

# UC Berkeley

## UC Berkeley Previously Published Works

### Title

The soybean rust pathogen *Phakopsora pachyrhizi* displays transposable element proliferation that correlates with broad host-range adaptation on legumes

### Permalink

<https://escholarship.org/uc/item/4km2f7w6>

### Authors

Gupta, Yogesh K  
Marcelino-Guimarães, Francismar C  
Lorrain, Cécile  
[et al.](#)

### Publication Date

2022-06-13

### DOI

10.1101/2022.06.13.495685

1 **The soybean rust pathogen *Phakopsora pachyrhizi* displays transposable element**  
2 **proliferation that correlates with broad host-range adaptation on legumes**

3

4 Yogesh K. Gupta<sup>1,2</sup>, Francismar C. Marcelino-Guimarães<sup>3</sup>, Cécile Lorrain<sup>4</sup>, Andrew Farmer<sup>5</sup>,  
5 Sajeet Haridas<sup>6</sup>, Everton Geraldo Capote Ferreira<sup>1,2,3</sup>, Valéria S. Lopes-Caitar<sup>3</sup>, Liliane Santana  
6 Oliveira<sup>3,7</sup>, Emmanuelle Morin<sup>8</sup>, Stephanie Widdison<sup>9</sup>, Connor Cameron<sup>5</sup>, Yoshihiro Inoue<sup>1,2</sup>,  
7 Kathrin Thor<sup>1,2</sup>, Kelly Robinson<sup>1,2</sup>, Elodie Drula<sup>10,11</sup>, Bernard Henrissat<sup>12,13</sup>, Kurt LaButti<sup>6</sup>, Aline  
8 Mara Rudsit Bini<sup>3,7</sup>, Eric Paget<sup>14</sup>, Vasanth Singan<sup>6</sup>, Christopher Daum<sup>6</sup>, Cécile Dorme<sup>14</sup>, Milan  
9 van Hoek<sup>15</sup>, Antoine Janssen<sup>15</sup>, Lucie Chandat<sup>14</sup>, Yannick Tarriotte<sup>14</sup>, Jake Richardson<sup>16</sup>,  
10 Bernardo do Vale Araújo Melo<sup>17</sup>, Alexander Wittenberg<sup>15</sup>, Harrie Schneiders<sup>15</sup>, Stephane  
11 Peyrard<sup>14</sup>, Larissa Goulart Zanardo<sup>17</sup>, Valéria Cristina Holtman<sup>17</sup>, Flavie Coulombier-Chauvel<sup>14</sup>,  
12 Tobias I. Link<sup>18</sup>, Dirk Balmer<sup>19</sup>, André N. Müller<sup>20</sup>, Sabine Kind<sup>20</sup>, Stefan Bohnert<sup>20</sup>, Louisa  
13 Wirtz<sup>20</sup>, Cindy Chen<sup>6</sup>, Mi Yan<sup>6</sup>, Vivian Ng<sup>6</sup>, Pierrick Gautier<sup>14</sup>, Maurício Conrado Meyer<sup>3</sup>, Ralf  
14 Thomas Voegelé<sup>18</sup>, Qingli Liu<sup>21</sup>, Igor V. Grigoriev<sup>6,22</sup>, Uwe Conrath<sup>20</sup>, Sérgio H.  
15 Brommonschenkel<sup>17</sup>, Marco Loehrer<sup>20</sup>, Ulrich Schaffrath<sup>20</sup>, Catherine Sirven<sup>14</sup>, Gabriel  
16 Scalliet<sup>19\*</sup>, Sébastien Duplessis<sup>8\*</sup>, H. Peter van Esse<sup>1,2§\*</sup>

17

18 **Author affiliations:**

- 19 (1) 2Blades, Evanston, Illinois, U.S.A.  
20 (2) The Sainsbury Laboratory, University of East Anglia, Norwich, U.K.  
21 (3) Brazilian Agricultural Research Corporation - National Soybean Research Center  
22 (Embrapa Soja), Paraná, Brazil  
23 (4) Pathogen Evolutionary Ecology, ETH Zürich, Zürich, Switzerland  
24 (5) National Center for Genome Resources, Santa Fe, New Mexico, U.S.A.  
25 (6) U.S. Department of Energy Joint Genome Institute, Lawrence Berkeley National  
26 Laboratory, Berkeley, California, U.S.A.  
27 (7) Department of Computer Science, Federal University of Technology of Paraná (UTFPR),  
28 Paraná, Brazil  
29 (8) Université de Lorraine, INRAE, IAM, Nancy, France  
30 (9) Syngenta Jealott's Hill Int. Research Centre, Bracknell Berkshire, U.K.  
31 (10) AFMB, Aix-Marseille Univ., INRAE, Marseille, France  
32 (11) Biodiversité et Biotechnologie Fongiques, INRAE, Marseille, France  
33 (12) Department of Biological Sciences, King Abdulaziz University, Jeddah, Saudi Arabia  
34 (13) DTU Bioengineering, Technical University of Denmark, Kgs. Lyngby, Denmark  
35 (14) Bayer SAS, Crop Science Division, Lyon, France  
36 (15) KeyGene N.V., Wageningen, The Netherlands  
37 (16) The John Innes Centre, Norwich, U.K.  
38 (17) Departamento de Fitopatologia, Universidade Federal de Viçosa, Viçosa, Brazil  
39 (18) Institute of Phytomedicine, University of Hohenheim, Stuttgart, Germany  
40 (19) Syngenta Crop Protection AG, Stein, Switzerland  
41 (20) Department of Plant Physiology, RWTH Aachen University, Aachen, Germany  
42 (21) Syngenta Crop Protection, LLC, Research Triangle Park, North Carolina, U.S.A.  
43 (22) Department of Plant and Microbial Biology, University of California Berkeley, Berkeley,  
44 California, U.S.A.

45

46 \* These authors contributed equally

47 § Author to whom correspondence should be addressed

## 48 **ABSTRACT**

49 Asian soybean rust, caused by *Phakopsora pachyrhizi*, is one of the world's most economically  
50 damaging agricultural diseases. Despite *P. pachyrhizi*'s impact, the exceptional size and  
51 complexity of its genome prevented generation of an accurate genome assembly. We  
52 simultaneously sequenced three *P. pachyrhizi* genomes uncovering a genome up to 1.25 Gb  
53 comprising two haplotypes with a transposable element (TE) content of ~93%. The  
54 proliferation of TEs within the genome occurred in several bursts and correlates with the  
55 radiation and speciation of the legumes. We present data of clear de-repression of TEs that  
56 mirrors expression of virulence-related candidate effectors. We can see a unique expansion  
57 in amino acid metabolism for this fungus. Our data shows that TEs play a dominant role in *P.*  
58 *pachyrhizi*'s genome and have a key impact on various processes such as host range  
59 adaptation, stress responses and genetic plasticity of the genome.

## 61 **INTRODUCTION**

62 Asian soybean rust caused by the obligate biotrophic fungus *Phakopsora pachyrhizi*, is a  
63 critical challenge for food security and one of the most damaging plant pathogens of this  
64 century (Fig. 1 A) (1). The disease is ubiquitously present in the soybean growing areas of Latin  
65 America, where 210 million metric tons of soybean are projected to be produced in 2022/23  
66 (<https://apps.fas.usda.gov/psdonline/app/index.html>), and on average representing a gross  
67 production value of U.S. \$ 115 billion per season ([https://www.ers.usda.gov/data-](https://www.ers.usda.gov/data-products/season-average-price-forecasts.aspx)  
68 [products/season-average-price-forecasts.aspx](https://www.ers.usda.gov/data-products/season-average-price-forecasts.aspx)). A low incidence of this devastating disease  
69 (0.05%) can already affect yields and, if not managed properly, yield losses are reported of up  
70 to 80% (2, 3). Chemical control in Brazil to manage the disease started in the 2002/03 growing  
71 season (3). In the following season, ~20 million hectares of soybeans were sprayed with  
72 fungicides to control this disease (Fig. 1 A) (3, 4). The cost of managing *P. pachyrhizi* exceeds  
73 \$2 billion USD per season in Brazil alone.

74 The pathogen is highly adaptive and individually deployed resistance genes have been  
75 rapidly overcome when respective cultivars have been released (5, 6). Similarly, the fungal  
76 tolerance to the main classes of site-specific fungicides is increasing, making chemical control  
77 less effective (7-9). Another remarkable feature for an obligate biotrophic pathogen is its wide  
78 host range, encompassing 153 species of legumes within 54 genera to date (10-12).  
79 Epidemiologically, this is relevant as it allows the pathogen to maintain itself in the absence  
80 of soybean on other legume hosts, such as overwintering on the invasive weed Kudzu in the  
81 United States (13).

82 Despite the importance of the pathogen, not much was known about its genetic  
83 makeup as the large genome size (an estimated 1 Gbp), coupled to a high repeat content,  
84 high levels of heterozygosity and the dikaryotic nature of the infectious urediospores of the  
85 fungus have frustrated whole genome assembly effort (14). In a community effort, here we  
86 provide first reference quality assemblies and genome annotations of three *P. pachyrhizi*  
87 isolates.

## 89 **RESULTS AND DISCUSSION**

### 90 **Two superfamilies of transposons dominate the *P. pachyrhizi* genome**

91 The high repeat content and dikaryotic nature of the *P. pachyrhizi* genome poses challenges  
92 to genome assembly methods (14). Recent improvements in sequencing technology and  
93 assembly methods have provided contiguous genome assemblies for several rust fungi (15-  
94 17). Here, we have expanded the effort and provided reference level genome assemblies of

95 three *P. pachyrhizi* isolates (K8108, MT2006, and UFV02) using long-read sequencing  
96 technologies. All three isolates were collected from different regions of South America. We  
97 have used PacBio sequencing for the K8108 and MT2006 isolates and Oxford Nanopore for  
98 the UFV02 isolate to generate three high-quality genomes (fig. S1). Due to longer read lengths  
99 from Oxford nanopore, the UFV02 assembly is more contiguous compared to K8108 and  
100 MT2006 and is used as a reference in the current study (Table 1). The total genome assembly  
101 size of up to 1.25 Gb comprising two haplotigs, makes the *P. pachyrhizi* genome one of the  
102 largest fungal genomes sequenced to date (Fig. 1B). Analysis of the TE content in the *P.*  
103 *pachyrhizi* genome indicates ~93% of the genome consist of repetitive elements, one of the  
104 highest TE contents reported for any organism to date (Fig. 1B and table S1). This high TE  
105 content may represent a key strategy to increase genetic variation in *P. pachyrhizi* (18). The  
106 largest class of TEs are class 1 retrotransposons, that account for 54.0% of the genome. The  
107 class II DNA transposons content is 34.0% (tables S1 and S2). This high percentage of class II  
108 DNA transposons appear to be present in three lineages of rust fungi, the Melampsoraceae  
109 (*Melampsora larici-populina*), Pucciniaceae (*Puccinia graminis* f. sp. *tritici*) and  
110 Phakopsoraceae (*P. pachyrhizi*) (Fig. 1B). The recently assembled large genome (haploid  
111 genome size, 1Gb) of the rust fungus *Austropuccinia psidii* in the family  
112 Sphaerophragmiaceae, however seems to mainly have expanded in retrotransposons (19).  
113 This illustrates that TEs exhibits different evolutionary trajectories in different rust  
114 taxonomical families. Over 80% of the *P. pachyrhizi* genome is comprised by only two  
115 superfamilies of TEs: long terminal repeat (LTR) and terminal inverted repeat (TIR) (Fig. 1B  
116 and table S2). The largest single family of TE are the Gypsy retrotransposons comprising 43%  
117 of the entire genome (Fig. 2A and table S2).

118 To understand the evolutionary dynamics of the different TE families present in the *P.*  
119 *pachyrhizi* genome, we compared the sequence similarities of TEs with their consensus  
120 sequences in the three genomes, which ranges from 65 to 100% sequence identity (fig. S2).  
121 Based on this threshold, TEs were categorised as (1) conserved TEs (copies with more than  
122 95% identity), (2) intermediate TEs (copies with 85 to 95% identity) and (3) divergent TEs  
123 (copies with less than 85% identity) (20). The divergent TEs represent 51.7% - 57.3% of the  
124 total TEs in *P. pachyrhizi*. The average Gypsy retrotransposon composition of the three  
125 isolates is 5.96% conserved, 9.62% intermediate and 15.85% divergent (fig. S3 and tables S3  
126 to S5). Similarly, average TIR composition of the three isolates is 5.3% conserved, 10%  
127 intermediate and 18.53% divergent (fig. S3 and tables S3 to S5). This suggests that i) multiple  
128 waves of TE proliferation have occurred during the history of the species, ii) the invasion of  
129 the two major TE families into the *P. pachyrhizi* genome is not a recent event, and iii) the  
130 presence of conserved TEs indicates ongoing bursts of expansion of TEs in the *P. pachyrhizi*  
131 genome. Therefore, the proportion and distribution of conserved and in a lesser extent of  
132 intermediate TEs indicate that different categories of TEs differentially shaped the genomic  
133 landscape of *P. pachyrhizi* during different times in its evolutionary history (Fig. 2B).

134 We set out to date the Gypsy and Copia TEs in *P. pachyrhizi*, using a TE insertion age  
135 analysis (21, 22). We observe that most TEs were dated less than 100 million years ago (Mya).  
136 We therefore decided to perform a more granulated study taking 1.0 million year intervals  
137 over this period. We observe the start of TEs expansion at around 65 Mya after which the TE  
138 content gradually accumulates (Fig. 2C). We can see a more rapid expansion of TEs in the last  
139 10 Mya, indeed over 40% of the Gypsy and Copia TEs in the genome seem to have arisen  
140 between today and 5 Mya (Fig. 2C). Strikingly, fossil records suggest that legumes started  
141 their main radiation event ~59 Mya (23-26). The climatic oscillations during the past 3 Myr

142 are well known as period of extremely rapid differentiation of species (27). Therefore, the  
143 rapid genome expansion through waves of TE proliferation in *P. pachyrhizi* correlates with the  
144 radiation and adaptation of legumes.

145

### 146 **A subset of TEs is highly expressed during early *in planta* stages of infection**

147 To build a high-quality resource that can facilitate future in-depth analyses, within the  
148 consortium we combined several robust, independently generated RNAseq datasets from all  
149 three isolates that include major soybean infection-stages and *in vitro* germination (Fig. 3, A  
150 and B). Altogether, eleven different stages are captured with seven having overlap of two or  
151 more isolates, representing a total of 72 different transcriptome data sets (Fig. 3C). These  
152 data were used to support the prediction of gene models with the *de novo* annotation  
153 pipeline of JGI MycoCosm (28). Those proteins secreted by the pathogen that impact the  
154 outcome of an interaction between host and pathogen are called effectors and are of  
155 particular interest (29, 30). We used a variety of complementary methods to identify 2,183,  
156 2,027 and 2,125 secreted proteins (the secretome) encoded within the genome assembly of  
157 K8108, MT2006 and UFV02 respectively (31-35) (tables S6-S8). This is a fivefold improvement  
158 when compared to previous transcriptomic studies (36-39). In *P. pachyrhizi*, depending on  
159 methodology 36.73 - 42.30% of these secreted proteins are predicted to be effectors (tables  
160 S6 to S8). We identified 437 common secreted proteins (shared by at least two isolates) that  
161 are differentially expressed at least in one time-point *in planta*, of which 246 are predicted to  
162 be effectors providing a robust set of proteins to investigate in follow-up functional studies  
163 (fig. S4 and table S9).

164 We performed expression analysis on the annotated TEs and observed that 6.66 -  
165 11.65% of TEs are expressed in the three isolates (tables S10 and S11). We compared the TE  
166 expression from different infection stages versus *in vitro* stages (Fig. 2A and tables S12 to S14)  
167 and used the *in planta* RNAseq data from the isolates K8108 and UFV02. A relatively small  
168 subset of TEs (0.03 – 0.25%) are expressed during the early infection stages between 10 to 72  
169 hours post inoculation (HPI) (figs. S5 and S6 and tables S12 and S14). Remarkably, for this  
170 subset we observed a 20 to 70-fold increase in the expression when compared to the spore  
171 and germinated-spore stages, with the expression levels reaching a peak at 24 HPI (figs. S5  
172 and S6). To estimate the impact of the insertion age of this *in planta* induced TE subset, we  
173 performed expression analysis on the conserved, intermediate, and divergent TEs. Although  
174 there is a slight overrepresentation of the conserved TEs, several intermediate TEs and  
175 divergent TEs are also highly expressed during 10 – 24 HPI (fig. S7).

176 To compare the expression profile of this subset of TEs to the predicted effectors, we  
177 used the 246 core effectors and compared these with 25 known and constitutively expressed  
178 housekeeping genes across three isolates. We found that both TE and effector expression  
179 peaked at 24 HPI (Fig. 3D). While expression of effectors remained higher than the 25 selected  
180 housekeeping genes during infection, expression of TEs started to be repressed after 72 HPI  
181 (Fig. 3D). This observation would corroborate the hypothesis of stress-driven TE de-repression  
182 observed in other patho-systems (40-42). However, it also shows that in *P. pachyrhizi* only a  
183 small percentage of the TEs are highly expressed during early infection stages.

184 In several different phytopathogenic species a distinct genomic organization or  
185 compartmentalization can be observed for effector proteins. For example, the bipartite  
186 genome architecture of *Phytophthora infestans* and *Leptosphaeria maculans* in which gene  
187 sparse, repeat rich compartments allow rapid adaptive evolution of effector genes (43). Other  
188 fungi display other organizations such as virulence chromosomes (44, 45) or lineage specific

189 regions (46, 47). However, when interrogating both genomic location and genomic  
190 distribution of the predicted candidate effector genes in *P. pachyrhizi*, we could not detect an  
191 analogous type of organization (fig. S8A-C). In addition, we did not observe evidence of  
192 specific association between TE superfamilies and secreted protein genes (fig. S9), as has  
193 been observed in other fungal species (43, 45, 48-50). Additional analyses comparing the  
194 distance between BUSCO (Benchmarking Universal Single-Copy Orthologue) genes and genes  
195 encoding secreted proteins also showed no specific association (fig. S8D). Therefore, despite  
196 the large genome size and high TE content of *P. pachyrhizi*, its genome appears to be  
197 organized in a similar fashion to other rust fungi with smaller genome sizes (16, 17, 19, 51).  
198 The lack of detection of a specific association between TE and genes in *P. pachyrhizi* may be  
199 due to the extreme nature of the TE invasion with 93% TE observed for this genome.

