factors compared models of dependency and non-independency of genomic traits on pathogenicity. B. Cells show a coefficient estimate obtained from the phylogenetic logistic regression run with phyloglm. Values with p-value > 0.05 are shown.

Figure 3–figure supplement 1. Ancestral states of 12 genomic traits mapped on the phylogeny. A. Assembly length in bp (genome). B. Number of genes (genes). C. Assembly size excluding repeats in bp (genome w/o repeats). D. Repeat content measured as a fraction of the whole genome (repeats). E. Proportion of GC bases (GC). F. Mean intron length in bp (intron length). G. Fraction of genes with introns (genes with introns). H. Mean number of introns in a gene (introns). I. Mean intergenic length in bp (intergenic length). J. Mean exon length in bp (exon length). K. Number of tRNA genes (tRNA). L. Number of pseudo tRNA genes (pseudo tRNA). Ancestral states were inferred using the fastAnc function in R package phytools.

Figure 4–figure supplement 1. Comparison of insect and non-insect associated clades in 5 groups (a-e). Comparison of genomic traits in current members of the clades, between insect-associated clades (blue) and non-insect-associated clades (yellow). Numbers in parentheses near the clade name on the tree indicate clade abundance. Stars show significant pairwise differences between blue and brown clades within each group (Wilcoxon rank-sum test, adjusted p<0.05). Clades in group C have too few species for testing.

Figure 4–figure supplement 2. Fold change of genomic traits in insect and non-insect associated clades compared to the ancestral state. Insect-associated clades are shown in blue, non-insect associated clades in yellow, and ancestral clades are indicated with a white dot on the tree. Shown are only genomic traits with narrow credible intervals for ancestral nodes. Clades O (group a), M1 (group b), H2.8 (group c), and H2.2 (group e) are insect mutualists or

symbionts, whereas clades H2.6 (group d) and clades H2.4 and H2.3 (group e) are insect pathogens.

Figure 5—figure supplement 1. Posterior probabilities of transition rates estimated with BayesTraits from one run of a dependent model of the evolution of insect association with 12 genomic traits. Asterisks indicate traits for which two dominant rates are not equal. I - insect association, NI - no insect association, low - genomic trait value below the median, high - genomic trait value above the median.

Figure 5—figure supplement 2. Performance of the Random Forest classifier trained to distinguish insect-associated from non-insect-associated species. A. ROC curves (where the positive label is the insect-associated) and performance metrics. B. Top feature importances (and credible intervals) according to random forest classifier. C and D are the same as A and B, but genomic traits were combined with internode phylogenetic distances.

Figure 5—figure supplement 3. Losses and gains of orthogroups estimated for 112 species with CAFE. The time scale is in millions of years. Blue stripes cluster insect-associated clades, and yellow stripes cluster other, non-insect-associated clades. The heatmap shows the cumulative sum of all losses and gains between the root and each leaf (species) on the tree.

Figure 6—figure supplement 1. Insect-associated pathogens lose genes involved in breaking host barriers. Heatmap shows the change ratio of genes/clusters relative to the ancestral state. Clades are shown in columns with the number of clade members in parentheses, functional classes are shown in rows. Dots indicate significant gain (red) or loss (blue) of genes/clusters across clade members estimated from 100 rounds of bootstrapping of 10 species in clades with ≥ 10 members. SMC - secondary metabolite clusters, M - Merops, TF - transcription factors.

Figure 7—figure supplement 1. Gene structure changes in orthologs and orthogroups of clades M1 and M2. A. Comparison of gene structures between one-to-one single-copy orthologs ($n=583$) from insect-associated (IA) clade M1 (blue) and the corresponding non-insect-associated (non-IA) clade M2 (yellow). Intron length was compared only between orthologs with at least one intron ($n=476$). Orthologue features were averaged across five species within a clade. Boxplots show medians, first and third quartiles, and lines span minimum to maximum values excluding outliers. p - p-values estimated with the paired two-sided Wilcoxon signed-rank test; d - mean differences between IA and non-IA orthologues. B. Average exon length and the number of exons of gene families (orthogroups) present in different frequencies across up to five clade members.

Source Data

Figure 1-Source Data 1. Transformed genomic traits for all species used to calculate correlations.

Figure 1-Source Data 2. Genomic traits for all samples.

Figure 1-Source Data 3. Contrasts of genomic traits.

Figure 2-Source Data 1. Phylogeny and ancestral states of pathogenicity.

Figure 2-Source Data 2. Pathogenicity, raw genomic features, binary genomic features and inter-node distances.

Figure 2-Source Data 3. Genomic traits transition rates estimated with BayesTraits for coevolution with pathogenicity.

Figure 2-Source Data 4. Pathogenicity, raw and binary genomic features for 10 subsampled datasets.

Figure 3-Source Data 1. Phylogeny and ancestral gene numbers.

Figure 3-Source Data 2. List of phylogenies and ancestral states of all genomic features.

Figure 4-Source Data 1. Genomic traits grouped by clade.

Figure 4-Source Data 2. Genomic traits and per-clade fold change estimates.

Figure 5-Source Data 1. Insect-association, raw genomic features, binary genomic features and inter-node distances.

Figure 5-Source Data 2. Genomic traits transition rates estimated with BayesTraits for coevolution with insect-association.

Figure 5-Source Data 3. Phylogeny with orthogroup contractions and expansions.

Figure 6-Source Data 1. Fold change estimates of gene counts since clade ancestors.

Figure 6-Source Data 2. Fold change estimates of gene counts since clade ancestors in pathogenic species only.

Figure 7-Source Data 1. Intron and exon features in one-to-one orthologs.

Figure 7-Source Data 2. Per-orthogroup mean intron and exon features and their clade occupancy.

Supplementary Files

Supplementary Table 1. Information on Sordariomycetes genome assemblies including assembly identifiers, quality statistics and estimated trait values.

Supplementary Table 2. Sequencing information for species sequenced in this study.

Supplementary Table 3. Log marginal likelihoods and bayes factors of compared models of coevolution of pathogenicity with genomic traits. P - pathogen, NP - non-pathogen, low - value below median, high - value above median over all species. In the column Dataset, "Complete" includes all species, and "Filtered" excludes samples from the same species. Bold values indicate log bayes factor ≥ 4 .

Supplementary Table 4. Results of fitting genomic traits to pathogenicity with phyloglm model. CI - confidence interval. In the column Dataset, "Complete" includes all species, and "Filtered" excludes samples from the same species. Bold values in columns "Coefficient" highlight positive coefficients, and in column "Adjusted p-value", p-values < 0.05 .

Supplementary Table 5. Discrete trait estimates for the ancestral nodes in

Sordariomycetes. Node assignment corresponds to the one in Figure 3. Trait = 0 - Absence of a trait (non-pathogenic, non -insect associated). Trait = 1 - Presence of a trait (pathogenic, insect associated).

Supplementary Table 6. Log marginal likelihoods and log bayes factors of compared models of coevolution of insect association with genomic traits. I - insect-associated, NI - non-insect-associated, low - value below median, high - value above median over all species. In the column Dataset, "Complete" includes all species, and "Filtered" excludes samples from the same species. Bold values indicate log bayes factor ≥ 4 .

Supplementary Table 7. Results of fitting genomic traits to insect-association trait with phyloglm model. CI - confidence interval. Complete dataset includes all species, and filtered dataset excludes samples from the same species. Bold values in columns "Coefficient" highlight positive coefficients., and in column "Adjusted p-value", p-values < 0.05 .

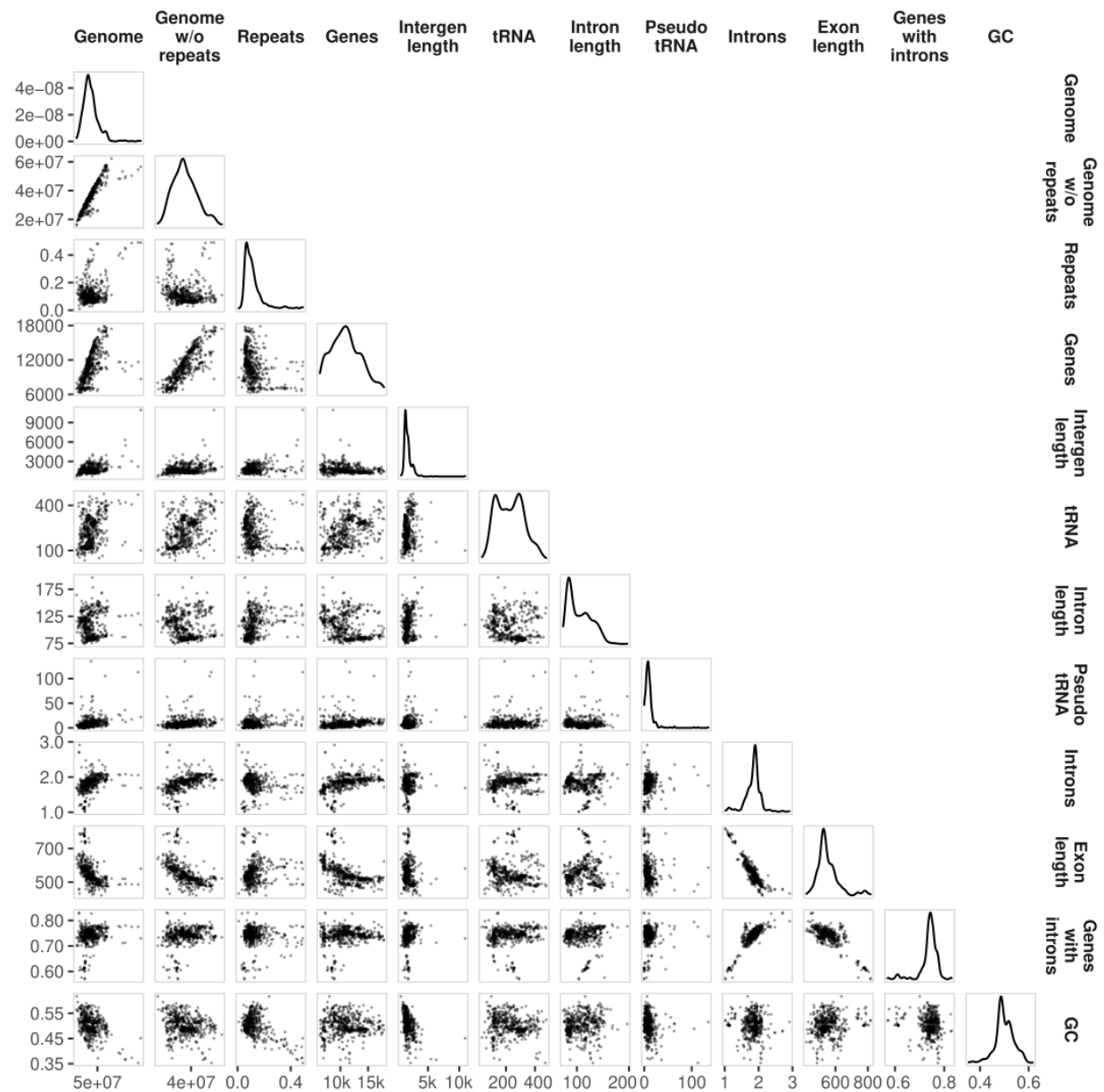


Figure 1—figure supplement 1. Correlations among 12 genomic traits. Squares within the matrix show scatterplots between pairs of genomic traits. K on the scale stands for x1000.

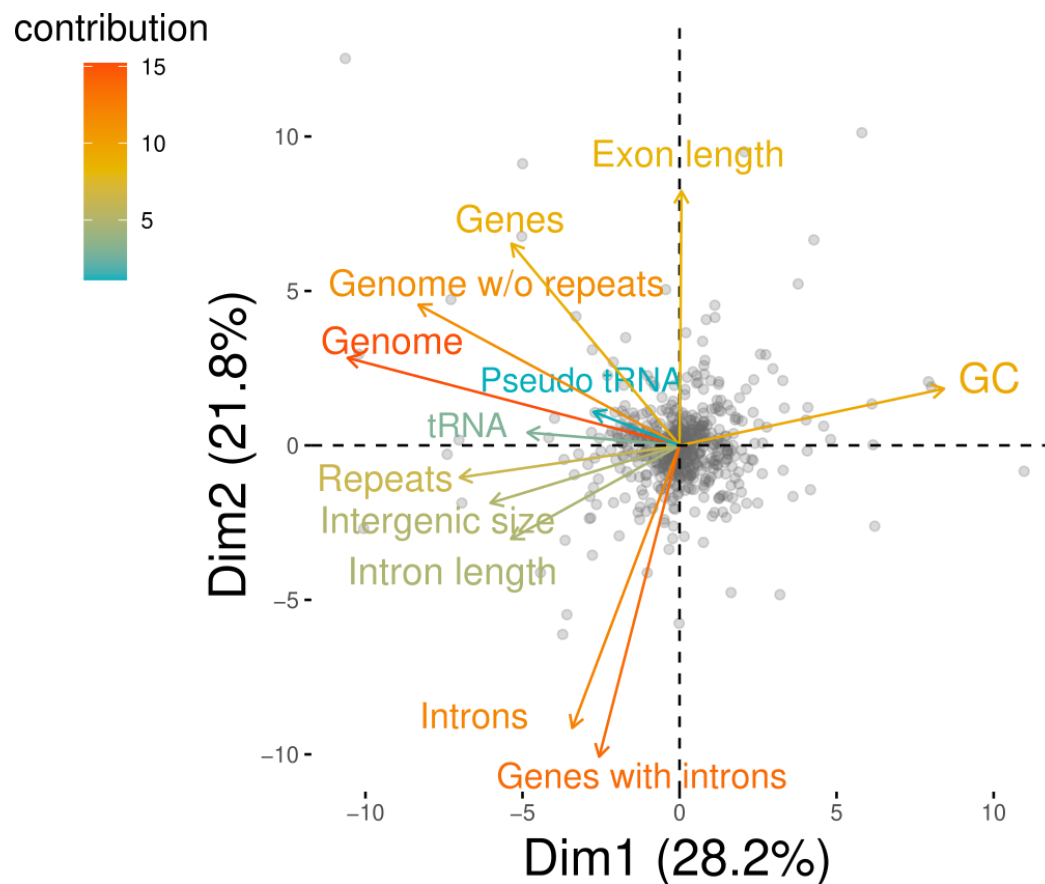


Figure 1—figure supplement 2. Principal component analysis on 12 genomic traits. The inset shows a biplot from the principal component analysis performed on contrasts of 12 genomic traits. The color of the variables indicates the contribution of each variable to the variance explained by principal components.

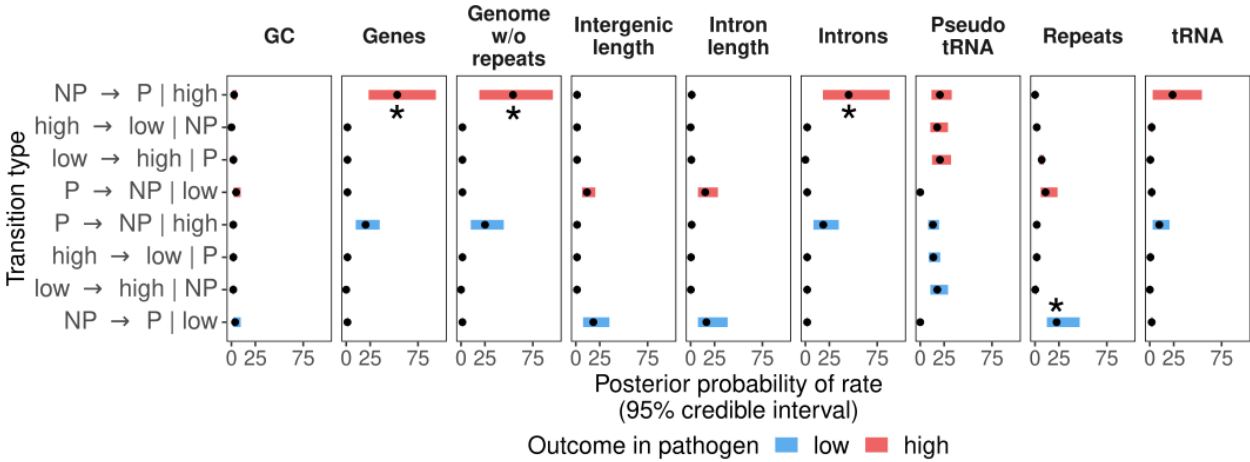


Figure 2–figure supplement 1. Posterior probabilities of transition rates estimated with BayesTraits from one run of a dependent model of the evolution of pathogenicity with each genomic trait. P - pathogen, NP - non-pathogen, low - genomic trait value below median, high - genomic trait value above the median. Asterisks indicate traits for which two dominant rates are not equal.