200

### 201 ***P. pachyrhizi* in South America is a single lineage with high levels of heterozygosity**

202 Rust fungi are dikaryotic, therefore variation can exist both between isolates and between  
203 the two nuclei present in each cell of a single isolate. Long-term asexual reproduction is  
204 predicted to promote divergence between alleles of loci (52), which in principle can increase  
205 indefinitely (53). Some rusts can reproduce both sexually and asexually leading to a mixed  
206 clonal/sexual reproduction. In the rust fungus *P. striiformis* f.sp. *tritici*, asexual lineages  
207 showed a higher degree of heterozygosity between two haploid nuclei when compared to the  
208 sexual lineages (54). In the case of *P. pachyrhizi*, there are clear indications that the population  
209 is propagating asexually in South America based on early studies using simple-sequence  
210 repeats (SSR) and internal transcribed spacer (ITS) sequences (55, 56). Our own data utilizing  
211 high coverage raw Illumina data corroborate these earlier studies as we observed high levels  
212 of heterozygosity; 2.47% for UFV02, 1.61% for K8108 and 1.43% in MT2006, respectively (fig.  
213 S1a). This was further corroborated by mapping the Illumina reads to the genome assembly.  
214 In total 1.2 million variants are found for each isolate, including 0.57 million SNPs, 70% of  
215 which are heterozygous (table S15).

216 We subsequently studied the structural variation (insertions and deletions, repeat  
217 expansion and contractions, tandem expansion and contractions) as well as the haplotype  
218 variation between the three isolates (table S16) (57). Remarkably, the structural variation  
219 between the haplotypes of UFV02 is 163.3 Mb, while the variation between the complete  
220 genomes of the three isolates is 8 to 13 Mb (Fig. 4A). Therefore, inter-haplotype variation is  
221 almost 20 times higher than the total variation between isolates. To look at this inter-  
222 haplotype variation in more detail, we selected contigs larger than 1 Mb to study large  
223 syntenic blocks between isolates and haplotigs. The largest of these contigs, the 1.3 Mb contig  
224 148 from UFV02 has synteny with contig 5809 from K8108, and contigs 220 and 362 from  
225 MT2006 (Fig. 4, C to E), but not with its haplotig genome counterpart within UFV02, which  
226 indicates lack of recombination between haplotypes. This corroborates earlier studies that in  
227 South America *P. pachyrhizi* reproduces only asexually (58).

228 Collection of the monopustule isolates K8108, MT2006, UFV02 is separated in both time  
229 and geographical location (i.e. K8108 from *Colonia*, Uruguay, 2015; MT2006 from *Mato*  
230 *Grosso do Sul*, Brazil, 2006; UFV02 from *Minas Gerais*, Brazil, 2006). To study SNP variation,  
231 we mapped the Illumina data of all three isolates to the reference assembly of UFV02. Given  
232 the high level of heterozygosity and TE content, we focussed our analysis on the now  
233 annotated exome space (table S15a). After removal of SNPs shared between either all three  
234 or two of the isolates, we identified only 3 non-synonymous mutations unique for K8180, 8  
235 non-synonymous mutations for MT2006 and 5 unique non-synonymous mutations for UFV02.

236 For these 16 predicted genes, we found evidence for expression in our transcriptome analyses  
237 for 10 genes. This total number of non-synonymous mutations within exons between the  
238 isolates may appear counterintuitive given the time and space differences between collection  
239 of these isolates. Nonetheless, it is likely that other single pustule isolates identified from  
240 another field would yield a similar number of mutations. Approximately, 6 million spores may  
241 be produced per plant in a single day resulting in  $3 * 10^{12}$  spores per hectare per day (59).  
242 Therefore, the ability to generate variation through mutation cannot be underestimated. We  
243 observed an enrichment of mutations in the upstream and downstream regions of protein  
244 coding genes (table S15b) similar to other rust fungi (60-62). In contrast to the low number of  
245 mutated exons, the number of uniquely expressed genes between the three isolates is  
246 relatively high when compared to the core set of differentially expressed genes (tables S17 to  
247 S19). This may reflect a mechanism in which transcriptional variation is generated via  
248 modification of promotor regions which would have the advantage that coding sequences  
249 that are not beneficial in a particular situation can be “shelved” for later use. This would result  
250 in a set of differentially transcribed genes for different isolates, and a core set of genes that  
251 are transcribed in each isolate.

252

### 253 **The *P. pachyrhizi* genome is expanded in genes related to amino acid metabolism and** 254 **energy production**

255 We subsequently set out to identify expanding and contracting gene families within  
256 *P. pachyrhizi*. To this end, a phylogenetic tree of 17 selected fungal species (table S20a) was  
257 built using 408 conserved orthologous markers. We estimated that the most recent common  
258 ancestor of *P. pachyrhizi* diverged 123.2 - 145.3 million years ago (fig. S10 and table S20b), a  
259 time frame that coincides with the evolution of the Pucciniales (63, 64). We derived gene  
260 families including orthologues and paralogues from a diverse set of plant-interacting fungi  
261 and identified gene gains and losses (i.e. family expansions and contractions) using  
262 computational analysis of gene family evolution (CAFÉ) (table S20a) (65). Genomes of rust  
263 fungi including *P. pachyrhizi* underwent more extensive gene losses than gains, as would be  
264 anticipated for obligate biotrophic parasites (fig. S11). In total, we identified 2,366 contracted  
265 families and 833 expanding families within UFV02 including 792 and 669 families with PFAM  
266 domains, respectively. The most striking and significant contraction in the *P. pachyrhizi*  
267 genome is related to DEAH helicase which is involved in many cellular processes, e.g., RNA  
268 metabolism and ribosome biogenesis (table S21). In contrast, significant expansions in 12  
269 gene families were found including genes encoding glutamate synthase, GMC (glucose-  
270 methanol-choline) oxidoreductase and CHROMO (CHRromatin Organisation MOdifier)  
271 domain containing proteins (table S22). Glutamate synthase plays a vital role in nitrogen  
272 metabolism, and its ortholog in the ascomycete *Magnaporthe oryzae* *MoGLT1* is required for  
273 conidiation and complete virulence on rice (66). GMC oxidoreductase exhibits important  
274 auxiliary activity 3 (AA3\_2) according to the Carbohydrate-Active enzymes (CAZy) database  
275 (67) and is required for the induction of asexual development in *Aspergillus nidulans* (68). An  
276 extensive approach was used for the global annotation of CAZyme genes in *P. pachyrhizi*  
277 genomes and after comparison with other fungal genomes, we also found clear expansions  
278 in glycoside hydrolases (GH) family 18 and glycosyltransferases (GT) family 1 (table S23). GH18  
279 chitinases are required for fungal cell wall degradation and remodelling, as well as multiple  
280 other physiological processes including nutrient uptake and pathogenicity (69, 70).

281 The Phakopsoraceae to which *P. pachyrhizi* belongs represents a new family branch in  
282 the order Pucciniales (71). With three *P. pachyrhizi* genome annotation replicates available,

283 next to the above CAFÉ-analysis, we can directly track gene family expansions and  
284 contractions in comparison to genomes previously sequenced. We therefore compared *P.*  
285 *pachyrhizi* to the taxonomical related families Coleosporiaceae, Melampsoraceae and  
286 Pucciniaceae, which in turn may reveal unique lifestyle adaptations (Table 2).

287 The largest uniquely expanded gene family (531-608 members) in *P. pachyrhizi*  
288 comprises sequences containing the Piwi (P-element Induced Wimpy testes in *Drosophila*)  
289 domain (Table 2). Typically, the Piwi domain is found in the Argonaute (AGO) complex where  
290 its function is to cleave ssRNA when guided by dsRNA (72). Interestingly, classes of longer-  
291 than-average miRNAs known as Piwi-interacting RNAs (piRNAs) that are 26-31 nucleotides  
292 long are known in animal systems. In *Drosophila*, these piRNAs function in nuclear RNA  
293 silencing where they associate specifically with repeat associated small interfering RNA  
294 (rasiRNAs) that originate from TEs (73). As in other fungal genomes, the canonical genes  
295 coding for large AGO proteins with canonical Argonaute, PAZ and Piwi domains can be  
296 observed in the genome annotation of the three *P. pachyrhizi* isolates. The hundreds of  
297 expanded predicted Piwi genes consist of short sequences of less than 500 nt containing only  
298 a partial Piwi domain aligning with the C-terminal part of the Piwi domain in the AGO protein.  
299 Some of these genes are pseudogenes marked by stop codons or encoding truncated protein  
300 forms, while others exhibit a partial Piwi domain starting with a methionine and eventually  
301 exhibiting a strong prediction for an N-terminal signal peptide. These expanded short Piwi  
302 genes are surrounded by TEs, several hundreds of which, but not all, are found in close  
303 proximity to specific TE consensus identified by the REPET analysis in the three *P. pachyrhizi*  
304 isolates (e.g. Gypsy, CACTA and TIR; fig. S12). However, no systematic and significant  
305 association could be made due to the numerous nested TEs present within the genome (74).  
306 Moreover, none of the expanded short Piwi domain genes are expressed in the conditions we  
307 tested. However, in many systems, Piwis and piRNAs play crucial roles during specific  
308 developmental stages where they influence epigenetic, germ cell, stem cell, transposon  
309 silencing, and translational regulation (75). Finally, the domain present in these short Piwi  
310 genes is partial and we do not know whether they retain any RNase activity. Therefore, we  
311 cannot validate at this stage the function of this family, which warrants further study and  
312 attention as it may represent either a new type of TE-associated regulator within *P.*  
313 *pachyrhizi*, or an extreme expansion of a control mechanism to deal with this highly repetitive  
314 genome.

315 Several families related with amino acid metabolism have expanded greatly when  
316 compared to the respective families in other rust fungi, most notably Asparagine synthase  
317 (KOG0573), which has ~75 copies in *P. pachyrhizi* compared to two copies in Pucciniaceae and  
318 one copy in Melampsoraceae (Table 2). Similarly, expanded gene families can be observed in  
319 citrate synthase (KOG2617), malate synthase (KOG1261), NAD-dependent malate  
320 dehydrogenase (KOG1494). These enzymes are involved in energy production and conversion  
321 via the citrate cycle required to produce certain amino acids and the reducing agent NADH  
322 (Table 2). Next to the molecular dialogue with effector proteins, plant-pathogen  
323 interactions are a “tug-of-war” of resources between the host and the pathogen (76). A key  
324 resource to secure in this process is nitrogen, a key raw material needed to produce proteins.  
325 Therefore, the expansion in amino acid metabolism may reflect an adaptation to become  
326 more effective at securing this resource. Alternatively, the expanded categories also may  
327 reflect the metabolic flexibility needed to facilitate the broad host range of *P. pachyrhizi*,  
328 which to date comprises 153 leguminous species in 56 genera (12).



329 Association with TEs are often a sign for adaptive evolution as they facilitate the genetic  
330 leaps required for rapid phenotypic diversification (41, 77-79). We therefore investigated  
331 whether the expansion in amino acid metabolism could reflect a more recent adaptation by  
332 studying the TEs in these genomic regions. Furthermore, as described above, a distinction can  
333 be made between more recent burst of TE activity (high conservation of the TEs) and older  
334 TE burst leading to degeneration of the TE sequence consensus (80). However, despite the  
335 presence of several copies of specific TE subfamilies (i.e. related to the same annotated TE  
336 consensus) in the vicinity of the surveyed expanded families such as amino acid metabolism,  
337 CAZymes and transporter related genes (figs. S13 and S14), no significant enrichment could  
338 be observed for any particular TE when compared to the overall TE content of the genome.  
339 This may reflect the challenge of making such clear associations due to the continuous  
340 transposition activity, which results in a high plasticity of the genomic landscape and a highly  
341 nested TE structure. Alternatively, it may suggest a more ancient origin of these unique  
342 expansions that have subsequently been masked by repetitive episodes of relaxed TE  
343 expression (figs. S15 and S16).

344

## 345 **CONCLUSION**

346 A comprehensive resource to fuel further studies on *P. pachyrhizi* was highly needed given  
347 the economic and social impact this pathogen can have for farmers and global food security.  
348 The *P. pachyrhizi* genome is one of the largest fungal genomes sequenced to date with a total  
349 assembly size of up to 1.25 Gb. The genome is highly repetitive with ~93% of the genome  
350 consisting of TEs, of which two superfamilies make up 80%. The three *P. pachyrhizi* isolates  
351 collected from South America represent a single clonal lineage with high levels of  
352 heterozygosity. Studying the TEs in detail, we demonstrate that the expansion of TEs within  
353 the genome correlates with the radiation and speciation of the legumes and did so in several  
354 bursts. Although TEs are tightly controlled during sporulation and appressoria formation, we  
355 can see a clear relaxation of repression during the *in planta* life stages of the pathogen. Due  
356 to the nested TEs, it is not possible at present to correlate specific TEs to specific expanded  
357 gene families. However, we can see that the *P. pachyrhizi* genome is expanded in genes  
358 related to amino acid metabolism and energy production which may represent key lifestyle  
359 adaptations. Overall, our data unveil that TEs that started their proliferation with the  
360 radiation of the Leguminosae may play a prominent role in the *P. pachyrhizi*'s genome and  
361 have a key impact on a variety of processes such as host range adaptation, stress responses  
362 and plasticity of the genome. The high-quality genome assembly and transcriptome data  
363 presented here are a key resource for the community. It represents a critical step for further  
364 in-depth studies of this pathogen to develop new methods of control and to better  
365 understand the molecular dialogue between *P. pachyrhizi* and its agriculturally relevant host,  
366 Soybean.

367

## 368 **MATERIALS AND METHODS**

### 369 **Fungal strain and propagation**

370 *P. pachyrhizi* isolates, K8108, MT2006 and UFV02 (81) are single uredosporal isolates collected  
371 from Uruguay (*Colonia* in 2015), Brazil (*Mato Grosso do Sul* in 2006) and Brazil (*Minas Gerais*  
372 in 2006), respectively. The isolates were propagated on susceptible soybean cultivars Abelina,  
373 Thorne, Toliman and Williams 82 by spraying suspension of urediniospores 1 mg ml<sup>-1</sup> in 0.01  
374 % (vol/vol) Tween-20 in distilled water onto 21-day-old soybean plants followed by 18 h  
375 incubation in an incubation chamber at saturated humidity and at 22°C in the dark. Infected

376 plants were kept at 22°C, 16-h day/8-h night cycle and 300  $\mu\text{mol s}^{-1} \text{m}^{-2}$  light. After 14 DPI  
377 (days post inoculation), the pustules were formed, and the urediospores were harvested  
378 using a Cyclone surface sampler (Burkard Manufacturing Co. Ltd.) and stored at -80°C.

379

### 380 **Genomic DNA extraction and genome sequencing**

381 The high molecular weight (HMW) genomic-DNA was extracted using a carboxyl-modified  
382 magnetic bead protocol (82) for K8108, a CTAB-based extraction for MT2006 (83), and a  
383 modified CTAB protocol for UFV02 (84).

384 For K8108, a 20-kb PacBio SMRTbell library was prepared by Genewiz (South  
385 Plainfield, NJ) with 15-kb Blue Pippin size selection being performed prior to sequencing on a  
386 PacBio Sequel system (Pacific Biosciences, Menlo Park, CA). The K8108 PacBio Sequel genomic  
387 reads yielding 69 Gbp of sequence data were error corrected using MECAT (85); following  
388 parameter optimization for contiguity and completeness the longest corrected reads yielding  
389 50x coverage were assembled with MECAT's mecat2canu adaptation of the Canu assembly  
390 workflow (86), using an estimated genome size of 500 Mbp and an estimated residual error  
391 rate of 0.02. The resulting assembly had further base pair-level error correction performed  
392 using the Arrow polishing tool from PacBio SMRTTools v5.1.0.26412 (87).

393 MT2006 genome was sequenced using Pacific Biosciences platform. The DNA sheared  
394 to >10kb using Covaris g-Tubes was treated with exonuclease to remove single-stranded ends  
395 and DNA damage repair mix followed by end repair and ligation of blunt adapters using  
396 SMRTbell Template Prep Kit 1.0 (Pacific Biosciences). The library was purified with AMPure  
397 PB beads and size selected with BluePippin (Sage Science) at >6 kb cutoff size. PacBio  
398 Sequencing primer was then annealed to the SMRTbell template library and sequencing  
399 polymerase was bound to them using Sequel Binding kit 2.0. The prepared SMRTbell template  
400 libraries were then sequenced on a Pacific Biosystem's Sequel sequencer using v2 sequencing  
401 primer, 1M v2 SMRT cells, and Version 2.0 sequencing chemistry with 1x360 and 1x600  
402 sequencing movie run times. The MT2006 genome was assembled with Falcon (88), improved  
403 with finisherSC version 2.0 (89), and polished with Arrow version SMRTLink v5.1.0.26412 (87).

404 For UFV02, the PromethION platform of Oxford nanopore technology (ONT) (Oxford,  
405 UK) was used for long-read sequencing at Keygene N.V. (Wageningen, The Netherlands). The  
406 libraries with long DNA fragments were constructed and sequenced on the PromethION  
407 platform. The raw sequencing data of 110 Gbp was generated and was base-called using ONT  
408 Albacore v2.1 available at <https://community.nanoporetech.com>. The UFV02 genome  
409 assembly, the longest 15, 20, 25, 30, 34, 40 and 56x nanopore reads were assembled using  
410 the Minimap2 and Miniasm pipeline (90). To improve the consensus, error correction was  
411 performed three times with Racon using all the nanopore reads (91). The resulting assembly  
412 was polished with 50x Illumina PCR-free 150 bp paired-end reads mapped with bwa (92) and  
413 Pilon (93), repeated three times. We assessed the BUSCO scores after each step to compare  
414 the improvement in the assemblies.