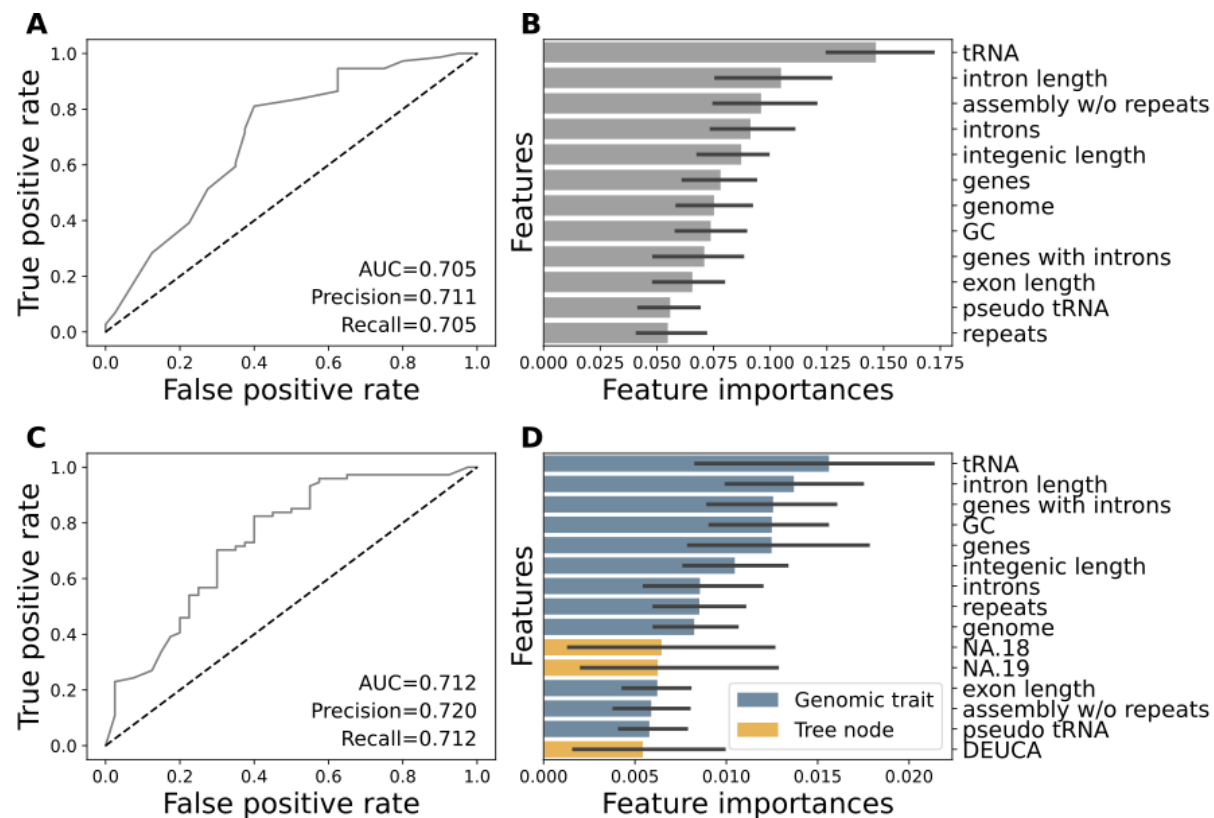


Figure 2—figure supplement 2. Results of the pathogen classification with Random Forest
A. ROC curves (where the positive label is the pathogen), and performance metrics of the best classifier (random forest) trained to distinguish pathogenic from non-pathogenic species. **C.** feature importances (and credible intervals) from most to least important for all genomic traits. **D** and **E.** Same as B and C, except that genomic traits were combined with internode phylogenetic distances.

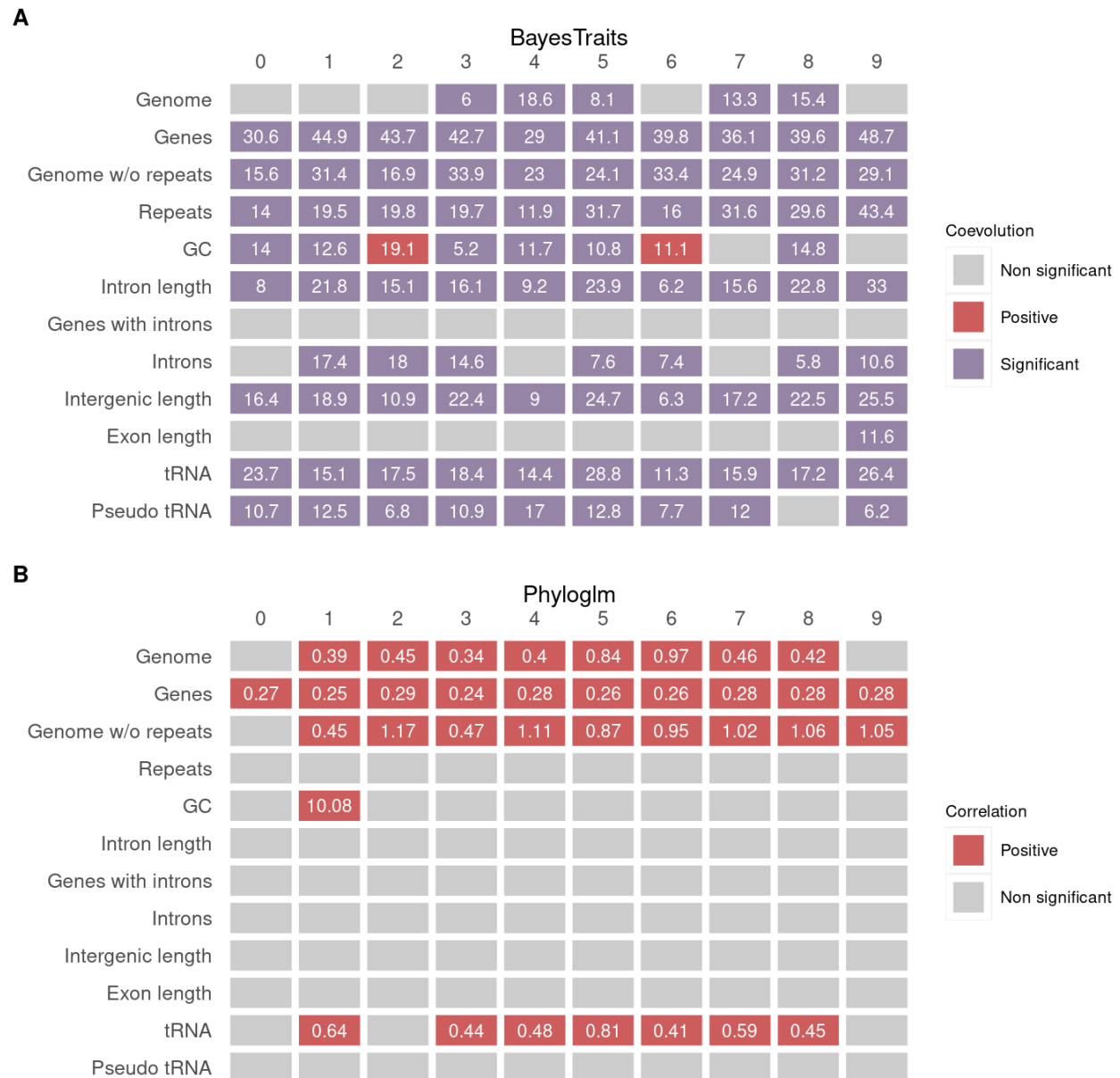


Figure 2—figure supplement 3. Coevolution of pathogenicity and genomic traits in 10 subsets, with the pathogenic species randomly selected to match the number of non-pathogenic species. A. Cells show average log bayes factors across 3 runs equal to 4 or more. Log bayes factors compared models of dependency and non-independency of genomic traits on pathogenicity. **B.** Cells show a coefficient estimate obtained from the phylogenetic logistic regression run with phyloglm. Values with p-value > 0.05 are shown.

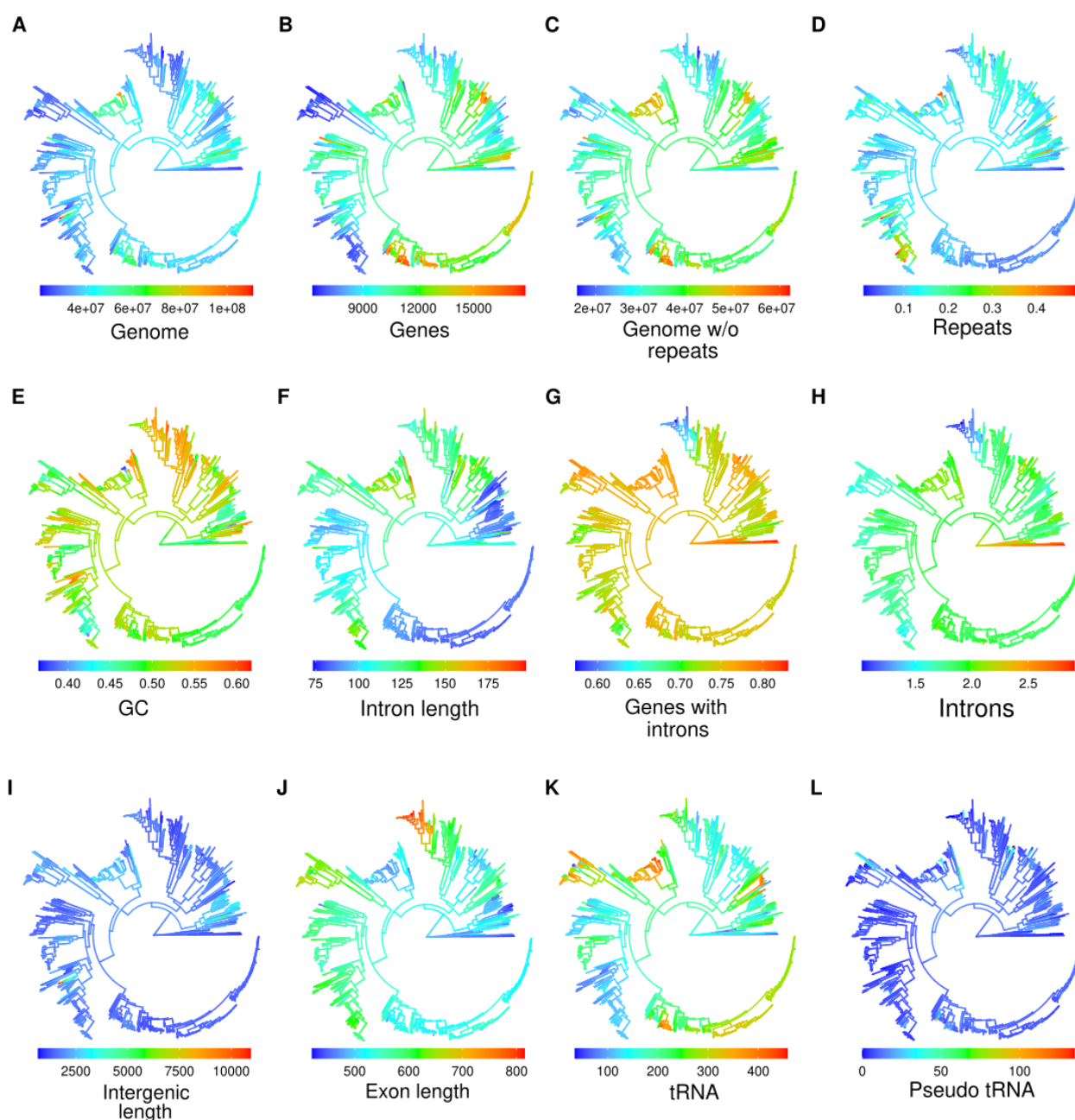


Figure 3—figure supplement 1. Ancestral states of 12 genomic traits mapped on the phylogeny. **A.** Assembly length in bp (genome). **B.** Number of genes (genes). **C.** Assembly size excluding repeats in bp (genome w/o repeats). **D.** Repeat content measured as a fraction of the whole genome (repeats). **E.** Proportion of GC bases (GC). **F.** Mean intron length in bp (intron length). **G.** Fraction of genes with introns (genes with introns). **H.** Mean number of introns in a gene (introns). **I.** Mean intergenic length in bp (intergenic length). **J.** Mean exon length in bp (exon length). **K.** Number of tRNA genes (tRNA). **L.** Number of pseudo tRNA genes (pseudo tRNA). Ancestral states were inferred using the fastAnc function in R package phytools.