415

### 416 **TE identification and classification**

417 To annotate the TEs in the *P. pachyrhizi* genome, we used the two pipelines of the REPET  
418 package (<https://urgi.versailles.inra.fr/Tools/REPET> (94, 95)). We identified the *de novo* TE  
419 consensus library from a subset of 300 Mbp of the longest contigs of each *P. pachyrhizi*  
420 isolates MT2006, K8108 and UFV02 following the large genomes repeats annotation  
421 recommendations (96). The subset strategy is based on the rationale that high abundance of  
422 TEs renders copy identification feasible in a subset of the whole genome. TE consensus

423 libraries were built independently for each isolate, using the TEdenovo pipeline with default  
424 parameters. Briefly, repeats were detected based on similarity approach using Blaster (95),  
425 then clustered combining Grouper (95), Piler (97), Recon (98) and aligned using MAP (99) to  
426 generate consensus sequences of repeats. Consensus sequences were classified and filtered  
427 based on Wicker *et al.* (2007) classification using PASTEClassifier (100, 101). Each TE  
428 consensus library was used to annotate independently the whole genome of *P. pachyrhizi*  
429 isolates with one first round of TEannot pipeline with default parameters. Briefly, the first  
430 round of TEannot consisted of a similarity search of TEs based on the consensus library  
431 combining Blaster (95), CENSOR (102) and RepeatMasker (<http://www.repeatmasker.org/>).  
432 Finally, the distant fragments were connected with a 'long joint procedure' to build copies of  
433 TEs. This first round of TEannot identified a library of TE consensus with at least one full length  
434 fragment (i.e. a TE copy that is aligned over more than 95% to its cognate consensus). The  
435 reduced full length fragment library was finally used for a second round of TEannot pipeline,  
436 to generate the final repeats annotation of *P. pachyrhizi* isolates.

437

### 438 **TE analysis**

439 The TE insertions are categorised based on the sequence identity 1) TEs with less than 85%  
440 sequence identity to the consensus, called old insertions, 2) TEs with 85-95% sequence  
441 identity are intermediate, and 3) TEs with more than 95% identity represent recent insertions  
442 (figs. S2 and S3) (20). All three isolates show common patterns of consensus identity and a  
443 majority of the TEs show an intermediate age of insertions (fig. S2). The retrotransposon  
444 superfamilies such as terminal-repeat retrotransposons in Miniature (TRIMs) are the most  
445 recent expansion and long interspersed nuclear element (LINE) and large retrotransposon  
446 derivative (LARD) superfamilies are the most ancient insertion in the *P. pachyrhizi* genome  
447 (fig. S3). To verify the relationship between secreted genes and TEs, we calculated the  
448 distance between these features using Bedtools (103) with Closest algorithm, which returns  
449 the smallest genomic distance between two features. From the results obtained, we  
450 calculated the number of TEs neighbouring each secreted gene, grouped them by each TE  
451 superfamily and built the graphs. The tools used for analysis and graphs construction were  
452 Pandas v.1.3.4 and Seaborn 0.11.2 libraries, together with Python 3.9.7.

453

### 454 **Insertion age of LTR-retrotransposons**

455 Full-length LTR-retrotransposons were identified from the genomes set assemblies using  
456 LTRharvest with default parameters, this tool belongs to the GenomeTools genome analysis  
457 software v1.6.1 (104). LTRs annotated as Gypsy or Copia were used for molecular dating,  
458 selection was based on a BLASTX against Repbase v20.11 (105). 3' and 5' LTR sequences were  
459 extracted and aligned with mafft v7.471 (106), alignments were used to calculate Kimura's 2P  
460 distances (107). The insertion age was determined using the formula  $T = K / 2r$ , with K the  
461 distance between the 2 LTRs and r the fungal substitution rate of  $1.05 \times 10^{-9}$  nucleotides per  
462 site per year (21, 22).

463

### 464 **Molecular dating and Phylogenetic analysis**

465 The phylogenetic tree was generated after alignment of 408 conserved orthologous markers  
466 identified from at least 13 out of 17 genomes using PHYling  
467 ([https://github.com/stajichlab/PHYling\\_unified](https://github.com/stajichlab/PHYling_unified)). The sequences were aligned and  
468 concatenated into a super-alignment with 408 partitions. The phylogenetic tree was built with  
469 RAXML-NG (v0.9.0) using a partitioned analysis and 200 bootstraps replicates. Molecular

470 dating was established with mcmctree from PAML v4.8. Calibration points were extracted to  
471 Puccinalies (64) and Sordariomycetes–Leotiomycetes (108). The 95% highest posterior  
472 density (HPD) values calibrated to the node.

473

#### 474 **Genome annotation**

475 The gene predications and annotations was performed in the genomes of the three isolates  
476 in parallel using the JGI Annotation Pipeline (28). TE masking was done during the JGI  
477 procedure, which detects and masks repeats and TEs. Later on, the extensive TE classification  
478 performed with REPET was imported and visualized as a supplementary track onto the  
479 genome portals. RNAseq data from each isolate (see section below) was used as intrinsic  
480 support information for the gene callers from the JGI pipeline. The gene prediction procedure  
481 identifies a series of gene models at each gene locus and proposes a best gene model which  
482 allows to define a filtered gene catalog. Translated proteins deduced from gene models are  
483 further used for functional annotation according to international reference databases. All the  
484 annotation information is collected into an open public JGI genome portal in the MycoCosm  
485 (<https://mycocosm.jgi.doe.gov/Phakopsora>) with dedicated tools for community-based  
486 annotation (28, 109). In total, 18,216, 19,618 and 22,467 gene models were predicted from  
487 K8108, MT2006 and UFV02, respectively (table S24); of which 10,492, 10,266 and 9,987 genes  
488 were functionally annotated. We have performed differential expression analyses using the  
489 germinated spores as a reference point in each of the three isolates (fig. S17 and tables S17  
490 to S19). A total of 3,608 common differentially expressed genes (DEGs) were identified in at  
491 least one condition shared between two or more isolates (fig. S18 and table S25).

492

#### 493 **Quality assessment of the whole-genome assemblies**

494 The whole-genome assemblies of *P. pachyrhizi* were evaluated using two different  
495 approaches. First, we used BUSCO version 5.0 (110), to assess the genic content based on  
496 near-universal single-copy orthologs with basidiomycetes\_odb9 database including 1335  
497 gene models. Second, K-mer's from different assemblies were compared using KAT version  
498 2.4.1 (111). Genome heterozygosity was estimated using GenomeScope 2.0 (112).

499

#### 500 **Identification of assembly haplotigs**

501 The haplotypes were phased using the purge-haplotig pipeline (113) using Illumina WGS data.  
502 The haplotigs were aligned with their corresponding primary contigs using Mummer-4.0 for  
503 UFV02 (114). Assemblytics was subsequently used to define six major types of structural  
504 variants (57), including insertions and deletions, repeat expansion and contractions, and  
505 tandem expansion and contractions.

506 The assembly was compared to itself using blastn (NCBI-BLAST+ 2.7.1) with  
507 max\_target\_seqs = 10 and culling\_limit = 10. After filtering for sequences matching  
508 themselves, overlaps among the remaining high scoring segment pairs (HSPs) of  $\geq 500$  bp  
509 and  $\geq 95\%$  identity were consolidated with an interval tree requiring 60% overlap, then  
510 chained using MCScanX\_h (115) to determine collinear series of matches, requiring 3 or more  
511 collinear blocks and choosing as a candidate haplotig sequences having at least 40% of their  
512 length subsumed by a chain corresponding to a longer contig sequence. For downstream  
513 analyses requiring a single haplotype representation, hard masking was applied to remove  
514 overlapped regions from the haplotigs using BEDtools v2.27.0 (103). To identify genes  
515 correspondence among the three isolates, we used LiftOff software (116). The genome  
516 assembly of each isolate was used as a reference to map the other two isolates' gene

517 catalogue with >95% coverage and identity of >95%. The correspondence was established  
518 based on the gene annotation coordinates of each reference genome and the mapping  
519 coordinates from liftoff results (table S26).

520

### 521 **Read mapping, variant calling and SNP effect prediction**

522 Illumina paired-end reads of the three isolates were trimmed with Trimmomatic v0.36 (117)  
523 to remove adapters, barcodes, and low-quality sequences with the following parameters:  
524 illuminaclip = TruSeq3-PE-2.fa:2:30:10, slidingwindow = 4:20, minlen = 36. Then, sequence  
525 data from all three isolates were aligned to the reference assembly of *P. pachyrhizi* UFV02  
526 v2.1 using BWA version 0.7.17 with the BWA-mem algorithm (92), with the options -M -R.  
527 Alignment files were converted to BAM files using SAMtools v1.9 (118), and duplicated reads  
528 were removed using the Picard package (<https://broadinstitute.github.io/picard/>). The GATK  
529 v3.8.1 software (119) was used to identify and realign poorly aligned reads around InDels  
530 using Realigner Target Creator and Indel Realigner tools, creating a merged bam file for all  
531 the three isolates. The subsequent realigned BAM file was used to calling SNPs and InDels  
532 using HaplotypeCaller in GATK and filtering steps were performed to keep only high-quality  
533 variants, as following: the thresholds setting as: "QUAL < 30.00 || MQ < 40.00 || SOR > 3.00  
534 || QD < 2.00 || FS > 60.00 || MQRankSum < -12.500 || ReadPosRankSum < -8.00 ||  
535 ReadPosRankSum > 8.00". The resulting SNPs and InDels were annotated with snpEffect v4.1  
536 (120).

537

### 538 **Infection and disease progression**

539 *P. pachyrhizi* is an obligate biotrophic fungus, which forms a functional appressorium to  
540 penetrate the host epidermal layer within 12 HPI (hours post inoculation) (121). The  
541 penetrated epidermal cell dies after fungus establishes the penetration hyphae (PH) and  
542 forms the primary invasive hyphae (PIH) in the mesophyll cells after 24 HPI (Figs. 3, A and B).  
543 The PIH differentiates and forms a haustorial mother cell, which establishes the haustorium  
544 in the spongy parenchyma cells. At 72 HPI, the fungus colonises the spongy and palisade  
545 parenchyma cells (spc and ppc) (122) (Figs. 3, A and B). At 168 HPI, the uredinium starts to  
546 develop in the palisade parenchyma. At 196 HPI, the epidermal layer is broken and the fully  
547 developed uredinia emerges. Each pustule forms thousands of urediniospores and carry on  
548 the infection (fig. S19).

549

### 550 **Sample preparation for RNAseq**

551 For expression analysis, 11 different stages were evaluated, with eight stages having overlap  
552 of two or more isolates. These stages were nominated 1-11 as illustrated in Fig. 3C. For K8108,  
553 seven *in vitro*, one *on planta* and eight *in planta* samples, each with three biological replicates,  
554 were generated and used to prepare RNA libraries. To get *in vitro* germ tubes and fungal  
555 penetration structures a polyethylene foil (dm freezer bag, Karlsruhe, Germany) was placed  
556 in glass plates and inoculated with a spore suspension (2 mg ml<sup>-1</sup>). Each biological replicate  
557 corresponded to 500 cm<sup>2</sup> foil and ~4 mg urediniospores. The plates were incubated at 22°C  
558 in the dark at saturated humidity for 0.5, 2, 4 or 8 h. After incubation, the spores were  
559 collected using a cell scraper. For the appressoria-enriched sample, urediniospore  
560 concentration was doubled and the plates rinsed with sterile water after 8 h of incubation  
561 prior to collection. The material was ground with mortar and pestle in liquid nitrogen. The  
562 time 0.5 h was considered as spore (Spore, Psp - stage 1), the 2 h as germinated spore  
563 (Germinated spore, PspG - stage 2), and the 8 h rinsed as appressoria enriched sample *in vitro*

564 (stage 3). The samples of spores collected after 4 and 8 h were not used for expression  
565 analysis. To obtain *in planta* fungal structures, three-week-old soybean plants (Williams 82)  
566 were inoculated as mentioned above. After 8 HPI, liquid latex (semi-transparent low  
567 ammonium, Latex-24, Germaringen, Germany) was sprayed (hand spray gun with gas unit,  
568 Preval, Bridgeview, USA) until complete leaf coverage. After drying off, latex was removed. It  
569 contained the appressoria and spores from the leaf surface but no plant tissue. This sample  
570 was considered as enriched in appressoria on plant and is exclusive for K8108 isolate (stage  
571 4). Three middle leaflets of different plants were bulked for each sample and ground in liquid  
572 nitrogen using mortar and pestle. The inoculated leaf samples were harvested at 10, 24, 72  
573 and 192 HPI (stages 5, 6, 8 and 10) for the *in planta* gene expression studies.

574 For MT2006, the germ tubes and appressorium were produced on polyethylene (PE)  
575 sheets where urediniospores were finely dusted with household sieves held in a double layer  
576 of sifting. The PE sheets were then sprayed with water using a chromatography vaporizer and  
577 were kept at 20°C, 95% humidity in the dark. For germ tubes the structures were scratched  
578 from the PE sheets after 3 h (stage 2) and for appressoria after 5 h (stage 3). Formation of  
579 both germ tubes and appressoria was checked microscopically. The *in vitro* samples were only  
580 used when there were at least 70% germ tubes or appressoria, respectively. The structures  
581 were dried by vacuum filtration and stored in 2-ml microcentrifuge tubes at -70°C after  
582 freezing in liquid nitrogen. The resting spores came directly from storage at -70 °C (stage 1).  
583 For the *in planta* samples 21 days old soybean cultivar Thorne were sprayed with a suspension  
584 containing 0.01% Tween-20, 0.08% milk-powder and 0.05% urediniospores. The inoculated  
585 plants were kept as mentioned previously. The samples were taken using a cork borer (18  
586 mm diameter) at 192 and 288 HPI (stages 10 and 11). Three leaf pieces were collected for  
587 each sample (three times and from three different plants) for every time-point and stored in  
588 liquid nitrogen and kept at -80°C.

589 For UFV02, the spore suspension of  $1 \times 10^6$  spores  $\text{ml}^{-1}$  concentration was prepared in  
590 0.01% v/v Tween-20. Four weeks old soybean plants were sprayed thoroughly on the abaxial  
591 surface of the leaves, and the plants were kept at saturated humidity in the dark for 24 h.  
592 After 24 h, plants were kept at 22°C and 16/8-h light/dark cycle. The leaf samples were  
593 collected from non-inoculated plants (0h) and infection-stages at 12, 24, 36, 72 and 168 HPI  
594 (stages 5, 6, 7, 8 and 9). Infection assay was performed in three biological replicates and three  
595 plants were used for each replicate. All the samples were stored in liquid nitrogen after  
596 collection and kept in -80°C for further processing (stage 1). Spores were harvested after 14  
597 days post-inoculation and used for the RNA extraction. The urediniospores were germinated  
598 *in vitro* on the water surface in a square petri dish and kept for 6 h at 24°C (stage 2). The  
599 germinated-urediniospores were collected in a falcon tube and snap freeze in liquid nitrogen.  
600 The samples were freeze-dried and kept at -80°C until further processing. The un-inoculated  
601 plants (0h) were not used in the expression analysis.

#### 602 **RNA isolation, sequencing, and transcriptome assembly**

603 All the samples were ground in liquid nitrogen, and the total RNA was extracted using the  
604 Direct-zol RNA Miniprep Plus Kit (ZymoResearch, Freiburg, Germany), the mirVana™ miRNA  
605 Isolation Kit (Ambion/life technologies, Calsbad, CA, USA), and TRIzol™ reagent (Invitrogen)  
606 according to the manufacturer's protocols for K8108, MT2006, and UFV02, respectively. The  
607 quality of RNA was assessed using TapeStation instrument (Agilent, Santa Clara, CA) or the  
608 Agilent 2100 bioanalyzer.

609 The RNA libraries from K8108 were normalized to 10 mM, pooled, and sequenced at  
610 150-bp paired-end on the HiSeq X instrument at Genewiz (South Plainfield, NJ), with ten

611 samples per lane. The transcriptome of MT2006 was sequenced with Illumina. Stranded cDNA  
612 libraries were generated using the Illumina Truseq Stranded mRNA Library Prep kit. mRNA  
613 was purified from 1 ug of total RNA using magnetic beads containing poly-T oligos. mRNA was  
614 fragmented and reversed transcribed using random hexamers and SSII (Invitrogen) followed  
615 by second strand synthesis. The fragmented cDNA was treated with end-pair, A-tailing,  
616 adapter ligation, and 8 cycles of PCR. The prepared libraries were quantified using KAPA  
617 Biosystem's next-generation sequencing library qPCR kit (Roche) and run on a Roche  
618 LightCycler 480 real-time PCR instrument. The quantified libraries were then multiplexed and  
619 the pool of libraries was prepared for sequencing on the Illumina HiSeq sequencing platform  
620 utilizing a TruSeq paired-end cluster kit, v4, and Illumina's cBot instrument to generate a  
621 clustered flow cell for sequencing. Sequencing of the flow cell was performed on the Illumina  
622 HiSeq 2500 sequencer using HiSeq TruSeq SBS sequencing kits, v4, following a 2x150 indexed  
623 run recipe. The RNA samples of UFV02 were sequenced at the Earlham Institute (Norwich,  
624 UK) on Illumina HiSeq 2500 platform with 250-bp paired-end reads. Eight different samples  
625 (as mentioned above) in three biological replicates were used for the RNA library preparation.  
626 All 24 libraries were multiplexed and sequenced on six lanes of HiSeq 2500.

627 The low-quality RNA-seq reads were processed and trimmed using Trimmomatic  
628 version 0.39 (117) with the parameters ILLUMINACLIP:2:30:10 LEADING:3 HEADCROP:10  
629 SLIDINGWINDOW:4:25 TRAILING:3 MINLEN:40 and read quality was assessed with FastQC  
630 version 0.11.5 (<https://www.bioinformatics.babraham.ac.uk/projects/fastqc/>). The high  
631 quality reads were filtered for any possible contamination among the fungi reads using  
632 Kraken2 software and parameter --unclassified-out for soybean genome and any possible  
633 contaminant species (123). After all filtering steps, reads from each library were mapped  
634 against the three isolates assemblies using STAR v2.7.6a (124). Parameters for mapping were  
635 (--outSAMtype BAM SortedByCoordinate, --outFilterMultimapNmax 100, --  
636 outFilterMismatchNmax 2, --outSAMattrIHstart 0, --winAnchorMultimapNmax 200, and --  
637 outWigType bedGraph). After mapping, duplicated reads were removed using Picard v.2.23.2.  
638 HTseq was used to count reads and DEseq2 to determine differential expression using spore  
639 germinated condition (spore – stage 1) in each transcriptome experimental design as the  
640 calibrator.