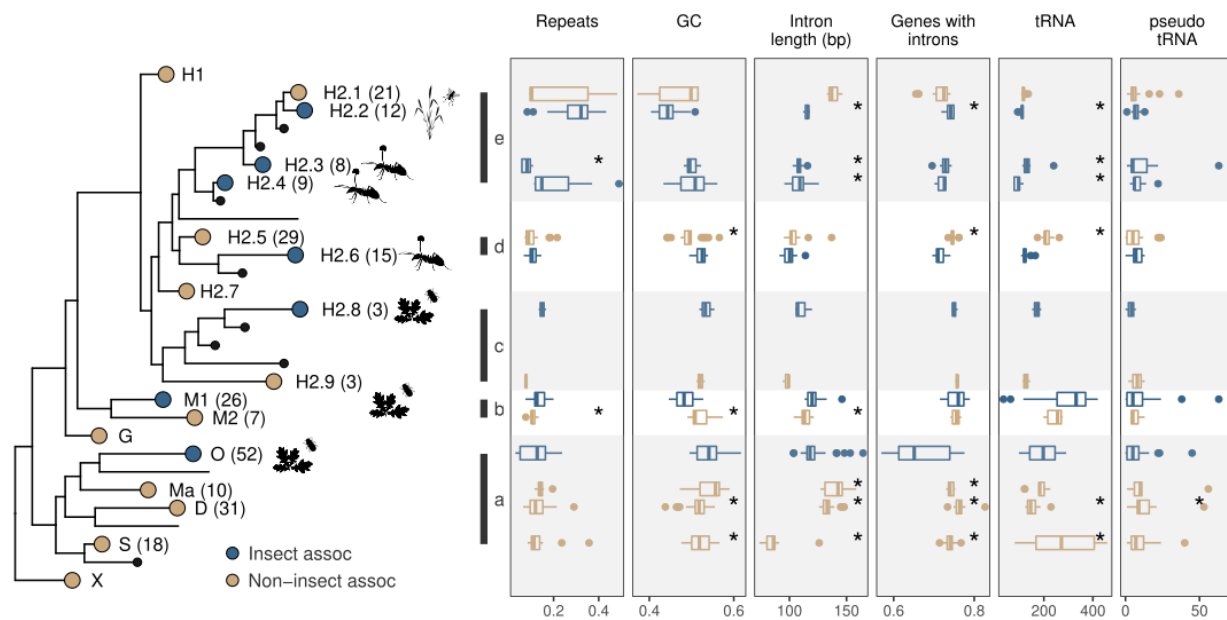


Figure 4-figure supplement 1. Comparison of insect and non-insect associated clades in 5 groups (a-e). Comparison of genomic traits in current members of the clades, between insect-associated clades (blue) and non-insect-associated clades (yellow). Numbers in parentheses near the clade name on the tree indicate clade abundance. Stars show significant pairwise differences between blue and brown clades within each group (Wilcoxon rank-sum test, adjusted $p < 0.05$). Clades in group C have too few species for testing.

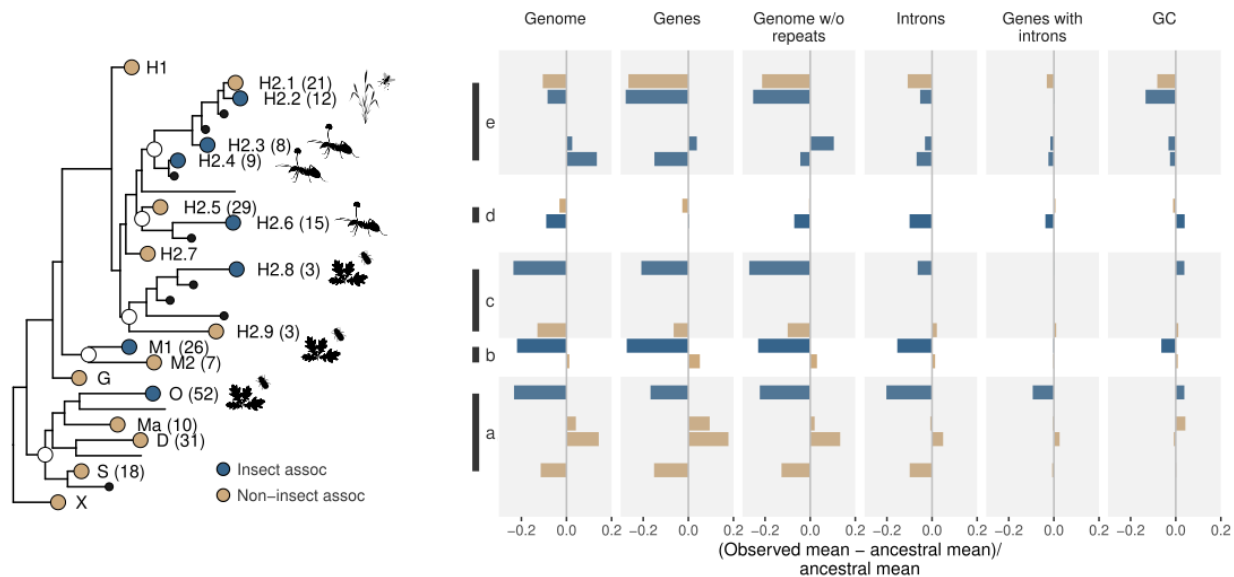


Figure 4—figure supplement 2. Fold change of genomic traits in insect and non-insect associated clades compared to the ancestral state. Insect-associated clades are shown in blue, non-insect associated clades in yellow, and ancestral clades are indicated with a white dot on the tree. Shown are only genomic traits with narrow credible intervals for ancestral nodes. Clades O (group a), M1 (group b), H2.8 (group c), and H2.2 (group e) are insect mutualists or symbionts, whereas clades H2.6 (group d) and clades H2.4 and H2.3 (group e) are insect pathogens.

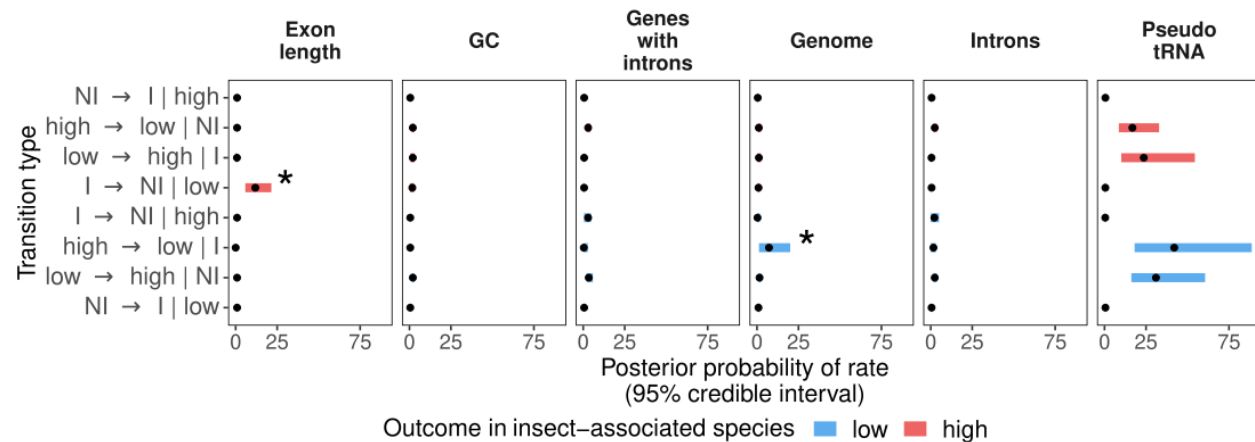


Figure 5–figure supplement 1. Posterior probabilities of transition rates estimated with BayesTraits from one run of a dependent model of the evolution of insect association with 12 genomic traits. Asterisks indicate traits for which two dominant rates are not equal. I - insect association, NI - no insect association, low - genomic trait value below the median, high - genomic trait value above the median.

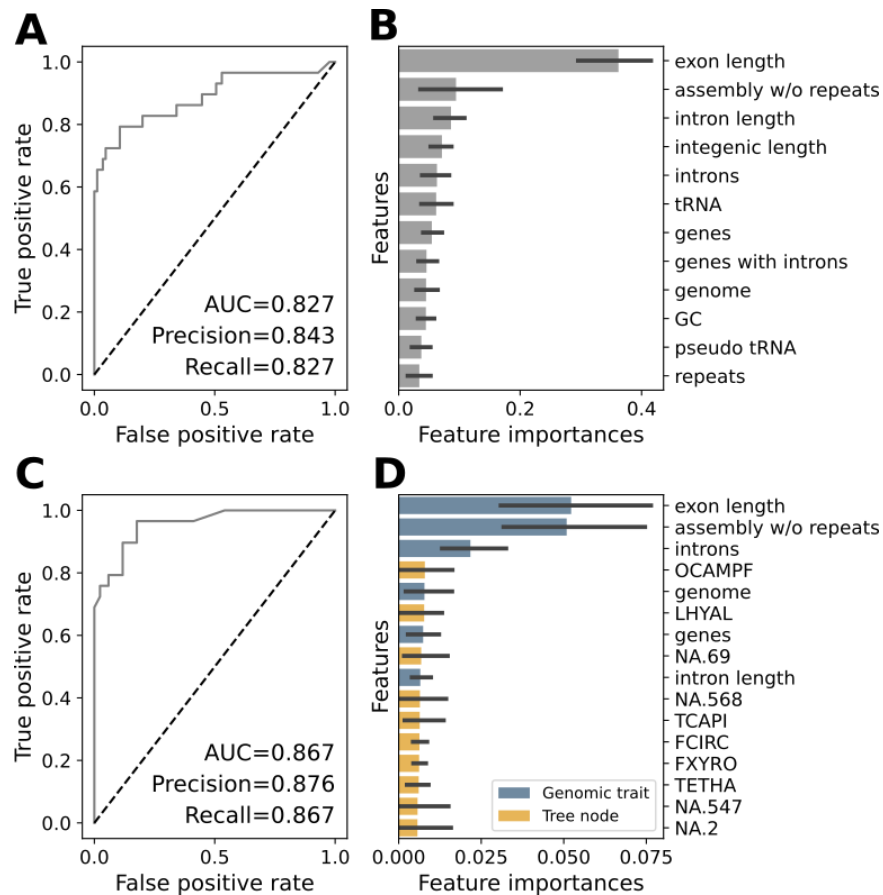


Figure 5—figure supplement 2. Performance of the Random Forest classifier trained to distinguish insect-associated from non-insect-associated species. A. ROC curves (where the positive label is the insect-associated) and performance metrics. **B.** Top feature importances (and credible intervals) according to random forest classifier. **C** and **D** are the same as **A** and **B**, but genomic traits were combined with internode phylogenetic distances.

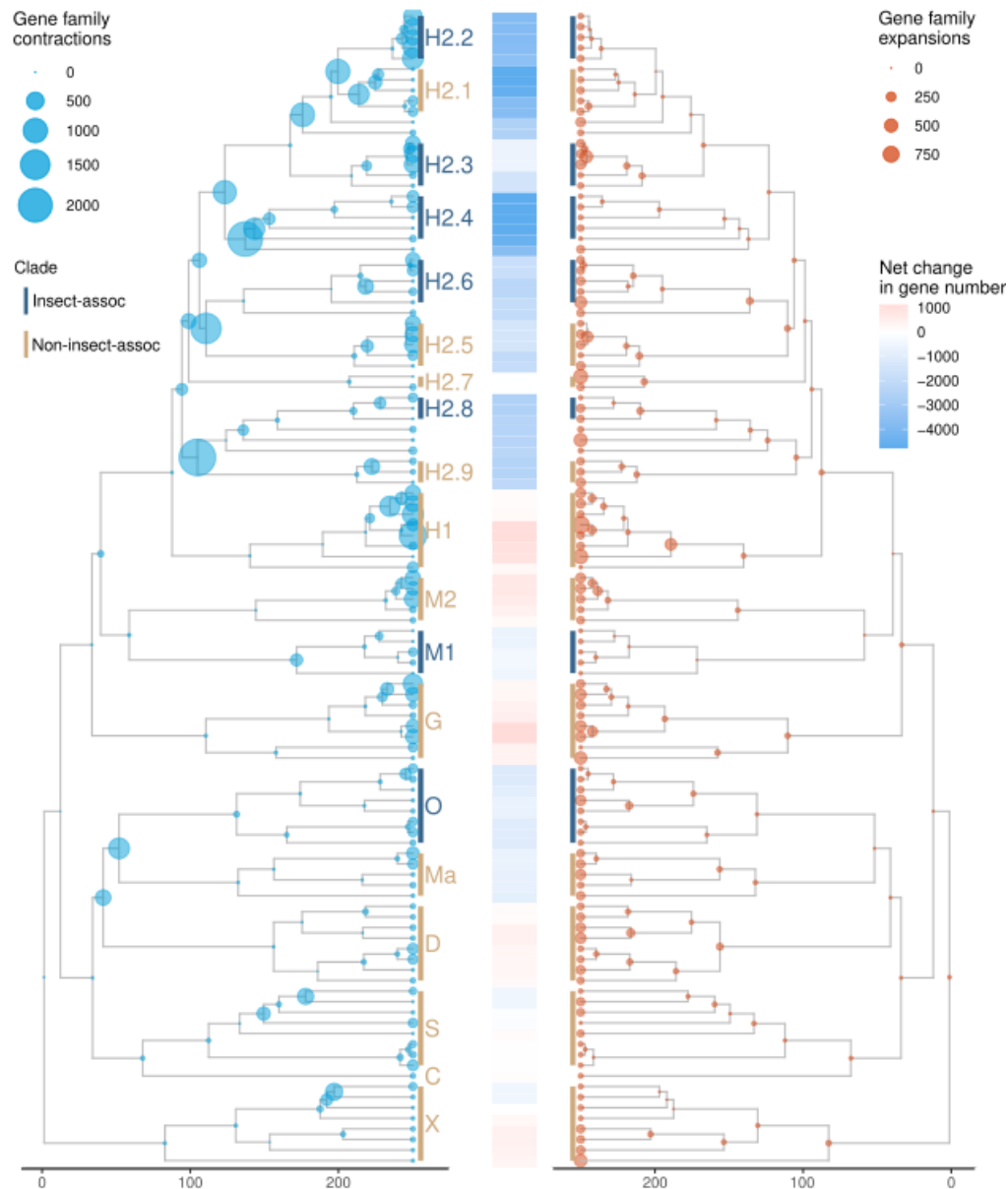


Figure 5—figure supplement 3. Losses and gains of orthogroups estimated for 112 species with CAFE. The time scale is in millions of years. Blue stripes cluster insect-associated clades, and yellow stripes cluster other, non-insect-associated clades. The heatmap shows the cumulative sum of all losses and gains between the root and each leaf (species) on the tree.

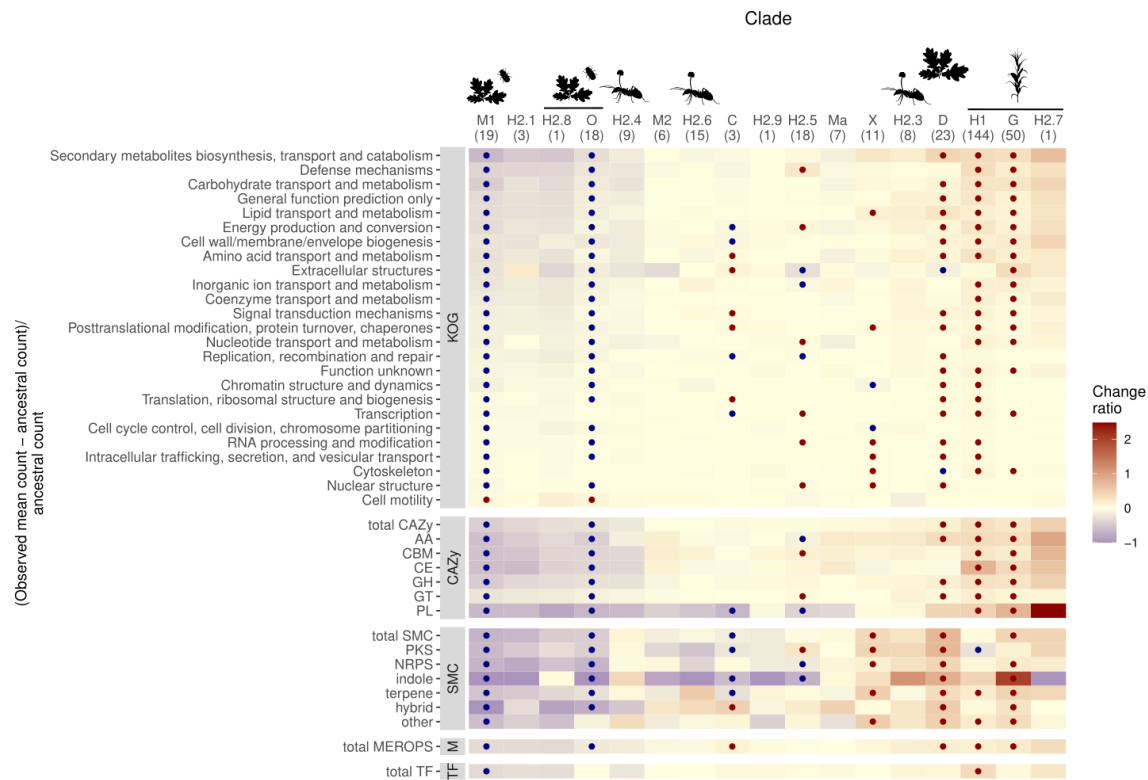


Figure 6—figure supplement 1. Insect-associated pathogens lose genes involved in breaking host barriers. Heatmap shows the change ratio of genes/clusters relative to the ancestral state. Clades are shown in columns with the number of clade members in parentheses, functional classes are shown in rows. Dots indicate significant gain (red) or loss (blue) of genes/clusters across clade members estimated from 100 rounds of bootstrapping of 10 species in clades with ≥ 10 members. SMC - secondary metabolite clusters, M - Merops, TF - transcription factors.