641 To validate gene annotation dedup-BAM files were analysed using StringTie v2.1.2  
642 (125) and the gtf files obtained were merged (-m 600 -c 5) for genes and (-m 200 -c 5) for TE  
643 (TE). The final gtf file was compared with each of the genome annotation file per isolate using  
644 gffcompare (126) software to validate the annotate genes and TEs. We detected 18,132,  
645 19,467, and 22,347 genes presenting transcriptional evidence in K8108, MT2006 and UFV  
646 genomes respectively, demonstrating high sensitivity (> 93.9%) and precision in a locus level  
647 (> 75.4%) in all three isolates (fig. S20 and table S27). For functional annotation, genes were  
648 considered expressed when each transcriptome reads were mapped against its respective  
649 reference genome considering the criteria of TPM (Transcripts Per Kilobase Million) values >  
650 0 in at least two biological replicates.

651 The BAM-dedup files obtained as above described were applied for TE expression  
652 analyses using Tetrascript software (127). TE read counts were normalized between  
653 replicates in different conditions using R/Bioconductor package EdgeR v.3.1(128, 129). Only  
654 TEs with a minimum of one read in at least two replicates were considered in this  
655 normalization step. Libraries were normalized with the TMM method (130) and CPM (counts  
656 per million) were generated with the EdgeR v.3.13. To better understand the expression  
657 distribution of TEs in the K8108, MT2006 and UFV02 genomes, we constructed boxplot plots

658 to visualize the variation of expression values (average CPM) in each of their conditions. For  
659 this, we calculated the arithmetic means, the standard deviation, and the quartile values of  
660 the TEs expression in each condition for the isolates K8108, MT2006 and UFV02.

#### 661 **Prediction and annotation of secreted proteins**

662 To predict classically secreted proteins, we initially searched for proteins containing a classic  
663 signal peptide and no transmembrane signal using SignalP (versions 3 and 5) (34), TMHMM  
664 (131) and Phobius (132) programs. For the identification of additional secreted proteins  
665 without a classic peptide signal and no transmembrane signal (non-classically secreted), we  
666 used EffectorP (versions 1 and 2) (31, 32) and TMHMM programs. In both approaches, we  
667 kept the proteins having a TM in the N-term region. The proteins selected by both approaches  
668 were analysed by PS-SCAN program (133) to remove putative endoplasmatic reticulum  
669 proteins. All programs were performed considering default parameters. The secreted  
670 proteins predicted in the previous step were annotated using Blast (134), RPSBlast, PredGPI  
671 (135), InterProScan (136) and hmmsearch (137) programs. Similarity searches using Blast  
672 program were performed against the NCBI non-redundant (nr), FunSecKb (138), Phi-base  
673 (139) and LED (140) databases, applying an e-value of  $10^{-5}$ . To search for domains in  
674 sequences, we used the programs RPSBlast and hmmsearch against the Conserved Domain  
675 Database (CDD) (141) and PFAM database (142) respectively, using an e-value of  $10^{-5}$  in both  
676 cases. Ortholog mapping was done through similarity searches with the hmmsearch program  
677 against profile HMMs obtained from eggNOG database (143). To predict the location of the  
678 predicted proteins in plant, ApoplastP (144), Localizer (145), targetP (146), WoLFPSORT (147)  
679 and DeepLoc (148) programs were performed using default parameters. To assign a final  
680 location for each protein, the following criteria were considered: if at least two programs  
681 found the same result, that result was considered as a predicted location. Otherwise, the term  
682 "Not classified" was assigned to the protein. To identify the motifs [Y/F/W]xC in the  
683 sequences, we used a proprietary script developed in Perl language. A summary of the  
684 prediction and annotation pipelines for the secreted proteins are illustrated in figs. S21 and  
685 S22.

686 For the prediction of putative effector proteins, we used the list of predicted secreted  
687 proteins containing a classical signal peptide. For the prediction of candidate effector proteins  
688 in each genome, we defined three different approaches. In the first one, sequences predicted  
689 as "Extracellular" or "Not Classified" by the location programs and with no annotation were  
690 selected as candidates to effector proteins. With this approach, we obtained 618, 531 and  
691 598 candidates to effector proteins in K8108, MT2006 and UFV02 respectively. In the second  
692 approach, we selected proteins which contain PFAM domains present in effector proteins  
693 (149). Applying this criterion, we selected 142, 128 and 55 candidates in K8108, MT2006 and  
694 UFV02, respectively. Finally, in the third approach, we ran EffectorP program to classify the  
695 effector candidates, and we obtained 802, 851 and 899 candidates in K8108, MT2006 and  
696 UFV02 genome respectively (tables S6 to S8).

#### 697 **Staining of leaf samples and microscopy**

698 Plants were inoculated by spray inoculation and leaves harvested at the indicated time points.  
699 Samples were destained in 1M KOH with 0.01% Silwet L-77 (Sigma Aldrich) for at least 12 h at  
700 37°C and stored in 50 mM Tris-HCl pH 7.5 at 4°C. Fungal staining was obtained with wheat  
701 germ agglutinin (WGA) FITC conjugate (Merck L4895), samples were incubated 30 min to  
702 overnight in a 20 µg/ml solution in Tris-HCl pH 7.5. Co-staining of plant tissue with propidium  
703 iodide (Sigma-Aldrich P4864) was performed according to the manufacturer's instructions.  
704 Images were obtained with a Leica SP5 confocal microscope (Leica Microsystems) with an



705 excitation of 488 nm and detection at 500-550 nm and 625-643 nm, respectively. Z-stacks  
706 were opened in the 3D viewer of the LAS X software (Leica Application Suite X 3.5.7.23225)  
707 and resulting images exported. Clipping was performed as indicated in the pictures. Shading  
708 was performed in some cases for better visualization.

709 For cryo-scanning electron microscopy, inoculated soybean leaves were cut and  
710 mounted on an aluminium stub with Tissue Tek OCT (Agar Scientific Ltd, Essex, UK) and plunge  
711 frozen in slushed liquid nitrogen to cryo-preserve the material before transfer to the cryo-  
712 stage of a PP3010 cryo-SEM preparation system (Quorum Technologies, Laughton, UK)  
713 attached to a Zeiss Gemini 300 field emission gun scanning electron microscope (Zeiss UK Ltd,  
714 Cambridge, UK). Surface frost was sublimated by warming the sample to -90 °C for 4 minutes,  
715 before the sample was cooled to -140 °C and sputter coated with platinum for 50 seconds at  
716 5 mA. The sample was loaded onto the cryo-stage of the main SEM chamber and held at -140  
717 °C during imaging at 3 kV using an Everhart-Thornley detector. False colouring of images was  
718 performed with Adobe Photoshop 22.4.2.  
719

**Table 1: *P. pachyrhizi* genome assembly metrics.**

	K8108	MT2006	UFV02
Assembly size (Gb)	1.083	1.0574	1.273
Total no of contigs	6,505	7,464	3,140
Contig N50 length (Kb)	278.753	222.464	677.464
Max contig length (Mb)	3.028	3.054	4.158
Min contig length (Kb)	16.399	21.118	11.733
Complete BUSCOs (%)	90.19	90.14	89.91
Complete single-copy BUSCO (%)	15.70	15.87	22.56
Complete duplicated BUSCO (%)	74.49	74.26	67.35
Fragmented BUSCO (%)	1.36	1.36	1.19
Missing BUSCO (%)	8.45	8.50	8.90
Total BUSCO	1,764	1,764	1,764

720

721

**Table 2: Expansion of gene families in the *P. pachyrhizi* genome.**

	Piwi	KOG 057	KOG 148	KOG 241	KOG 039	KOG 246	KOG 068	KOG 261	KOG 126	KOG 149
		3	1	0	9	7	3	7	1	4
<i>P. pachyrhizi</i> UFV02	531	78	28	62	48	12	10	15	26	13
<i>P. pachyrhizi</i> MT2006	568	77	25	22	44	8	5	12	29	8
<i>P. pachyrhizi</i> K8108	608	74	34	78	18	11	8	11	24	13
<i>C. quercuum</i> f. sp. <i>fusiforme</i> G11	3	1	2	3	2	1	3	2	1	2
<i>M. larici-</i> <i>populina</i>	3	1	2	2	2	5	4	2	1	3
<i>M. allii-</i> <i>populina</i> 12AY07	6	1	3	3	2	1	5	2	1	2
<i>P. graminis</i> f. sp. <i>tritici</i>	3	1	2	2	2	2	3	2	1	2
<i>P. striiformis</i> f. sp. <i>tritici</i> 104 E137 A-	7	2	5	4	2	4	8	4	3	4
<i>P. 17oronate</i> <i>avenae</i> 12SD80	5	2	4	2	8	4	5	5	2	2
<i>P. triticina</i> 1-1 BBBD Race 1	3	2	3	2	1	2	5	2	1	2

722

723

724

725 **DATA AVAILABILITY**

726 The raw sequencing data of MT2006, K8108 and UFV02 isolates has been deposited at NCBI  
727 under the accession numbers PRJNA368291, PRJEB46918, and PRJEB44222, respectively.

728

729 **CONFLICT OF INTEREST**

730 Connor Cameron, Andrew Farmer, Dirk Balmer, Stephanie Widdison, Qingli Liu and Gabriel  
731 Scalliet were employees of Syngenta or affiliates during the course of the research project.  
732 Work on the soybean isolate K8108 in the Conrath and Schaffrath lab was supported, in part,  
733 by Syngenta Crop Protection.

734 2Blades has two collaborations with Bayer crop science on Asian soybean rust.

735

736 **CONTRIBUTIONS**

737 Y.K.G, F.C.M.G., C.L., A.F., S.H., E.G.C.F, V.S.L., L.S.O., E.M., S.W., C.C., Y.I., K.T., K.R., E.D., B.H.,  
738 K.L., A.M.R.B., E.P., V.S., C.D., C.D., M.v.H, A.J., L.C., Y.T., J.R., B.d.V.A.M., A.W., H.S., S.P.,  
739 L.G.Z., V.C.H., F.C., T.I.L., D.B., A.M., S.K., S.B., L.W., C.C, M.Y., Q.L., M.L., S.H.B., and S.D.  
740 performed research. Y.K.G, C.L., A.F., S.H., E.G.C.F, V.S.L., L.S.O., A.M.R.B., E.M., S.W., C.C.,  
741 Y.I., E.D., B.H., A.J., A.W., B.d.V.A.M., L.G.Z., T.I.L., M.L., S.H.B., and S.D. analyzed the data.  
742 Y.K.G., F.C.M.G., M.L, S.D., and H.P.v.E. edited the manuscript. Y.K.G., F.C.M.G., V.S.L., L.S.O.,  
743 M.L., U.S., S.D., and H.P.v.E. wrote the paper. F.C.M.G., V.N., P.G., R.T.V., I.V.G., U.C., G.S.,  
744 C.S., S.D., and H.P.v.E. directed aspects of the project. For detail see table S28.

745

746 **ACKNOWLEDGEMENTS**

747 Sequencing and RNAseq analyses of UFV02 was supported by 2Blades. We thank Dan  
748 MacLean and Ram Krishna Shrestha for bioinformatics support. Bioinformatics infrastructure  
749 was supported in part by NBI Research Computing. We thank Matthew Moscou for many  
750 fruitful discussions. We acknowledge Heike Popovitsch for technical support. The work  
751 (Proposal 10.46936/10.25585/60000959) conducted by the U.S. Department of Energy Joint  
752 Genome Institute (<https://ror.org/04xm1d337>), a DOE Office of Science User Facility, is  
753 supported by the Office of Science of the U.S. Department of Energy under Contract No. DE-  
754 AC02-05CH11231. We thank Robert Dietrich and Lucio Garcia (Syngenta RTP) for their  
755 technical support with the sequencing of K8108 genome and transcriptome.

756

757

758 **REFERENCES:**

- 759 1. S. Savary *et al.*, The global burden of pathogens and pests on major food crops. *Nature*  
760 *Ecology & Evolution* **3**, 430-439 (2019).
- 761 2. H. Scherm, R. S. C. Christiano, P. D. Esker, E. M. Del Ponte, C. V. Godoy, Quantitative  
762 review of fungicide efficacy trials for managing soybean rust in Brazil. *Crop Protection*  
763 **28**, 774-782 (2009).
- 764 3. J. T. Yorinori *et al.*, Epidemics of Soybean Rust (*Phakopsora pachyrhizi*) in Brazil and  
765 Paraguay from 2001 to 2003. *Plant Dis* **89**, 675-677 (2005).
- 766 4. E. Melo Reis, E. Deuner, M. Zanatta, *In vivo* sensitivity of *Phakopsora pachyrhizi* to DMI  
767 and QoI fungicides. *Summa Phytopathologica* **41**, 21-24 (2015).
- 768 5. H. Akamatsu *et al.*, Pathogenic diversity of soybean rust in Argentina, Brazil, and  
769 Paraguay. *Journal of General Plant Pathology* **79**, 28-40 (2013).

- 770 6. C. Paul, G. L. Hartman, J. J. Marois, D. L. Wright, D. R. Walker, First report of  
771 *Phakopsora pachyrhizi* adapting to soybean genotypes with Rpp1 or Rpp6 rust  
772 resistance genes in field plots in the United States. *Plant Disease* **97**, 1379-1379 (2013).  
773 7. C. V. Godoy *et al.*, Asian soybean rust in Brazil: past, present, and future. *Pesquisa*  
774 *Agropecuária Brasileira* **51**, 407-421 (2016).  
775 8. M. A. Müller, G. Stammler, L. L. May De Mio, Multiple resistance to DMI, Qol and SDHI  
776 fungicides in field isolates of *Phakopsora pachyrhizi*. *Crop Protection* **145**, 105618  
777 (2021).  
778 9. J. P. Barro *et al.*, Performance of dual and triple fungicide premixes for managing  
779 soybean rust across years and regions in Brazil: A meta-analysis. *Plant Pathology* **70**,  
780 1920-1935 (2021).  
781 10. Y. Ono, P. Buritica, J. F. Hennen, Delimitation of *Phakopsora*, *Physopella* and  
782 *Cerotelium* and their species on Leguminosae. *Mycological Research* **96**, 825-850  
783 (1992).  
784 11. M. R. Bonde *et al.*, Comparative susceptibilities of legume species to infection by  
785 *Phakopsora pachyrhizi*. *Plant Disease* **92**, 30-36 (2008).  
786 12. T. L. Slaminko, M. R. Miles, R. D. Frederick, M. R. Bonde, G. L. Hartman, New legume  
787 hosts of *Phakopsora pachyrhizi* based on greenhouse evaluations. *Plant Disease* **92**,  
788 767-771 (2008).  
789 13. C. L. Harmon, P. F. Harmon, T. A. Mueller, J. J. Marois, G. L. Hartman, First report of  
790 *Phakopsora pachyrhizi* telia on kudzu in the United States. *Plant Disease* **90**, 380-380  
791 (2006).  
792 14. M. Loehrer *et al.*, On the current status of *Phakopsora pachyrhizi* genome sequencing.  
793 *Front Plant Sci* **5**, 377-377 (2014).  
794 15. F. Li *et al.*, Emergence of the Ug99 lineage of the wheat stem rust pathogen through  
795 somatic hybridisation. *Nature Communications* **10**, (2019).  
796 16. B. Schwessinger *et al.*, A near-complete haplotype-phased genome of the dikaryotic  
797 wheat stripe rust fungus *Puccinia striiformis* f. sp. *tritici* reveals high interhaplotype  
798 diversity. *mBio* **9**, e02275-02217 (2018).  
799 17. M. E. Miller *et al.*, *De Novo* assembly and phasing of dikaryotic genomes from two  
800 isolates of *Puccinia coronata* f. sp. *avenae*, the causal agent of oat crown rust. *mBio* **9**,  
801 e01650-01617 (2018).  
802 18. U. Oggenfuss *et al.*, A population-level invasion by transposable elements triggers  
803 genome expansion in a fungal pathogen. *eLife* **10**, e69249 (2021).  
804 19. P. A. Tobias *et al.*, *Austropuccinia psidii*, causing myrtle rust, has a gigabase-sized  
805 genome shaped by transposable elements. *G3 (Bethesda)* **11**, (2020).  
806 20. F. Maumus, H. Quesneville, Ancestral repeats have shaped epigenome and genome  
807 composition for millions of years in *Arabidopsis thaliana*. *Nature Communications* **5**,  
808 4104 (2014).  
809 21. R. Castanera *et al.*, Transposable elements versus the fungal genome: impact on  
810 whole-Genome architecture and transcriptional profiles. *PLOS Genetics* **12**, e1006108  
811 (2016).  
812 22. B. Dhillon, N. Gill, R. C. Hamelin, S. B. Goodwin, The landscape of transposable  
813 elements in the finished genome of the fungal wheat pathogen *Mycosphaerella*  
814 *graminicola*. *BMC Genomics* **15**, 1132 (2014).  
815 23. J. Schmutz *et al.*, Genome sequence of the palaeopolyploid soybean. *Nature* **463**, 178-  
816 183 (2010).