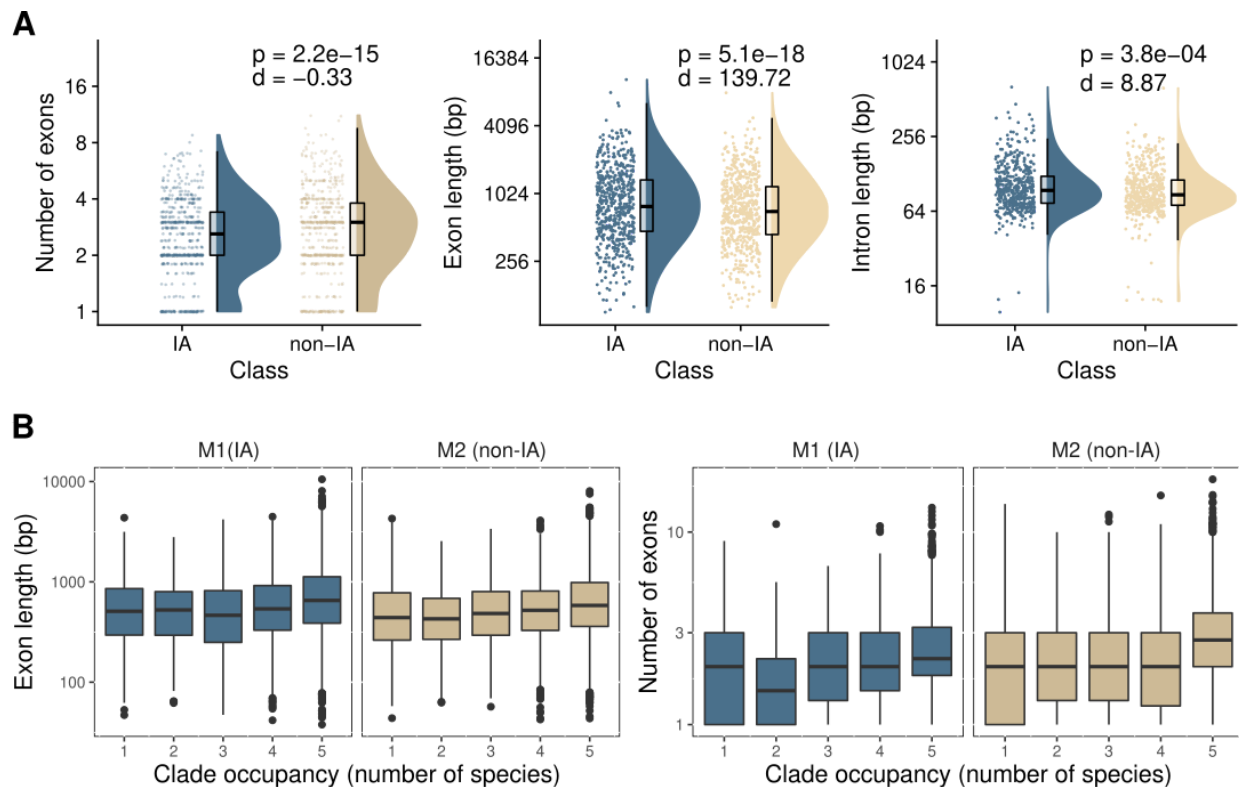


Figure 7—figure supplement 1. Gene structure changes in orthologs and orthogroups of clades M1 and M2. **A.** Comparison of gene structures between one-to-one single-copy orthologs ($n=583$) from insect-associated (IA) clade M1 (blue) and the corresponding non-insect-associated (non-IA) clade M2 (yellow). Intron length was compared only between orthologs with at least one intron ($n=476$). Orthologue features were averaged across five species within a clade. Boxplots show medians, first and third quartiles, and lines span minimum to maximum values excluding outliers. p - p -values estimated with the paired two-sided Wilcoxon signed-rank test; d - mean differences between IA and non-IA orthologues. **B.** Average exon length and the number of exons of gene families (orthogroups) present in different frequencies across up to five clade members.

References

- Andrews, S. (2010). *FastQC: A Quality Control Tool for High Throughput Sequence Data*.
<http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>
- Badet, T., Peyraud, R., Mbengue, M., Navaud, O., Derbyshire, M., Oliver, R. P., Barbacci, A., & Raffaele, S. (2017). Codon optimization underpins generalist parasitism in fungi. *eLife*, 6.
<https://doi.org/10.7554/eLife.22472>
- Bankevich, A., Nurk, S., Antipov, D., Gurevich, A. A., Dvorkin, M., Kulikov, A. S., Lesin, V. M., Nikolenko, S. I., Pham, S., Prjibelski, A. D., Pyshkin, A. V., Sirotkin, A. V., Vyahhi, N., Tesler, G., Alekseyev, M. A., & Pevzner, P. A. (2012). SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing. *Journal of Computational Biology: A Journal of Computational Molecular Cell Biology*, 19(5), 455–477.
- Bao, J., Chen, M., Zhong, Z., Tang, W., Lin, L., Zhang, X., Jiang, H., Zhang, D., Miao, C., Tang, H., Zhang, J., Lu, G., Ming, R., Norvienyeku, J., Wang, B., & Wang, Z. (2017). PacBio Sequencing Reveals Transposable Elements as a Key Contributor to Genomic Plasticity and Virulence Variation in *Magnaporthe oryzae*. *Molecular Plant*, 10(11), 1465–1468.
- Baroncelli, R., Amby, D. B., Zapparata, A., Sarrocco, S., Vannacci, G., Le Floch, G., Harrison, R. J., Holub, E., Sukno, S. A., Sreenivasaprasad, S., & Thon, M. R. (2016). Gene family expansions and contractions are associated with host range in plant pathogens of the genus *Colletotrichum*. *BMC Genomics*, 17, 555.
- Bilinski, P., Albert, P. S., Berg, J. J., Birchler, J. A., Grote, M. N., Lorant, A., Quezada, J., Swarts, K., Yang, J., & Ross-Ibarra, J. (2018). Parallel altitudinal clines reveal trends in adaptive evolution of genome size in *Zea mays*. *PLoS Genetics*, 14(5), e1007162.
- Blin, K., Pascal Andreu, V., de Los Santos, E. L. C., Del Carratore, F., Lee, S. Y., Medema, M. H., & Weber, T. (2019). The antiSMASH database version 2: a comprehensive resource on secondary metabolite biosynthetic gene clusters. *Nucleic Acids Research*, 47(D1),

D625–D630.

Blommaert, J. (2020). Genome size evolution: towards new model systems for old questions.

Proceedings. Biological Sciences / The Royal Society, 287(1933), 20201441.

Bolger, A. M., Lohse, M., & Usadel, B. (2014). Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics*, 30(15), 2114–2120.

Bourguignon, T., Kinjo, Y., Villa-Martín, P., Coleman, N. V., Tang, Q., Arab, D. A., Wang, Z.,

Tokuda, G., Hongoh, Y., Ohkuma, M., Ho, S. Y. W., Pigolotti, S., & Lo, N. (2020). Increased Mutation Rate Is Linked to Genome Reduction in Prokaryotes. *Current Biology: CB*, 30(19), 3848–3855.e4.

Buchfink, B., Xie, C., & Huson, D. H. (2015). Fast and sensitive protein alignment using DIAMOND. *Nature Methods*, 12(1), 59–60.

Bushnell, B., Rood, J., & Singer, E. (2017). BBMerge - Accurate paired shotgun read merging via overlap. *PloS One*, 12(10), e0185056.

Cantarel, B. L., Coutinho, P. M., Rancurel, C., Bernard, T., Lombard, V., & Henrissat, B. (2009).

The Carbohydrate-Active EnZymes database (CAZy): an expert resource for Glycogenomics. *Nucleic Acids Research*, 37(Database issue), D233–D238.

Capella-Gutiérrez, S., Silla-Martínez, J. M., & Gabaldón, T. (2009). trimAl: a tool for automated alignment trimming in large-scale phylogenetic analyses. *Bioinformatics*, 25(15), 1972–1973.

Chan, P. P., Lin, B. Y., Mak, A. J., & Lowe, T. M. (2021). tRNAscan-SE 2.0: improved detection and functional classification of transfer RNA genes. *Nucleic Acids Research*, 49(16), 9077–9096.

Chen, H., Raffaele, S., & Dong, S. (2021). Silent control: microbial plant pathogens evade host immunity without coding sequence changes. *FEMS Microbiology Reviews*, 45(4).