- 817 24. Z.-Q. Shao *et al.*, Long-term evolution of nucleotide-binding site-leucine-rich repeat  
818 genes: understanding gained from and beyond the Legume Family. *Plant Physiology*  
819 **166**, 217-234 (2014).
- 820 25. F. Zheng *et al.*, Molecular phylogeny and dynamic evolution of disease resistance  
821 genes in the legume family. *BMC Genomics* **17**, 402 (2016).
- 822 26. M. Lavin, P. S. Herendeen, M. F. Wojciechowski, Evolutionary rates analysis of  
823 Leguminosae implicates a rapid diversification of lineages during the tertiary. *Syst Biol*  
824 **54**, 575-594 (2005).
- 825 27. G. Hewitt, The genetic legacy of the Quaternary ice ages. *Nature* **405**, 907-913 (2000).
- 826 28. I. V. Grigoriev *et al.*, MycoCosm portal: gearing up for 1000 fungal genomes. *Nucleic*  
827 *Acids Res* **42**, D699-704 (2014).
- 828 29. P. N. Dodds, J. P. Rathjen, Plant immunity: towards an integrated view of plant-  
829 pathogen interactions. *Nat Rev Genet* **11**, 539-548 (2010).
- 830 30. R. de Jonge, M. D. Bolton, B. P. Thomma, How filamentous pathogens co-opt plants:  
831 the ins and outs of fungal effectors. *Curr Opin Plant Biol* **14**, 400-406 (2011).
- 832 31. J. Sperschneider, P. N. Dodds, D. M. Gardiner, K. B. Singh, J. M. Taylor, Improved  
833 prediction of fungal effector proteins from secretomes with EffectorP 2.0. *Molecular*  
834 *Plant Pathology* **19**, 2094-2110 (2018).
- 835 32. J. Sperschneider *et al.*, EffectorP: predicting fungal effector proteins from secretomes  
836 using machine learning. *New Phytologist* **210**, 743-761 (2016).
- 837 33. L. Käll, A. Krogh, E. L. Sonnhammer, A combined transmembrane topology and signal  
838 peptide prediction method. *J Mol Biol* **338**, 1027-1036 (2004).
- 839 34. J. J. Almagro Armenteros *et al.*, SignalP 5.0 improves signal peptide predictions using  
840 deep neural networks. *Nature Biotechnology* **37**, 420-423 (2019).
- 841 35. J. Dyrlov Bendtsen, H. Nielsen, G. von Heijne, S. Brunak, Improved prediction of signal  
842 peptides: SignalP 3.0. *Journal of Molecular Biology* **340**, 783-795 (2004).
- 843 36. T. I. Link *et al.*, The haustorial transcriptomes of *Uromyces appendiculatus* and  
844 *Phakopsora pachyrhizi* and their candidate effector families. *Mol Plant Pathol* **15**, 379-  
845 393 (2014).
- 846 37. S. G. Kunjeti *et al.*, Identification of *Phakopsora pachyrhizi* candidate effectors with  
847 virulence activity in a distantly related pathosystem. *Front Plant Sci* **7**, 269-269 (2016).
- 848 38. M. C. de Carvalho *et al.*, Prediction of the in planta *Phakopsora pachyrhizi* secretome  
849 and potential effector families. *Mol Plant Pathol* **18**, 363-377 (2017).
- 850 39. M. Qi *et al.*, Suppression or activation of immune responses by predicted secreted  
851 proteins of the soybean rust pathogen *Phakopsora pachyrhizi*. *Mol Plant Microbe*  
852 *Interact* **31**, 163-174 (2018).
- 853 40. S. Fouché *et al.*, Stress-driven transposable element de-repression dynamics and  
854 virulence evolution in a fungal pathogen. *Molecular Biology and Evolution* **37**, 221-239  
855 (2019).
- 856 41. S. Fouché, U. Oggenfuss, E. Chanclud, D. Croll, A devil's bargain with transposable  
857 elements in plant pathogens. *Trends Genet*, (2021).
- 858 42. D. E. Torres, B. P. H. J. Thomma, M. F. Seidl, Transposable elements contribute to  
859 genome dynamics and gene expression variation in the fungal plant pathogen  
860 *Verticillium dahliae*. *Genome Biology and Evolution* **13**, (2021).
- 861 43. S. Raffaele *et al.*, Genome evolution following host jumps in the Irish potato famine  
862 pathogen lineage. *Science* **330**, 1540-1543 (2010).

- 863 44. H. C. van der Does, M. Rep, Virulence genes and the evolution of host specificity in  
864 plant-pathogenic fungi. *Mol Plant Microbe Interact* **20**, 1175-1182 (2007).
- 865 45. J. Li, L. Fokkens, L. J. Conneely, M. Rep, Partial pathogenicity chromosomes in *Fusarium*  
866 *oxysporum* are sufficient to cause disease and can be horizontally transferred. *Environ*  
867 *Microbiol* **22**, 4985-5004 (2020).
- 868 46. R. Harting *et al.*, A 20-kb lineage-specific genomic region tames virulence in  
869 pathogenic amphidiploid *Verticillium longisporum*. *Molecular Plant Pathology* **22**, 939-  
870 953 (2021).
- 871 47. R. de Jonge *et al.*, Tomato immune receptor Ve1 recognizes effector of multiple fungal  
872 pathogens uncovered by genome and RNA sequencing. *PNAS* **109**, 5110-5115 (2012).
- 873 48. D. Croll, B. A. McDonald, The accessory genome as a cradle for adaptive evolution in  
874 pathogens. *PLoS Pathog* **8**, e1002608 (2012).
- 875 49. S. M. Schmidt *et al.*, MITEs in the promoters of effector genes allow prediction of novel  
876 virulence genes in *Fusarium oxysporum*. *BMC Genomics* **14**, 119 (2013).
- 877 50. R. de Jonge *et al.*, Extensive chromosomal reshuffling drives evolution of virulence in  
878 an asexual pathogen. *Genome Res* **23**, 1271-1282 (2013).
- 879 51. C. Lorrain, K. C. Gonçalves dos Santos, H. Germain, A. Hecker, S. Duplessis, Advances  
880 in understanding obligate biotrophy in rust fungi. *New Phytologist* **222**, 1190-1206  
881 (2019).
- 882 52. O. P. Judson, B. B. Normark, Ancient asexual scandals. *Trends Ecol Evol* **11**, 41-46  
883 (1996).
- 884 53. F. Balloux, L. Lehmann, T. de Meeûs, The population genetics of clonal and partially  
885 clonal diploids. *Genetics* **164**, 1635-1644 (2003).
- 886 54. B. Schwessinger *et al.*, Distinct life histories impact dikaryotic genome evolution in the  
887 rust fungus *Puccinia striiformis* causing stripe rust in wheat. *Genome Biology and*  
888 *Evolution* **12**, 597-617 (2020).
- 889 55. V. R. Jorge *et al.*, The origin and genetic diversity of the causal agent of Asian soybean  
890 rust, *Phakopsora pachyrhizi*, in South America. *Plant Pathology* **64**, 729-737 (2015).
- 891 56. L. M. Darben *et al.*, Characterization of genetic diversity and pathogenicity of  
892 *Phakopsora pachyrhizi* mono-uredinial isolates collected in Brazil. *European Journal of*  
893 *Plant Pathology* **156**, 355-372 (2020).
- 894 57. M. Nattestad, M. C. Schatz, Assemblytics: a web analytics tool for the detection of  
895 variants from an assembly. *Bioinformatics* **32**, 3021-3023 (2016).
- 896 58. K. Goellner *et al.*, *Phakopsora pachyrhizi*, the causal agent of Asian soybean rust.  
897 *Molecular plant pathology* **11**, 169-177 (2010).
- 898 59. S. A. Isard, S. H. Gage, P. Comtois, J. M. Russo, Principles of the atmospheric pathway  
899 for invasive species applied to soybean rust. *BioScience* **55**, 851-861 (2005).
- 900 60. W. Zheng *et al.*, High genome heterozygosity and endemic genetic recombination in  
901 the wheat stripe rust fungus. *Nature Communications* **4**, 2673 (2013).
- 902 61. J. Chen *et al.*, *De novo* genome assembly and comparative genomics of the barley leaf  
903 rust pathogen *Puccinia hordei* identifies candidates for three avirulence genes. *G3*  
904 *(Bethesda)* **9**, 3263-3271 (2019).
- 905 62. C. A. Cuomo *et al.*, Comparative analysis highlights variable genome content of wheat  
906 rusts and divergence of the mating loci. *G3 (Bethesda)* **7**, 361-376 (2017).
- 907 63. A. R. McTaggart *et al.*, Host jumps shaped the diversity of extant rust fungi  
908 (Pucciniales). *New Phytologist* **209**, 1149-1158 (2016).

- 909 64. M. C. Aime, C. D. Bell, A. W. Wilson, Deconstructing the evolutionary complexity  
910 between rust fungi (Pucciniales) and their plant hosts. *Studies in Mycology* **89**, 143-  
911 152 (2018).
- 912 65. M. V. Han, G. W. Thomas, J. Lugo-Martinez, M. W. Hahn, Estimating gene gain and loss  
913 rates in the presence of error in genome assembly and annotation using CAFE 3. *Mol*  
914 *Biol Evol* **30**, 1987-1997 (2013).
- 915 66. W. Zhou *et al.*, Glutamate synthase MoGlt1-mediated glutamate homeostasis is  
916 important for autophagy, virulence and conidiation in the rice blast fungus. *Mol Plant*  
917 *Pathol* **19**, 564-578 (2018).
- 918 67. V. Lombard, H. Golaconda Ramulu, E. Drula, P. M. Coutinho, B. Henrissat, The  
919 carbohydrate-active enzymes database (CAZy) in 2013. *Nucleic Acids Res* **42**, D490-495  
920 (2014).
- 921 68. O. Etxebeste *et al.*, GmcA is a putative Glucose-Methanol-Choline Oxidoreductase  
922 required for the induction of asexual development in *Aspergillus nidulans*. *PLOS ONE*  
923 **7**, e40292 (2012).
- 924 69. W. Chen, X. Jiang, Q. Yang, Glycoside hydrolase family 18 chitinases: The known and  
925 the unknown. *Biotechnology Advances* **43**, 107553 (2020).
- 926 70. T. Langner, V. Göhre, Fungal chitinases: function, regulation, and potential roles in  
927 plant/pathogen interactions. *Curr Genet* **62**, 243-254 (2016).
- 928 71. M. C. Aime, A. R. McTaggart, A higher-rank classification for rust fungi, with notes on  
929 genera. *Fungal Syst Evol* **7**, 21-47 (2021).
- 930 72. N. Darricarrère, N. Liu, T. Watanabe, H. Lin, Function of Piwi, a nuclear Piwi/Argonaute  
931 protein, is independent of its slicer activity. *PNAS* **110**, 1297-1302 (2013).
- 932 73. K. Saito *et al.*, Specific association of Piwi with rasiRNAs derived from retrotransposon  
933 and heterochromatic regions in the *Drosophila* genome. *Genes Dev* **20**, 2214-2222  
934 (2006).
- 935 74. G. Bourque *et al.*, Ten things you should know about transposable elements. *Genome*  
936 *Biology* **19**, 199 (2018).
- 937 75. T. Thomson, H. Lin, The biogenesis and function of PIWI proteins and piRNAs: progress  
938 and prospect. *Annu Rev Cell Dev Biol* **25**, 355-376 (2009).
- 939 76. M. D. Bolton, Primary metabolism and plant defense--fuel for the fire. *Mol Plant*  
940 *Microbe Interact* **22**, 487-497 (2009).
- 941 77. L. Schrader, J. Schmitz, The impact of transposable elements in adaptive evolution.  
942 *Molecular Ecology* **28**, 1537-1549 (2019).
- 943 78. M. F. Seidl, B. Thomma, Transposable elements direct the coevolution between plants  
944 and microbes. *Trends Genet* **33**, 842-851 (2017).
- 945 79. I. K. Jordan, N. J. Bowen, Computational analysis of transposable element sequences.  
946 *Methods Mol Biol* **260**, 59-71 (2004).
- 947 80. C. Lorrain, A. Feurtey, M. Möller, J. Haueisen, E. Stukenbrock, Dynamics of  
948 transposable elements in recently diverged fungal pathogens: lineage-specific  
949 transposable element content and efficiency of genome defenses. *G3 (Bethesda)* **11**,  
950 (2021).
- 951 81. C. G. Kawashima *et al.*, A pigeonpea gene confers resistance to Asian soybean rust in  
952 soybean. *Nature Biotechnology* **34**, 661-665 (2016).
- 953 82. B. Mayjonade *et al.*, Extraction of high-molecular-weight genomic DNA for long-read  
954 sequencing of single molecules. *BioTechniques* **61**, 203-205 (2016).

- 955 83. A. Persoons *et al.*, Patterns of genomic variation in the poplar rust fungus *Melampsora*  
956 *larici-populina* identify pathogenesis-related factors. *Front Plant Sci* **5**, (2014).
- 957 84. B. Schwessinger, J. P. Rathjen, in *Wheat Rust Diseases: Methods and Protocols*, S.  
958 Periyannan, Ed. (Springer New York, New York, NY, 2017), pp. 49-57.
- 959 85. C.-L. Xiao *et al.*, MECAT: fast mapping, error correction, and de novo assembly for  
960 single-molecule sequencing reads. *Nature Methods* **14**, 1072-1074 (2017).
- 961 86. S. Koren *et al.*, Canu: scalable and accurate long-read assembly via adaptive k-mer  
962 weighting and repeat separation. *Genome Research* **27**, 722-736 (2017).
- 963 87. C.-S. Chin *et al.*, Nonhybrid, finished microbial genome assemblies from long-read  
964 SMRT sequencing data. *Nature Methods* **10**, 563-569 (2013).
- 965 88. C.-S. Chin *et al.*, Phased diploid genome assembly with single-molecule real-time  
966 sequencing. *Nature Methods* **13**, 1050-1054 (2016).
- 967 89. K.-K. Lam, K. LaButti, A. Khalak, D. Tse, FinisherSC: a repeat-aware tool for upgrading  
968 de novo assembly using long reads. *Bioinformatics* **31**, 3207-3209 (2015).
- 969 90. H. Li, Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**,  
970 3094-3100 (2018).
- 971 91. R. Vaser, I. Sovic, N. Nagarajan, M. Sikic, Fast and accurate de novo genome assembly  
972 from long uncorrected reads. *Genome Res* **27**, 737-746 (2017).
- 973 92. H. Li, R. Durbin, Fast and accurate short read alignment with Burrows-Wheeler  
974 transform. *Bioinformatics* **25**, 1754-1760 (2009).
- 975 93. B. J. Walker *et al.*, Pilon: an integrated tool for comprehensive microbial variant  
976 detection and genome assembly improvement. *PLoS ONE* **9**, e112963 (2014).
- 977 94. T. Flutre, E. Duprat, C. Feuillet, H. Quesneville, Considering transposable element  
978 diversification in *de novo* annotation approaches. *PLOS ONE* **6**, e16526 (2011).
- 979 95. H. Quesneville *et al.*, Combined evidence annotation of transposable elements in  
980 genome sequences. *PLOS Computational Biology* **1**, e22 (2005).
- 981 96. V. Jamilloux, J. Daron, F. Choulet, H. Quesneville, *De Novo* annotation of transposable  
982 elements: tackling the fat genome issue. *Proceedings of the IEEE* **105**, 978-978 (2017).
- 983 97. R. C. Edgar, E. W. Myers, PILER: identification and classification of genomic repeats.  
984 *Bioinformatics* **21**, i152-158 (2005).
- 985 98. Z. Bao, S. R. Eddy, Automated de novo identification of repeat sequence families in  
986 sequenced genomes. *Genome Res* **12**, 1269-1276 (2002).
- 987 99. X. Huang, On global sequence alignment. *Bioinformatics* **10**, 227-235 (1994).
- 988 100. T. Wicker *et al.*, A unified classification system for eukaryotic transposable elements.  
989 *Nature Reviews Genetics* **8**, 973-982 (2007).
- 990 101. C. Hoede *et al.*, PASTEC: an automatic transposable element classification tool. *PLoS*  
991 *One* **9**, e91929 (2014).
- 992 102. O. Kohany, A. J. Gentles, L. Hankus, J. Jurka, Annotation, submission and screening of  
993 repetitive elements in Repbase: RepbaseSubmitter and Censor. *BMC Bioinformatics* **7**,  
994 474 (2006).
- 995 103. A. R. Quinlan, I. M. Hall, BEDTools: a flexible suite of utilities for comparing genomic  
996 features. *Bioinformatics* **26**, 841-842 (2010).
- 997 104. D. Ellinghaus, S. Kurtz, U. Willhoeft, LTRharvest, an efficient and flexible software for  
998 de novo detection of LTR retrotransposons. *BMC Bioinformatics* **9**, 18 (2008).
- 999 105. J. Jurka, Repbase Update: a database and an electronic journal of repetitive elements.  
1000 *Trends in Genetics* **16**, 418-420 (2000).



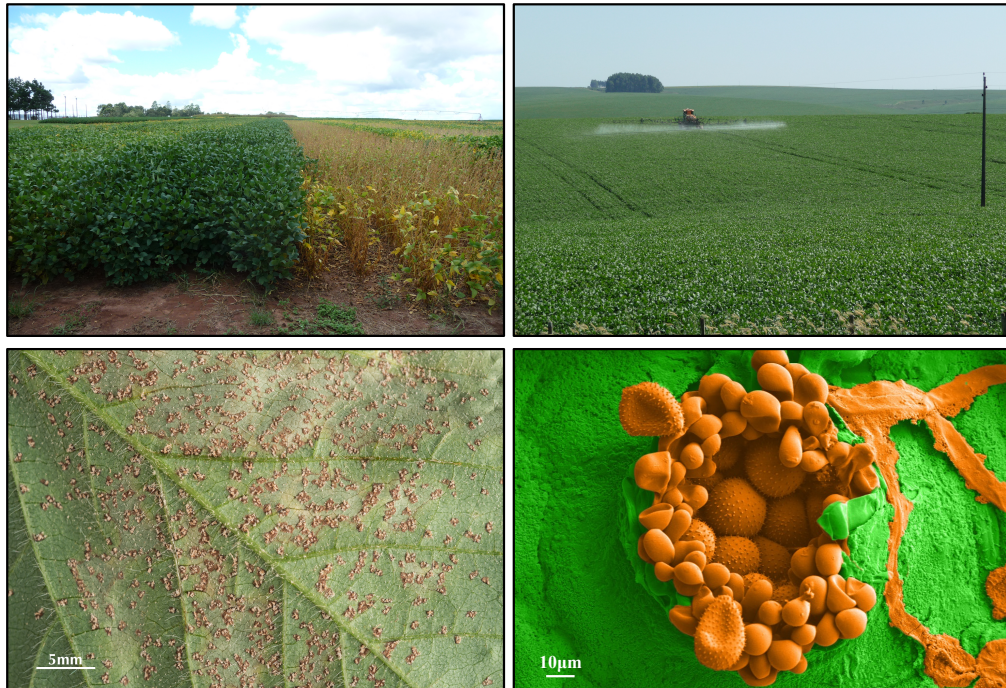
- 1001 106. K. Katoh, D. M. Standley, MAFFT multiple sequence alignment software version 7:  
1002 improvements in performance and usability. *Molecular Biology and Evolution* **30**, 772-  
1003 780 (2013).
- 1004 107. M. Kimura, A simple method for estimating evolutionary rates of base substitutions  
1005 through comparative studies of nucleotide sequences. *J Mol Evol* **16**, 111-120 (1980).
- 1006 108. M. Prieto, M. Wedin, Dating the diversification of the major lineages of ascomycota  
1007 (fungi). *PLOS ONE* **8**, e65576 (2013).
- 1008 109. A. Kuo, B. Bushnell, I. V. Grigoriev, in *Advances in Botanical Research*, F. M. Martin, Ed.  
1009 (Academic Press, 2014), vol. 70, pp. 1-52.
- 1010 110. F. A. Simao, R. M. Waterhouse, P. Ioannidis, E. V. Kriventseva, E. M. Zdobnov, BUSCO:  
1011 assessing genome assembly and annotation completeness with single-copy orthologs.  
1012 *Bioinformatics* **31**, 3210-3212 (2015).
- 1013 111. D. Mapleson, G. Garcia Accinelli, G. Kettleborough, J. Wright, B. J. Clavijo, KAT: a K-mer  
1014 analysis toolkit to quality control NGS datasets and genome assemblies.  
1015 *Bioinformatics* **33**, 574-576 (2016).
- 1016 112. T. R. Ranallo-Benavidez, K. S. Jaron, M. C. Schatz, GenomeScope 2.0 and Smudgeplots:  
1017 Reference-free profiling of polyploid genomes. *bioRxiv*, 747568 (2019).
- 1018 113. M. J. Roach, S. A. Schmidt, A. R. Borneman, Purge Haplotigs: allelic contig  
1019 reassignment for third-gen diploid genome assemblies. *BMC Bioinformatics* **19**, 460  
1020 (2018).
- 1021 114. G. Marçais *et al.*, MUMmer4: A fast and versatile genome alignment system. *PLOS*  
1022 *Computational Biology* **14**, e1005944 (2018).
- 1023 115. Y. Wang *et al.*, MCScanX: a toolkit for detection and evolutionary analysis of gene  
1024 synteny and collinearity. *Nucleic Acids Research* **40**, e49-e49 (2012).
- 1025 116. A. Shumate, S. L. Salzberg, Liftoff: accurate mapping of gene annotations.  
1026 *Bioinformatics* **37**, 1639-1643 (2021).
- 1027 117. A. M. Bolger, M. Lohse, B. Usadel, Trimmomatic: a flexible trimmer for Illumina  
1028 sequence data. *Bioinformatics* **30**, 2114-2120 (2014).
- 1029 118. H. Li *et al.*, The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**,  
1030 2078-2079 (2009).
- 1031 119. A. McKenna *et al.*, The Genome Analysis Toolkit: a MapReduce framework for  
1032 analyzing next-generation DNA sequencing data. *Genome Res* **20**, 1297-1303 (2010).
- 1033 120. P. Cingolani *et al.*, A program for annotating and predicting the effects of single  
1034 nucleotide polymorphisms, SnpEff: SNPs in the genome of *Drosophila melanogaster*  
1035 strain w1118; iso-2; iso-3. *Fly (Austin)* **6**, 80-92 (2012).
- 1036 121. M. Loehrer *et al.*, *In vivo* assessment by Mach-Zehnder double-beam interferometry  
1037 of the invasive force exerted by the Asian soybean rust fungus (*Phakopsora*  
1038 *pachyrhizi*). *New Phytologist* **203**, 620-631 (2014).
- 1039 122. A. Heller, Host-parasite interaction during subepidermal sporulation and pustule  
1040 opening in rust fungi (Pucciniales). *Protospasma* **257**, 783-792 (2020).
- 1041 123. D. E. Wood, J. Lu, B. Langmead, Improved metagenomic analysis with Kraken 2.  
1042 *Genome Biology* **20**, 257 (2019).
- 1043 124. A. Dobin *et al.*, STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**, 15-21  
1044 (2012).
- 1045 125. S. Kovaka *et al.*, Transcriptome assembly from long-read RNA-seq alignments with  
1046 StringTie2. *Genome Biology* **20**, 278 (2019).
- 1047 126. G. Pertea, M. Pertea, GFF Utilities: GffRead and GffCompare. *F1000Res* **9**, (2020).

- 1048 127. Y. Jin, O. H. Tam, E. Paniagua, M. Hammell, TETranscripts: a package for including  
1049 transposable elements in differential expression analysis of RNA-seq datasets.  
1050 *Bioinformatics* **31**, 3593-3599 (2015).
- 1051 128. M. D. Robinson, D. J. McCarthy, G. K. Smyth, edgeR: a Bioconductor package for  
1052 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139-  
1053 140 (2010).
- 1054 129. D. J. McCarthy, Y. Chen, G. K. Smyth, Differential expression analysis of multifactor  
1055 RNA-Seq experiments with respect to biological variation. *Nucleic Acids Res* **40**, 4288-  
1056 4297 (2012).
- 1057 130. M. D. Robinson, A. Oshlack, A scaling normalization method for differential expression  
1058 analysis of RNA-seq data. *Genome Biology* **11**, R25 (2010).
- 1059 131. A. Krogh, B. Larsson, G. von Heijne, E. L. Sonnhammer, Predicting transmembrane  
1060 protein topology with a hidden Markov model: application to complete genomes. *J*  
1061 *Mol Biol* **305**, 567-580 (2001).
- 1062 132. L. Käll, A. Krogh, E. L. Sonnhammer, Advantages of combined transmembrane  
1063 topology and signal peptide prediction--the Phobius web server. *Nucleic Acids Res* **35**,  
1064 W429-432 (2007).
- 1065 133. A. Gattiker, E. Gasteiger, A. Bairoch, ScanProsite: A reference implementation of a  
1066 PROSITE scanning tool. *Applied bioinformatics* **1**, 107-108 (2002).
- 1067 134. S. F. Altschul, W. Gish, W. Miller, E. W. Myers, D. J. Lipman, Basic local alignment  
1068 search tool. *J Mol Biol* **215**, 403-410 (1990).
- 1069 135. A. Pierleoni, P. L. Martelli, R. Casadio, PredGPI: a GPI-anchor predictor. *BMC*  
1070 *Bioinformatics* **9**, 392 (2008).
- 1071 136. M. Blum *et al.*, The InterPro protein families and domains database: 20 years on.  
1072 *Nucleic Acids Research* **49**, D344-D354 (2020).
- 1073 137. S. R. Eddy, Accelerated profile HMM searches. *PLoS Comput Biol* **7**, e1002195 (2011).
- 1074 138. G. Lum, X. J. Min, FunSeckB: the fungal secretome knowledgeBase. *Database (Oxford)*  
1075 **2011**, bar001 (2011).
- 1076 139. M. Urban *et al.*, PHI-base: the pathogen–host interactions database. *Nucleic Acids*  
1077 *Research* **48**, D613-D620 (2019).
- 1078 140. M. Fischer, J. Pleiss, The Lipase Engineering Database: a navigation and analysis tool  
1079 for protein families. *Nucleic Acids Res* **31**, 319-321 (2003).
- 1080 141. A. Marchler-Bauer *et al.*, CDD: NCBI's conserved domain database. *Nucleic Acids Res*  
1081 **43**, D222-226 (2015).
- 1082 142. J. Mistry *et al.*, Pfam: The protein families database in 2021. *Nucleic Acids Research*  
1083 **49**, D412-D419 (2021).
- 1084 143. J. Huerta-Cepas *et al.*, eggNOG 5.0: a hierarchical, functionally and phylogenetically  
1085 annotated orthology resource based on 5090 organisms and 2502 viruses. *Nucleic*  
1086 *Acids Research* **47**, D309-D314 (2018).
- 1087 144. J. Sperschneider, P. N. Dodds, K. B. Singh, J. M. Taylor, ApoplastP: prediction of  
1088 effectors and plant proteins in the apoplast using machine learning. *New Phytologist*  
1089 **217**, 1764-1778 (2018).
- 1090 145. J. Sperschneider *et al.*, LOCALIZER: subcellular localization prediction of both plant and  
1091 effector proteins in the plant cell. *Scientific Reports* **7**, 44598 (2017).
- 1092 146. J. J. Almagro Armenteros *et al.*, Detecting sequence signals in targeting peptides using  
1093 deep learning. *Life Sci Alliance* **2**, e201900429 (2019).

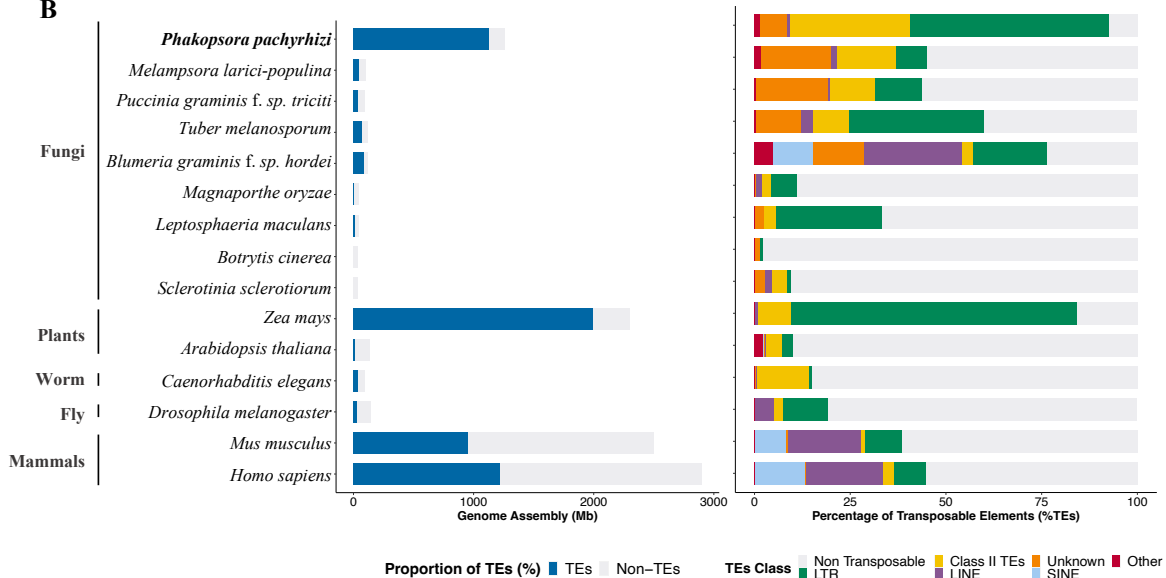
- 1094 147. P. Horton *et al.*, WoLF PSORT: protein localization predictor. *Nucleic acids research* **35**,  
1095 W585-W587 (2007).
- 1096 148. J. J. Almagro Armenteros, C. K. Sønderby, S. K. Sønderby, H. Nielsen, O. Winther,  
1097 DeepLoc: prediction of protein subcellular localization using deep learning.  
1098 *Bioinformatics* **33**, 3387-3395 (2017).
- 1099 149. D. G. O. Saunders *et al.*, Using hierarchical clustering of secreted protein families to  
1100 classify and rank candidate effectors of rust fungi. *PLOS ONE* **7**, e29847 (2012).
- 1101 150. D. Müllner, fastcluster: fast hierarchical, agglomerative clustering routines for R and  
1102 python. *Journal of Statistical Software* **53**, 1 - 18 (2013).
- 1103 151. A. L. Pendleton *et al.*, Duplications and losses in gene families of rust pathogens  
1104 highlight putative effectors. *Front Plant Sci.* **5**, 299 (2014).
- 1105 152. A. Persoons *et al.*, Genomic signatures of a major adaptive event in the pathogenic  
1106 fungus *Melampsora larici-populina*. *Genome Biol Evol.* **14**, evab279 (2022).
- 1107 153. S. Duplessis *et. al.*, Obligate biotrophy features unraveled by the genomic analysis of  
1108 rust fungi. *PNAS* **108**, 9166-9171 (2011).
- 1109 154. M. Toome *et al.*, Genome sequencing provides insight into the reproductive biology,  
1110 nutritional mode and ploidy of the fern pathogen *Mixia osmundae*. *New Phytologist*,  
1111 **202**, 554–564 (2013).
- 1112 155. M. H. Perlin *et al.*, Sex and parasites: genomic and transcriptomic analysis of  
1113 *Microbotryum lychnidis-dioicae*, the biotrophic and plant-castrating anther smut  
1114 fungus. *BMC Genomics* **16**, 461 (2015).
- 1115 156. J. Kämper *et al.*, Insights from the genome of the biotrophic fungal plant pathogen  
1116 *Ustilago maydis*. *Nature* **444**, 97-101 (2006).
- 1117 157. J. Schirawski *et al.*, Pathogenicity determinants in smut fungi revealed by genome  
1118 comparison. *Science* **330**, 1546-8 (2010).
- 1119 158. F. Martin *et al.*, The genome of *Laccaria bicolor* provides insights into mycorrhizal  
1120 symbiosis. *Nature* **452**, 88-92 (2008).
- 1121 159. Å Olson *et al.*, Insight into trade-off between wood decay and parasitism from the  
1122 genome of a fungal forest pathogen. *New Phytol.* **194**, 1001-1013 (2012).
- 1123 160. L. Frantzeskakis *et al.*, Signatures of host specialization and a recent transposable  
1124 element burst in the dynamic one-speed genome of the fungal barley powdery mildew  
1125 pathogen. *BMC Genomics* **19**, 381 (2018).
- 1126 161. R. A. Dean *et al.*, The genome sequence of the rice blast fungus *Magnaporthe grisea*.  
1127 *Nature* **434**, 980-6 (2005).  
1128  
1129

1130 FIGURES

A



B



1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

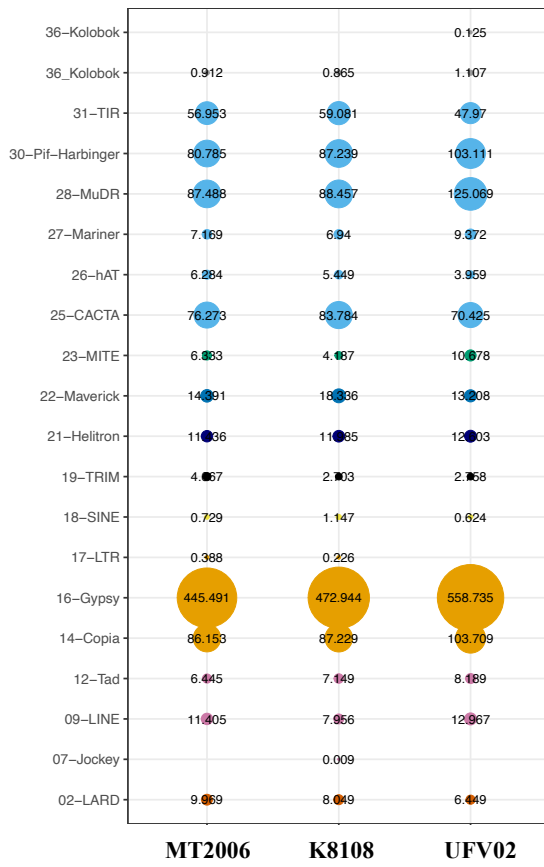
1142

1143

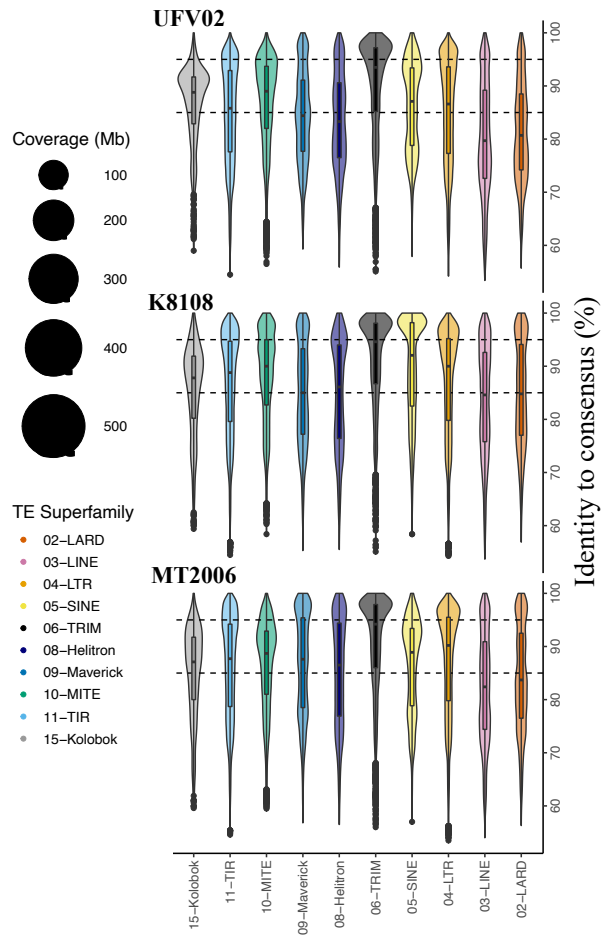
**Fig. 1. Impact of *P. pachyrhizi* incidence in a soybean field, comparative genome assembly size, and TE content.**

(A) Soybean field sprayed with fungicide (left) and unsprayed (right) in Brazil (top left). Soybean field being sprayed with fungicide (top right). Soybean leaf with a high level of *P. pachyrhizi* urediniospores, Tan reaction (bottom left). Electron micrograph of *P. pachyrhizi* infected leaf tissue, showing paraphyses and urediniospores highlighted in pseudo-color with orange, and leaf tissue in green, respectively (bottom right). (B) Transposable elements (TEs) content in different species of fungi (mostly plant pathogens), plants, and animals. The left histogram shows TEs proportion (%) per genome size, blue representing TEs content and grey non-TEs content; while the right histogram shows different classes of TEs in each genome.