<https://doi.org/10.1093/femsre/fuab002>

Chernomor, O., von Haeseler, A., & Minh, B. Q. (2016). Terrace Aware Data Structure for

- Phylogenomic Inference from Supermatrices. *Systematic Biology*, 65(6), 997–1008.
- Chin, C.-S., Peluso, P., Sedlazeck, F. J., Nattestad, M., Concepcion, G. T., Clum, A., Dunn, C., O'Malley, R., Figueroa-Balderas, R., Morales-Cruz, A., Cramer, G. R., Delledonne, M., Luo, C., Ecker, J. R., Cantu, D., Rank, D. R., & Schatz, M. C. (2016). Phased diploid genome assembly with single-molecule real-time sequencing. *Nature Methods*, 13(12), 1050–1054.
- Christendat, D. (2013). *Ophiostoma ulmi resource browser*.
<http://www.moseslab.csb.utoronto.ca/o.ulmi/>
- Comeau, A. M., Dufour, J., Bouvet, G. F., Jacobi, V., Nigg, M., Henrissat, B., Laroche, J., Levesque, R. C., & Bernier, L. (2014). Functional annotation of the *Ophiostoma novo-ulmi* genome: insights into the phytopathogenicity of the fungal agent of Dutch elm disease. *Genome Biology and Evolution*, 7(2), 410–430.
- Croll, D., & McDonald, B. A. (2012). The accessory genome as a cradle for adaptive evolution in pathogens. *PLoS Pathogens*, 8(4), e1002608.
- Darling, A. C. E., Mau, B., Blattner, F. R., & Perna, N. T. (2004). Mauve: multiple alignment of conserved genomic sequence with rearrangements. *Genome Research*, 14(7), 1394–1403.
- Decena-Segarra, L. P., Bizjak-Mali, L., Kladnik, A., Sessions, S. K., & Rovito, S. M. (2020). Miniaturization, Genome Size, and Biological Size in a Diverse Clade of Salamanders. *The American Naturalist*, 196(5), 634–648.
- Dierckxsens, N., Mardulyn, P., & Smits, G. (2017). NOVOPlasty: de novo assembly of organelle genomes from whole genome data. *Nucleic Acids Research*, 45(4), e18.
- Dobin, A., Davis, C. A., Schlesinger, F., Drenkow, J., Zaleski, C., Jha, S., Batut, P., Chaisson, M., & Gingeras, T. R. (2013). STAR: ultrafast universal RNA-seq aligner. *Bioinformatics*, 29(1), 15–21.
- Emms, D. M., & Kelly, S. (2019). OrthoFinder: phylogenetic orthology inference for comparative genomics. *Genome Biology*, 20(1), 238.
- Flynn, J. M., Hubley, R., Goubert, C., Rosen, J., Clark, A. G., Feschotte, C., & Smit, A. F. (2020).

RepeatModeler2 for automated genomic discovery of transposable element families.

Proceedings of the National Academy of Sciences of the United States of America, 117(17), 9451–9457.

Franco, M. E. E., Wisecaver, J. H., Arnold, A. E., Ju, Y.-M., Slot, J. C., Ahrendt, S., Moore, L. P., Eastman, K. E., Scott, K., Konkel, Z., Mondo, S. J., Kuo, A., Hayes, R. D., Haridas, S., Andreopoulos, B., Riley, R., LaButti, K., Pangilinan, J., Lipzen, A., ... U'Ren, J. M. (2022).

Ecological generalism drives hyperdiversity of secondary metabolite gene clusters in xylarialean endophytes. *The New Phytologist*, 233(3), 1317–1330.

Giovannoni, S. J., Cameron Thrash, J., & Temperton, B. (2014). Implications of streamlining theory for microbial ecology. *The ISME Journal*, 8(8), 1553–1565.

Gregory, T. R. (2002). A bird's-eye view of the C-value enigma: genome size, cell size, and metabolic rate in the class aves. *Evolution; International Journal of Organic Evolution*, 56(1), 121–130.

Grotkopp, E., Rejmánek, M., Sanderson, M. J., & Rost, T. L. (2004). Evolution of genome size in pines (Pinus) and its life-history correlates: supertree analyses. *Evolution; International Journal of Organic Evolution*, 58(8), 1705–1729.

Grützmann, K., Szafranski, K., Pohl, M., Voigt, K., Petzold, A., & Schuster, S. (2014). Fungal alternative splicing is associated with multicellular complexity and virulence: a genome-wide multi-species study. *DNA Research: An International Journal for Rapid Publication of Reports on Genes and Genomes*, 21(1), 27–39.

Han, M. V., Thomas, G. W. C., Lugo-Martinez, J., & Hahn, M. W. (2013). Estimating gene gain and loss rates in the presence of error in genome assembly and annotation using CAFE 3. *Molecular Biology and Evolution*, 30(8), 1987–1997.

Hedges, B., & Kumar, S. (2005). *TIMETREE 5: The timescale of life*. <http://www.timetree.org/>

Hessenauer, P., Fijarczyk, A., Martin, H., Prunier, J., Charron, G., Chapuis, J., Bernier, L., Tanguay, P., Hamelin, R. C., & Landry, C. R. (2020). Hybridization and introgression drive

- genome evolution of Dutch elm disease pathogens. *Nature Ecology & Evolution*, 4(4), 626–638.
- Hinsch, J., Galuszka, P., & Tudzynski, P. (2016). Functional characterization of the first filamentous fungal tRNA-isopentenyltransferase and its role in the virulence of *Claviceps purpurea*. *The New Phytologist*, 211(3), 980–992.
- Hoang, D. T., Chernomor, O., von Haeseler, A., Minh, B. Q., & Vinh, L. S. (2018). UFBoot2: Improving the Ultrafast Bootstrap Approximation. *Molecular Biology and Evolution*, 35(2), 518–522.
- Hoff, K. J., Lomsadze, A., Borodovsky, M., & Stanke, M. (2019). Whole-Genome Annotation with BRAKER. In M. Kollmar (Ed.), *Gene Prediction: Methods and Protocols* (pp. 65–95). Springer New York.
- Ho, L. si T., & Ané, C. (2014). A linear-time algorithm for Gaussian and non-Gaussian trait evolution models. *Systematic Biology*, 63(3), 397–408.
- Hongsanan, S., Maharachchikumbura, S. S. N., Hyde, K. D., Samarakoon, M. C., Jeewon, R., Zhao, Q., Al-Sadi, A. M., & Bahkali, A. H. (2017). An updated phylogeny of Sordariomycetes based on phylogenetic and molecular clock evidence. *Fungal Diversity*, 84(1), 25–41.
- Hultgren, K. M., Jeffery, N. W., Moran, A., & Gregory, T. R. (2018). Latitudinal variation in genome size in crustaceans. *Biological Journal of the Linnean Society. Linnean Society of London*, 123(2), 348–359.
- Junier, T., & Zdobnov, E. M. (2010). The Newick utilities: high-throughput phylogenetic tree processing in the UNIX shell. *Bioinformatics*, 26(13), 1669–1670.
- Kalyaanamoorthy, S., Minh, B. Q., Wong, T. K. F., von Haeseler, A., & Jermini, L. S. (2017). ModelFinder: fast model selection for accurate phylogenetic estimates. *Nature Methods*, 14(6), 587–589.
- Katinka, M. D., Duprat, S., Cornillot, E., Méténier, G., Thomarat, F., Prensier, G., Barbe, V.,

- Peyretailade, E., Brottier, P., Wincker, P., Delbac, F., El Alaoui, H., Peyret, P., Saurin, W., Gouy, M., Weissenbach, J., & Vivarès, C. P. (2001). Genome sequence and gene compaction of the eukaryote parasite *Encephalitozoon cuniculi*. *Nature*, 414(6862), 450–453.
- Katoh, K., & Standley, D. M. (2013). MAFFT multiple sequence alignment software version 7: improvements in performance and usability. *Molecular Biology and Evolution*, 30(4), 772–780.
- Kelkar, Y. D., & Ochman, H. (2012). Causes and consequences of genome expansion in fungi. *Genome Biology and Evolution*, 4(1), 13–23.
- Kuo, C.-H., Moran, N. A., & Ochman, H. (2009). The consequences of genetic drift for bacterial genome complexity. *Genome Research*, 19(8), 1450–1454.
- Li, H. (2018). Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics*, 34(18), 3094–3100.
- Li, H., & Durbin, R. (2009). Fast and accurate short read alignment with Burrows-Wheeler transform. *Bioinformatics*, 25(14), 1754–1760.
- Lim, C. S., Weinstein, B. N., Roy, S. W., & Brown, C. M. (2021). Analysis of Fungal Genomes Reveals Commonalities of Intron Gain or Loss and Functions in Intron-Poor Species. *Molecular Biology and Evolution*, 38(10), 4166–4186.
- Lomsadze, A., Ter-Hovhannisyan, V., Chernoff, Y. O., & Borodovsky, M. (2005). Gene identification in novel eukaryotic genomes by self-training algorithm. *Nucleic Acids Research*, 33(20), 6494–6506.
- Lynch, M. (2006). The origins of eukaryotic gene structure. *Molecular Biology and Evolution*, 23(2), 450–468.
- Lynch, M., & Conery, J. S. (2003). The origins of genome complexity. *Science*, 302(5649), 1401–1404.
- Malerba, M. E., Ghedini, G., & Marshall, D. J. (2020). Genome Size Affects Fitness in the