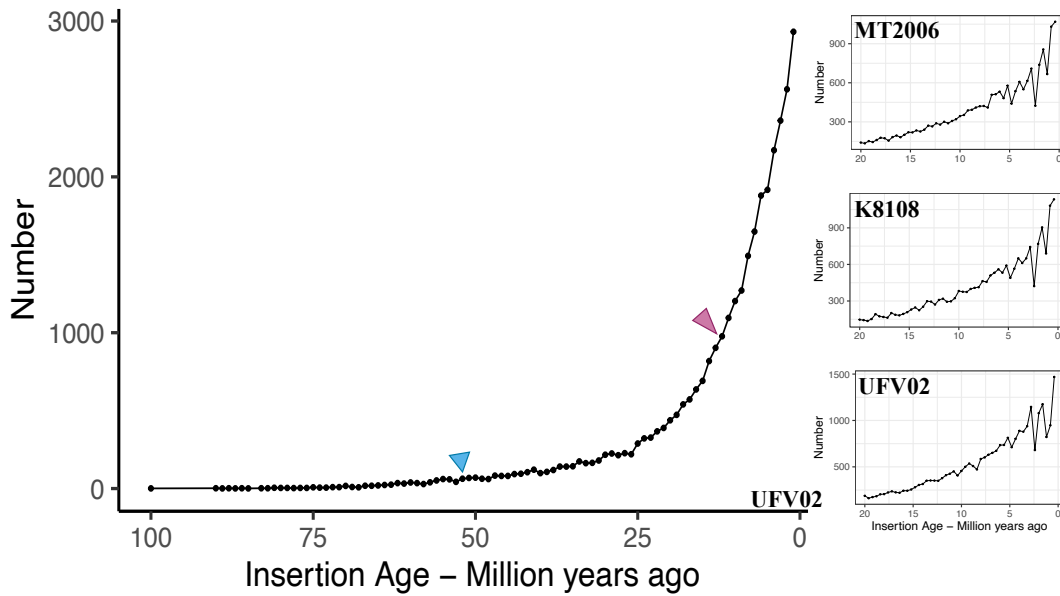
**A**



**B**

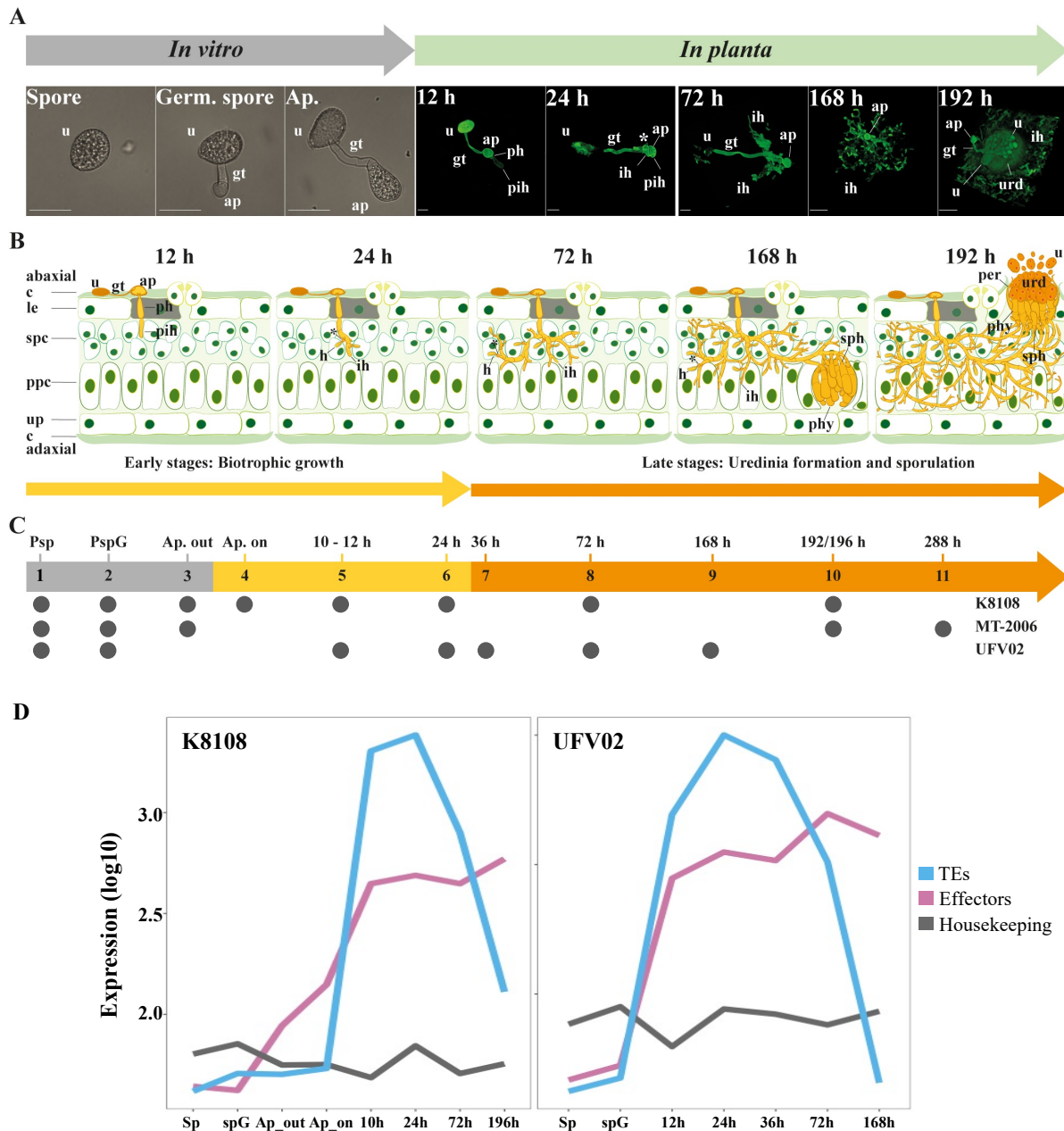


**C**



1144  
1145

1146 **Fig. 2. Transposable element superfamilies in *P. pachyrhizi* genomes.**  
1147 **(A)** Genome coverage of different TE superfamilies in three *P. pachyrhizi* genomes. **(B)** TE  
1148 superfamilies are categorized based on the consensus identity, (1) conserved TEs, copies with  
1149 more than 95% identity (2) intermediate TEs, copies with 85 to 95% identity and (3) divergent  
1150 TEs, copies with less than 85% identity. **(C)** The number of LTR retrotransposons in UFV02  
1151 based on the insertion age (Million years ago, Mya) with 1.0 million year intervals (left). The  
1152 legume speciation event around 53 Mya showed in blue triangle and ~13 Mya whole genome  
1153 duplication event in *Glycine* spp. marked with pink triangle (Schmutz *et al.*, 2010). In the right,  
1154 the three plot shows recent burst of TEs between 0-20 Mya in three genomes of MT2006,  
1155 K8108 and UFV02, respectively.  
1156



1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

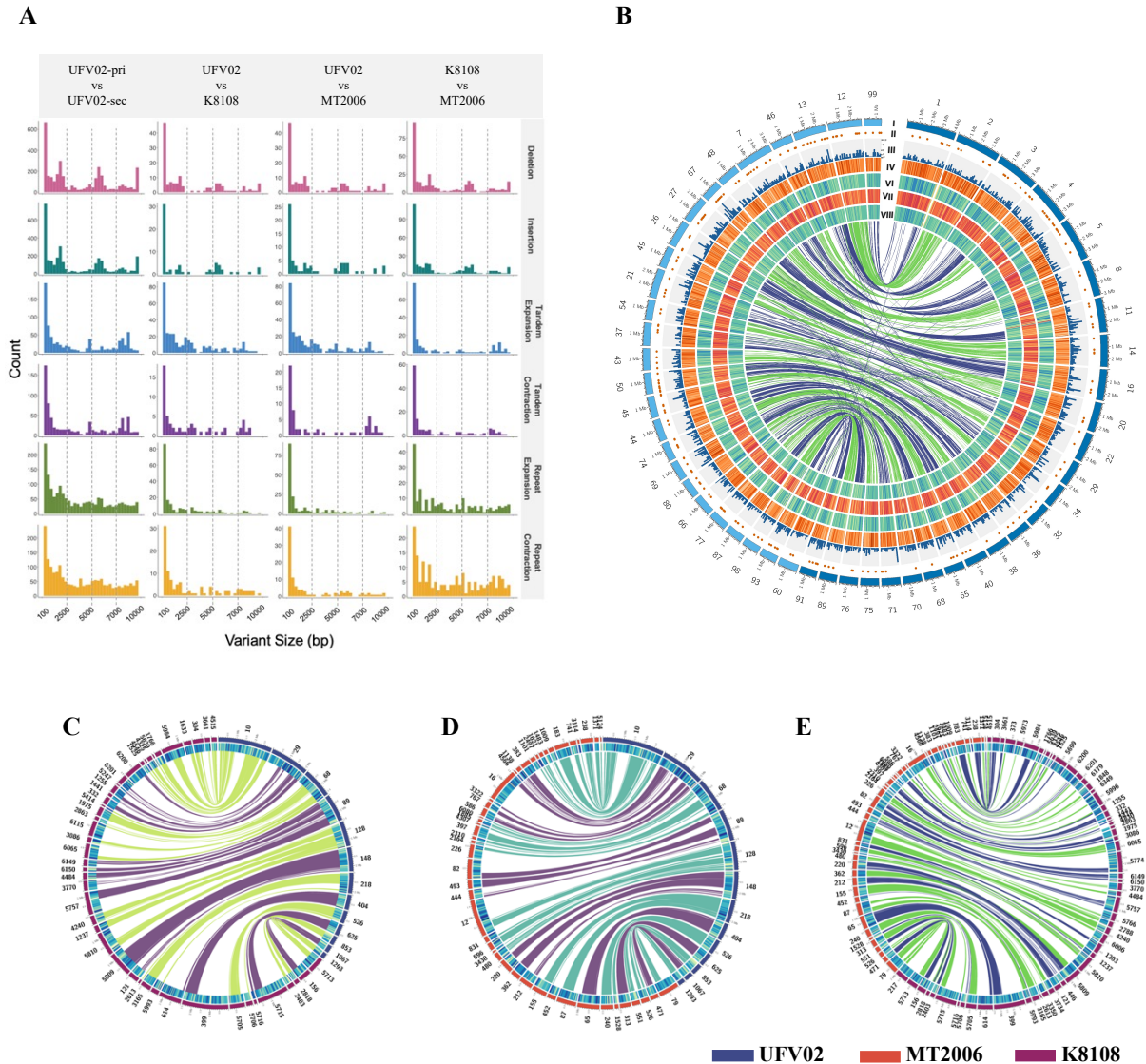
1168

1169

**Fig. 3. Infection cycle of *P. pachyrhizi* and gene expression on the critical infection stages.**

**(A)** Developmental phases of *P. pachyrhizi* infection *in vitro* and *in planta* on susceptible soybean plants. **(B)** Schematic of critical infection stages shown in the panel (A). **(C)** RNA sequencing on the critical time-points from three isolates. The time-points included in this study are assigned as small black circle for the three isolates. **(D)** Expression profiling of the effectors and a subset of TEs compared to the housekeeping genes during different stages of infection in K8108 and UFV02 isolates.

**Abbreviations:** urediospores (u), germ tube (gt), appressorium (ap), penetration hypha (ph), primary invasive hypha (pih), haustorial mother cell (\*), haustorium (h), invasive hyphae (ih), sporogenous hyphae (sph), paraphyses (phy), peridium (per), uredinium (urd), cuticle (c), lower epidermal cells (lec), spongy parenchyma cells (spc), palisade parenchyma cells (ppc), upper epidermal cells (ue).



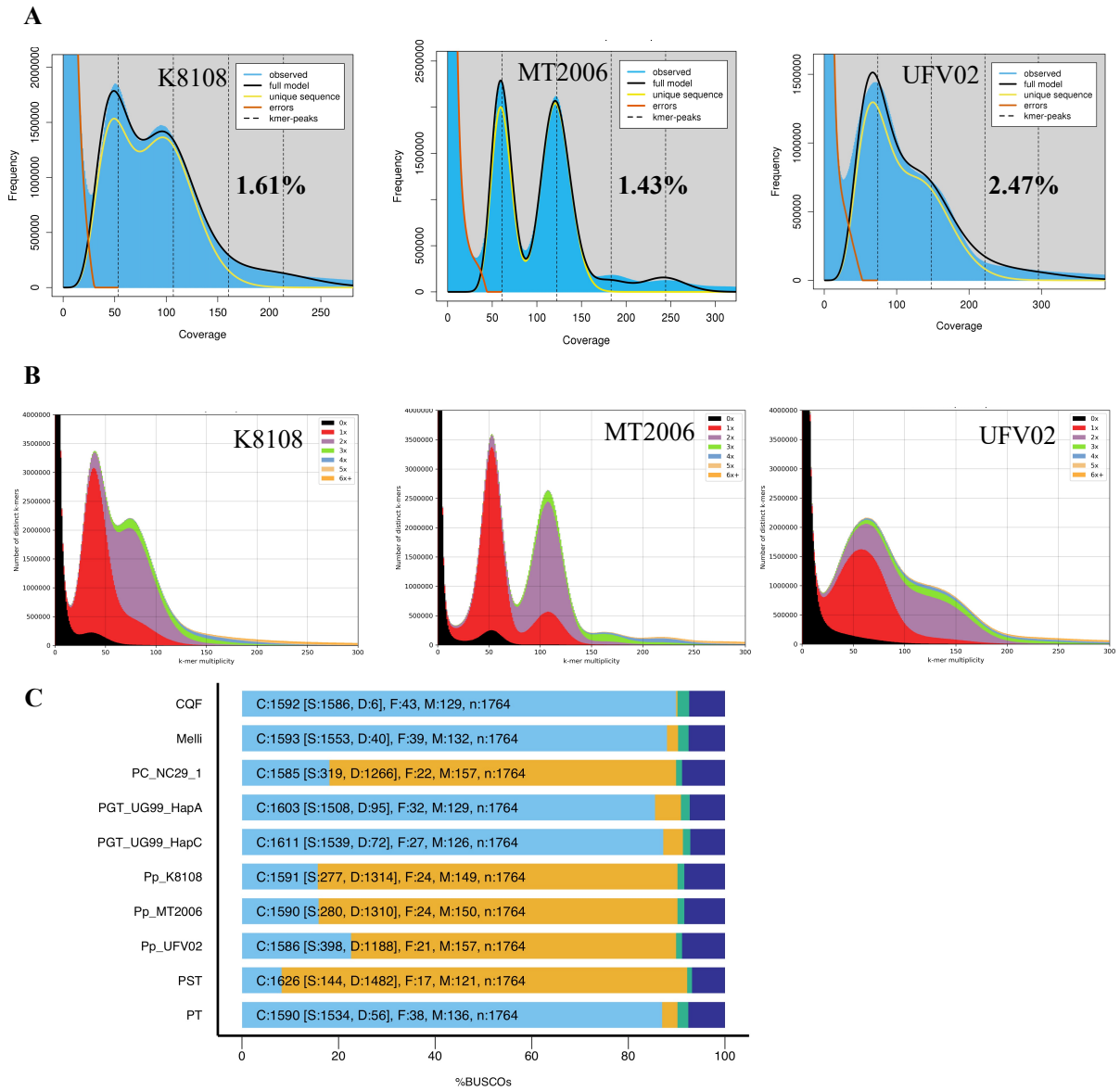
1170  
1171  
1172  
1173  
1174  
1175  
1176  
1177  
1178  
1179  
1180

**Fig. 4. Structural variation between *P. pachyrhizi* haplotypes is higher than variation between isolates.**

**(A)** Density plots with different structural variation between haplotypes and across isolates. **(B)** Circos plot representing inter-haplotype variation in PpUFV02 isolate. Layers from outside: I dark blue represent primary haplotigs and light blue secondary haplotigs; II secreted protein; III gene density (100 kb); IV TE density (50 kb); V SNP density K8108 isolate (25 kb); VI SNP density MT2006 isolate; VII SNP density UFV02 isolate (25 kb). **(C-E)** Circos plot showing inter-isolate variation. Layers from outside: I contigs from isolates represent in different colors; II TE density.