- Eukaryotic Alga *Dunaliella tertiolecta*. *Current Biology: CB*, 30(17), 3450–3456.e3.
- Mat Razali, N., Cheah, B. H., & Nadarajah, K. (2019). Transposable Elements Adaptive Role in Genome Plasticity, Pathogenicity and Evolution in Fungal Phytopathogens. *International Journal of Molecular Sciences*, 20(14). <https://doi.org/10.3390/ijms20143597>
- Mira, A., Ochman, H., & Moran, N. A. (2001). Deletional bias and the evolution of bacterial genomes. *Trends in Genetics: TIG*, 17(10), 589–596.
- Mistry, J., Chuguransky, S., Williams, L., Qureshi, M., Salazar, G. A., Sonnhammer, E. L. L., Tosatto, S. C. E., Paladin, L., Raj, S., Richardson, L. J., Finn, R. D., & Bateman, A. (2021). Pfam: The protein families database in 2021. *Nucleic Acids Research*, 49(D1), D412–D419.
- Mistry, J., Finn, R. D., Eddy, S. R., Bateman, A., & Punta, M. (2013). Challenges in homology search: HMMER3 and convergent evolution of coiled-coil regions. *Nucleic Acids Research*, 41(12), e121.
- Miyauchi, S., Kiss, E., Kuo, A., Drula, E., Kohler, A., Sánchez-García, M., Morin, E., Andreopoulos, B., Barry, K. W., Bonito, G., Buée, M., Carver, A., Chen, C., Cichocki, N., Clum, A., Culley, D., Crous, P. W., Fauchery, L., Girlanda, M., ... Martin, F. M. (2020). Large-scale genome sequencing of mycorrhizal fungi provides insights into the early evolution of symbiotic traits. *Nature Communications*, 11(1), 5125.
- Mohanta, T. K., & Bae, H. (2015). The diversity of fungal genome. *Biological Procedures Online*, 17, 8.
- Möller, M., & Stukenbrock, E. H. (2017). Evolution and genome architecture in fungal plant pathogens. *Nature Reviews. Microbiology*, 15(12), 756–771.
- Moran, N. A., McCutcheon, J. P., & Nakabachi, A. (2008). Genomics and evolution of heritable bacterial symbionts. *Annual Review of Genetics*, 42, 165–190.
- Morrison, E. N., Emery, R. J. N., & Saville, B. J. (2017). Fungal derived cytokinins are necessary for normal *Ustilago maydis* infection of maize. *Plant Pathology*, 66(5), 726–742.
- Muszewska, A., Taylor, J. W., Szczesny, P., & Grynberg, M. (2011). Independent subtilases

- expansions in fungi associated with animals. *Molecular Biology and Evolution*, 28(12), 3395–3404.
- Muzafar, S., Sharma, R. D., Chauhan, N., & Prasad, R. (2021). Intron distribution and emerging role of alternative splicing in fungi. *FEMS Microbiology Letters*, 368(19).
<https://doi.org/10.1093/femsle/fnab135>
- Nguyen, L.-T., Schmidt, H. A., von Haeseler, A., & Minh, B. Q. (2015). IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies. *Molecular Biology and Evolution*, 32(1), 268–274.
- Organ, C. L., Shedlock, A. M., Meade, A., Pagel, M., & Edwards, S. V. (2007). Origin of avian genome size and structure in non-avian dinosaurs. *Nature*, 446(7132), 180–184.
- Pagel, M., & Meade, A. (2006). Bayesian analysis of correlated evolution of discrete characters by reversible-jump Markov chain Monte Carlo. *The American Naturalist*, 167(6), 808–825.
- Paradis, E., Claude, J., & Strimmer, K. (2004). APE: Analyses of Phylogenetics and Evolution in R language. *Bioinformatics*, 20(2), 289–290.
- Petrov, D. A. (2001). Evolution of genome size: new approaches to an old problem. *Trends in Genetics: TIG*, 17(1), 23–28.
- Petrov, D. A. (2002). Mutational equilibrium model of genome size evolution. *Theoretical Population Biology*, 61(4), 531–544.
- Raffaele, S., & Kamoun, S. (2012). Genome evolution in filamentous plant pathogens: why bigger can be better. *Nature Reviews. Microbiology*, 10(6), 417–430.
- Rawlings, N. D., Barrett, A. J., Thomas, P. D., Huang, X., Bateman, A., & Finn, R. D. (2018). The MEROPS database of proteolytic enzymes, their substrates and inhibitors in 2017 and a comparison with peptidases in the PANTHER database. *Nucleic Acids Research*, 46(D1), D624–D632.
- Revell, L. J. (2011). phytools: an R package for phylogenetic comparative biology (and other things). *Methods in Ecology and Evolution*, 3(2), 217–223.

- Rhodes, J., Abdolrasouli, A., Dunne, K., Sewell, T. R., Zhang, Y., Ballard, E., Brackin, A. P., van Rhijn, N., Chown, H., Tsitsopoulou, A., Posso, R. B., Chotirmall, S. H., McElvaney, N. G., Murphy, P. G., Talento, A. F., Renwick, J., Dyer, P. S., Szekely, A., Bowyer, P., ... Fisher, M. C. (2022). Population genomics confirms acquisition of drug-resistant *Aspergillus fumigatus* infection by humans from the environment. *Nature Microbiology*, 7(5), 663–674.
- Sanderson, M. J. (2003). r8s: inferring absolute rates of molecular evolution and divergence times in the absence of a molecular clock. *Bioinformatics*, 19(2), 301–302.
- Scharf, D. H., Heinekamp, T., & Brakhage, A. A. (2014). Human and plant fungal pathogens: the role of secondary metabolites. *PLoS Pathogens*, 10(1), e1003859.
- Sipos, G., Prasanna, A. N., Walter, M. C., O'Connor, E., Bálint, B., Krizsán, K., Kiss, B., Hess, J., Varga, T., Slot, J., Riley, R., Bóka, B., Rigling, D., Barry, K., Lee, J., Mihaltcheva, S., LaButti, K., Lipzen, A., Waldron, R., ... Nagy, L. G. (2017). Genome expansion and lineage-specific genetic innovations in the forest pathogenic fungi *Armillaria*. *Nature Ecology & Evolution*, 1(12), 1931–1941.
- Stajich, J. E. (2017). Fungal Genomes and Insights into the Evolution of the Kingdom. *Microbiology Spectrum*, 5(4). <https://doi.org/10.1128/microbiolspec.FUNK-0055-2016>
- Stanke, M., & Morgenstern, B. (2005). AUGUSTUS: a web server for gene prediction in eukaryotes that allows user-defined constraints. *Nucleic Acids Research*, 33(Web Server issue), W465–W467.
- Stukenbrock, E. H., & Croll, D. (2014). The evolving fungal genome. *Fungal Biology Reviews*, 28(1), 1–12.
- Tarailo-Graovac, M., & Chen, N. (2009). Using RepeatMasker to identify repetitive elements in genomic sequences. *Current Protocols in Bioinformatics / Editorial Board, Andreas D. Baxevanis ... [et Al.]*, Chapter 4, Unit 4.10.
- Tatusov, R. L., Fedorova, N. D., Jackson, J. D., Jacobs, A. R., Kiryutin, B., Koonin, E. V., Krylov, D. M., Mazumder, R., Mekhedov, S. L., Nikolskaya, A. N., Rao, B. S., Smirnov, S., Sverdlov,

- A. V., Vasudevan, S., Wolf, Y. I., Yin, J. J., & Natale, D. A. (2003). The COG database: an updated version includes eukaryotes. *BMC Bioinformatics*, 4, 41.
- Venditti, C., Meade, A., & Pagel, M. (2011). Multiple routes to mammalian diversity. *Nature*, 479(7373), 393–396.
- Wacker, T., Helmstetter, N., Wilson, D., Fisher, M. C., Studholme, D. J., & Farrer, R. A. (2021). Two-speed genome expansion drives the evolution of pathogenicity in animal fungal pathogens. In *bioRxiv* (p. 2021.11.03.467166). <https://doi.org/10.1101/2021.11.03.467166>
- Walker, B. J., Abeel, T., Shea, T., Priest, M., Abouelliel, A., Sakthikumar, S., Cuomo, C. A., Zeng, Q., Wortman, J., Young, S. K., & Earl, A. M. (2014). Pilon: an integrated tool for comprehensive microbial variant detection and genome assembly improvement. *PloS One*, 9(11), e112963.
- Waterhouse, R. M., Seppey, M., Simão, F. A., Manni, M., Ioannidis, P., Klioutchnikov, G., Kriventseva, E. V., & Zdobnov, E. M. (2018). BUSCO Applications from Quality Assessments to Gene Prediction and Phylogenomics. *Molecular Biology and Evolution*, 35(3), 543–548.
- Wong, B., Leal, I., Feau, N., Dale, A., Uzunovic, A., & Hamelin, R. C. (2020). Molecular assays to detect the presence and viability of *Phytophthora ramorum* and *Grosmannia clavigera*. *PloS One*, 15(2), e0221742.
- Xiao, C.-L., Chen, Y., Xie, S.-Q., Chen, K.-N., Wang, Y., Han, Y., Luo, F., & Xie, Z. (2017). MECAT: fast mapping, error correction, and de novo assembly for single-molecule sequencing reads. *Nature Methods*, 14(11), 1072–1074.
- Yu, G. (2020). Using ggtree to Visualize Data on Tree-Like Structures. *Current Protocols in Bioinformatics / Editorial Board, Andreas D. Baxevanis ... [et Al.]*, 69(1), e96.
- Yu, J. P., Liu, W., Mai, C. L., & Liao, W. B. (2020). Genome size variation is associated with life-history traits in birds. *Journal of Zoology*, 310(4), 255–260.
- Zhang, C., Rabiee, M., Sayyari, E., & Mirarab, S. (2018). ASTRAL-III: polynomial time species

tree reconstruction from partially resolved gene trees. *BMC Bioinformatics*, 19(Suppl 6), 153.