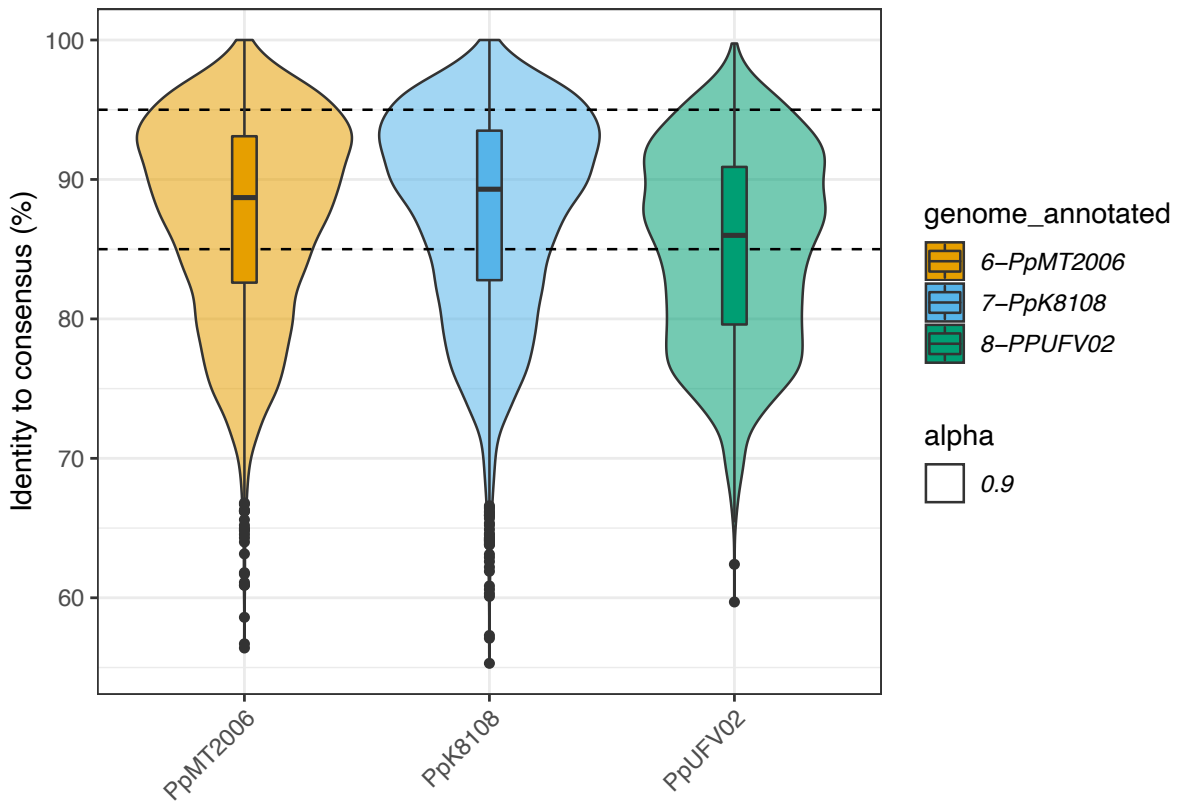
1181 SUPPLEMENTARY FIGURES



1182  
1183

1184 **Fig. S1. Comparison between the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**  
1185 **A.** K-mer frequency plots generated with WGS Illumina data of three isolates. The k-mer  
1186 frequency was estimated using Jellyfish and GenomeScope2. The x-axis shows k-mer  
1187 coverage and y-axis show the frequency. Two peaks in K-mer frequency profile shows a high  
1188 level of heterozygosity in *P. pachyrhizi*. The level of heterozygosity (shown in the bold letters)  
1189 varied between 1.43 to 2.47%. **B.** K-mer spectra plot comparing k-mer content of Illumina  
1190 read to k-mer content of the respective genomes, where different colors represent the  
1191 number of times k-mers from the reads found in the genome assembly. Black: indicates k-  
1192 mer content present in the raw reads but missing the genome assembly. Red: K-mers present  
1193 in the reads and once in the assembly. Purple: K-mers present in the reads and twice in the  
1194 genome assembly. Other colors indicate k-mers present in the genome more than twice. **C.**  
1195 BUSCO analysis of three *P. pachyrhizi* genomes and comparison with the genomes of other  
1196 published rust fungi. The basidiomycota database (n=1764) was used for the BUSCO analysis.  
1197 **Abbreviations:** *Cronartium quercuum* f. sp. *fusiforme* G11 (CQF), *Melampsora lini* CH5 (Melli),  
1198 *Puccinia coronata* f. sp. *avenae* 12NC29 (PC\_NC29\_1), *Puccinia graminis* f. sp. *tritici* UG99  
1199 haplotype A (PGT\_UG99\_HapA), *Puccinia graminis* f. sp. *tritici* UG99 haplotype C  
1200 (PGT\_UG99\_HapC), *P. pachyrhizi* K8108 (Pp\_K8108), *P. pachyrhizi* MT2006 (Pp\_MT2006), *P.*  
1201 *pachyrhizi* UFV02 (Pp\_UFV02), *Puccinia striiformis* f. sp. *tritici* PST-130 (PST), *Puccinia triticina*  
1202 Pt76 (PT).  
1203

1204



1205

1206

1207

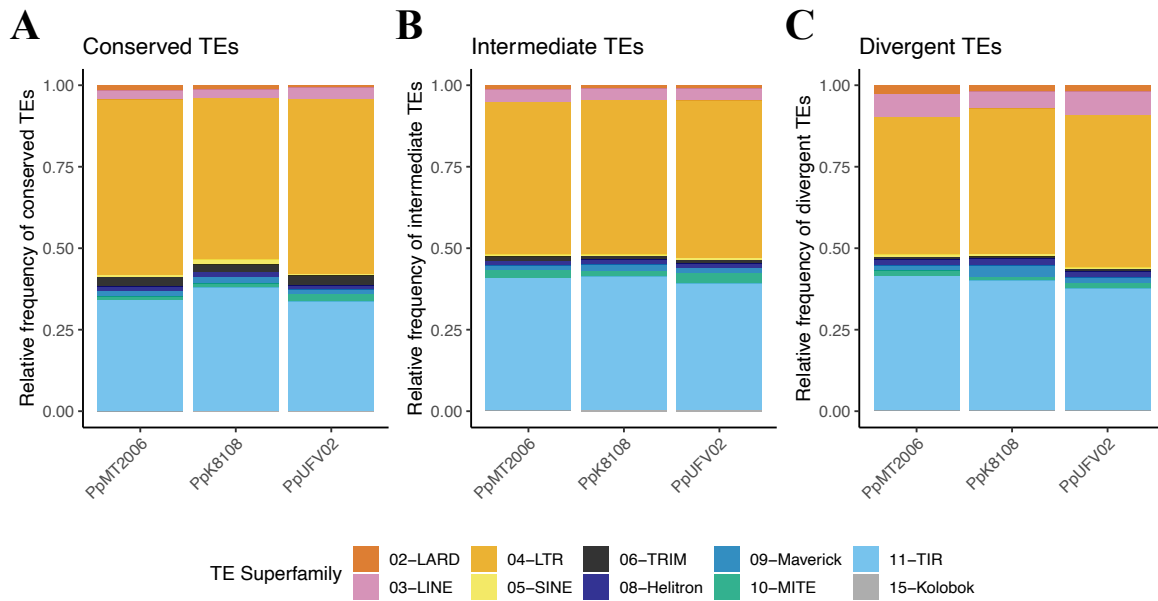
1208

1209

1210

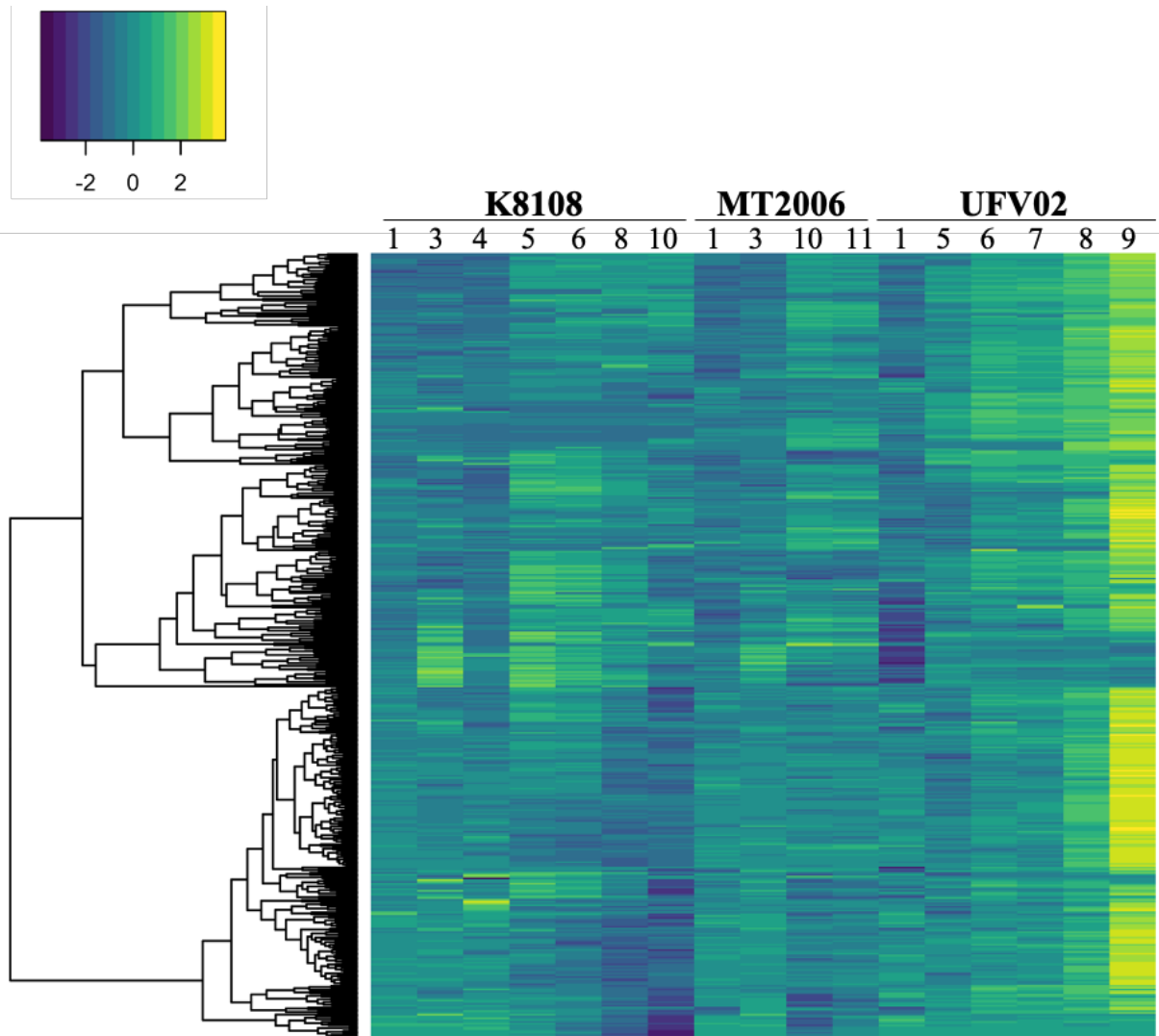
1211

**Fig. S2. TE consensus identity in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.** Based on the sequence identity, TEs were categorized as (1) conserved TEs (copies with more than 95% identity), (2) intermediate TEs (copies with 85 to 95% identity) and (3) divergent TEs (copies with less than 85% identity). The dotted line represents the cutoff for the sequence identity.



1212  
1213  
1214  
1215

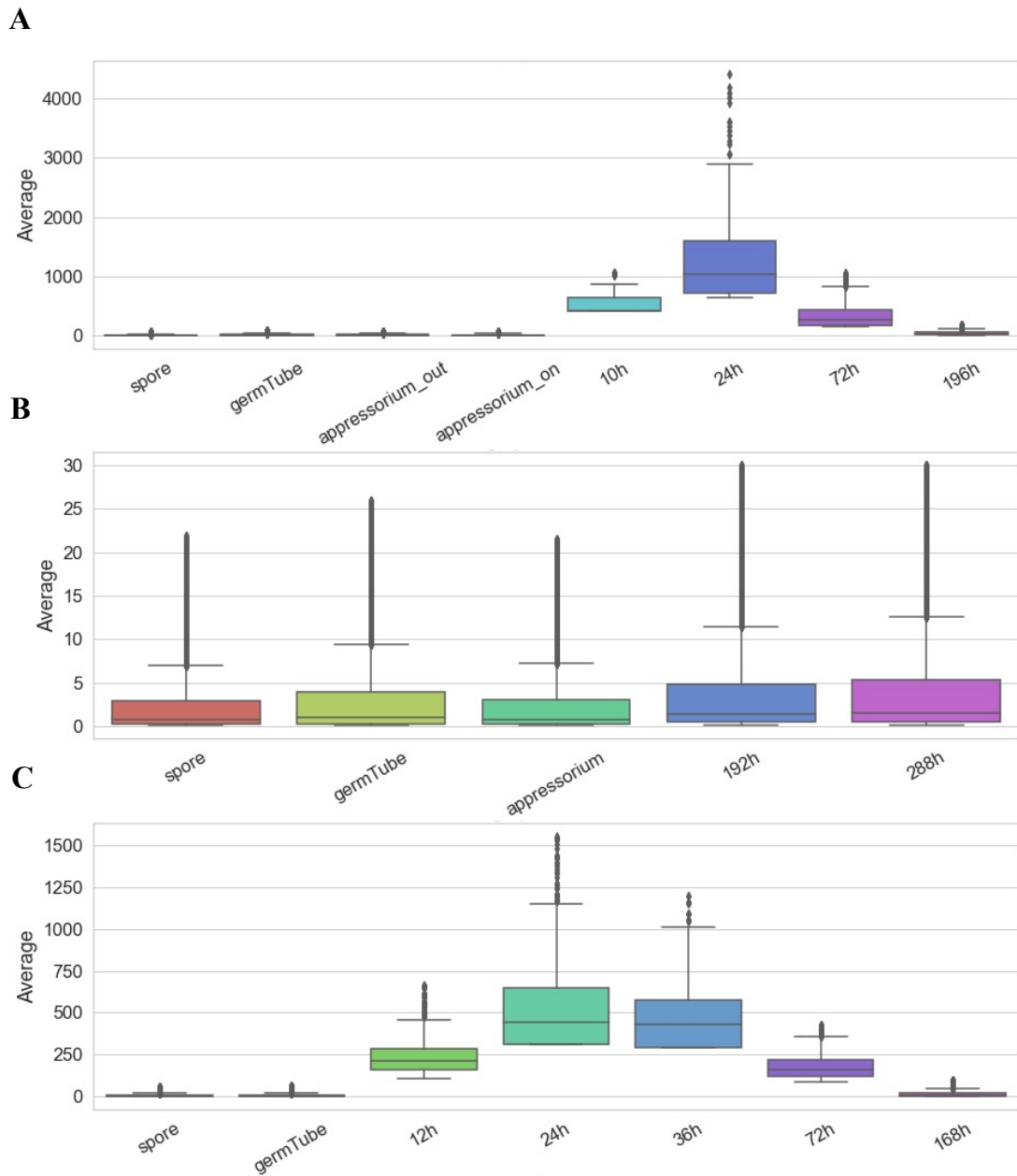
**Fig. S3. Relative frequency of TE superfamilies in categories such as, Conserved TEs (A), Intermediate TEs (B), and Divergent TEs (C) in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**



1216  
1217  
1218  
1219  
1220  
1221  
1222  
1223

**Fig. S4. Heatmap of differentially expressed secreted genes from the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**

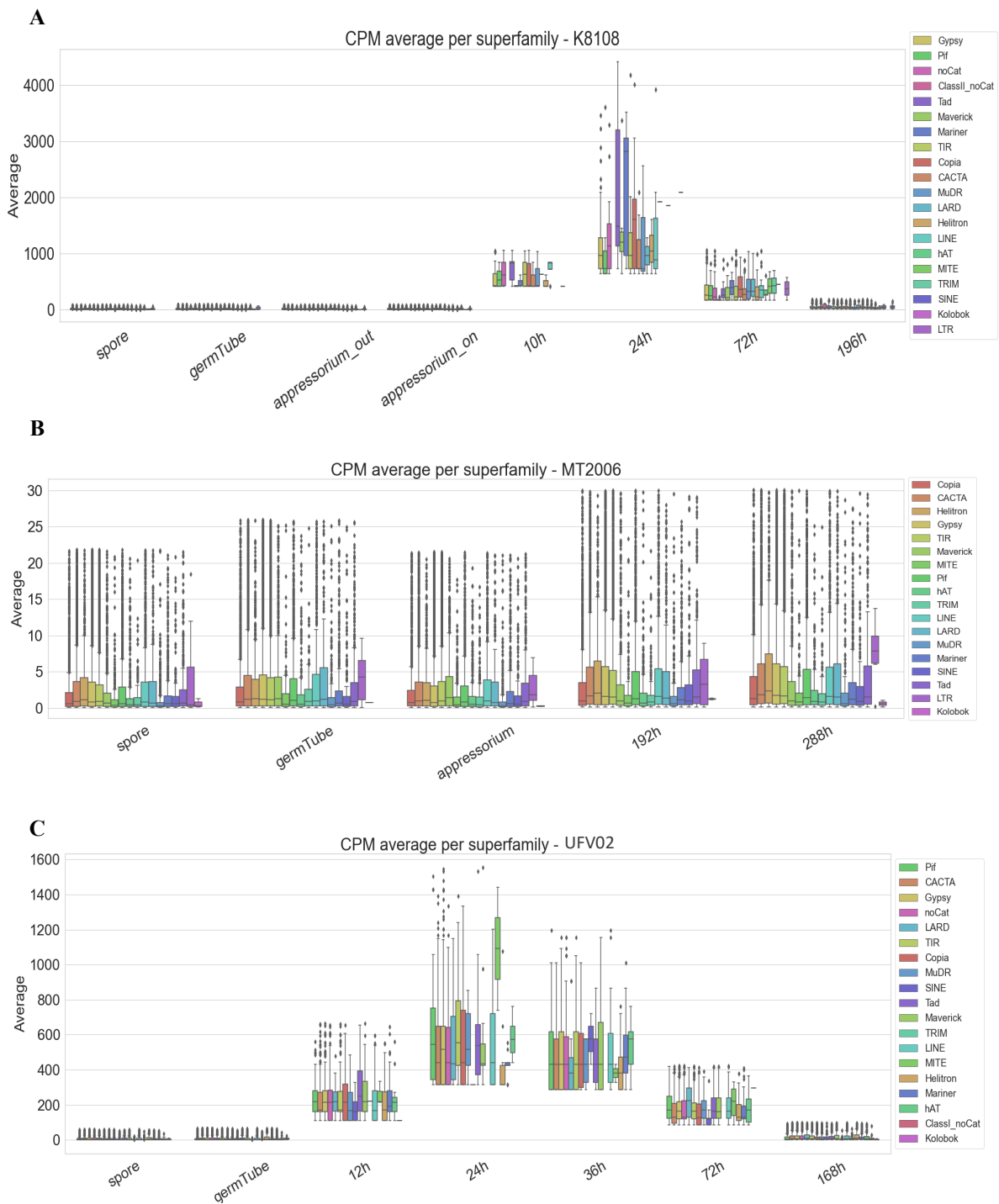
DEGs were hierarchical clustered by treatment, applying hclust method using R package (142). The germinated-spore (2) was used as calibrator. The other conditions are (1) Spore; (3) appressorium *in vitro*; (4) appressorium *in planta*; (5) 10-12 HPI; (6) 24 HPI; (7) 36 HPI; (8) 72 HPI; (9) 168 HPI; (10) 192-196 HPI; (11) 288 HPI.



1224  
1225  
1226  
1227

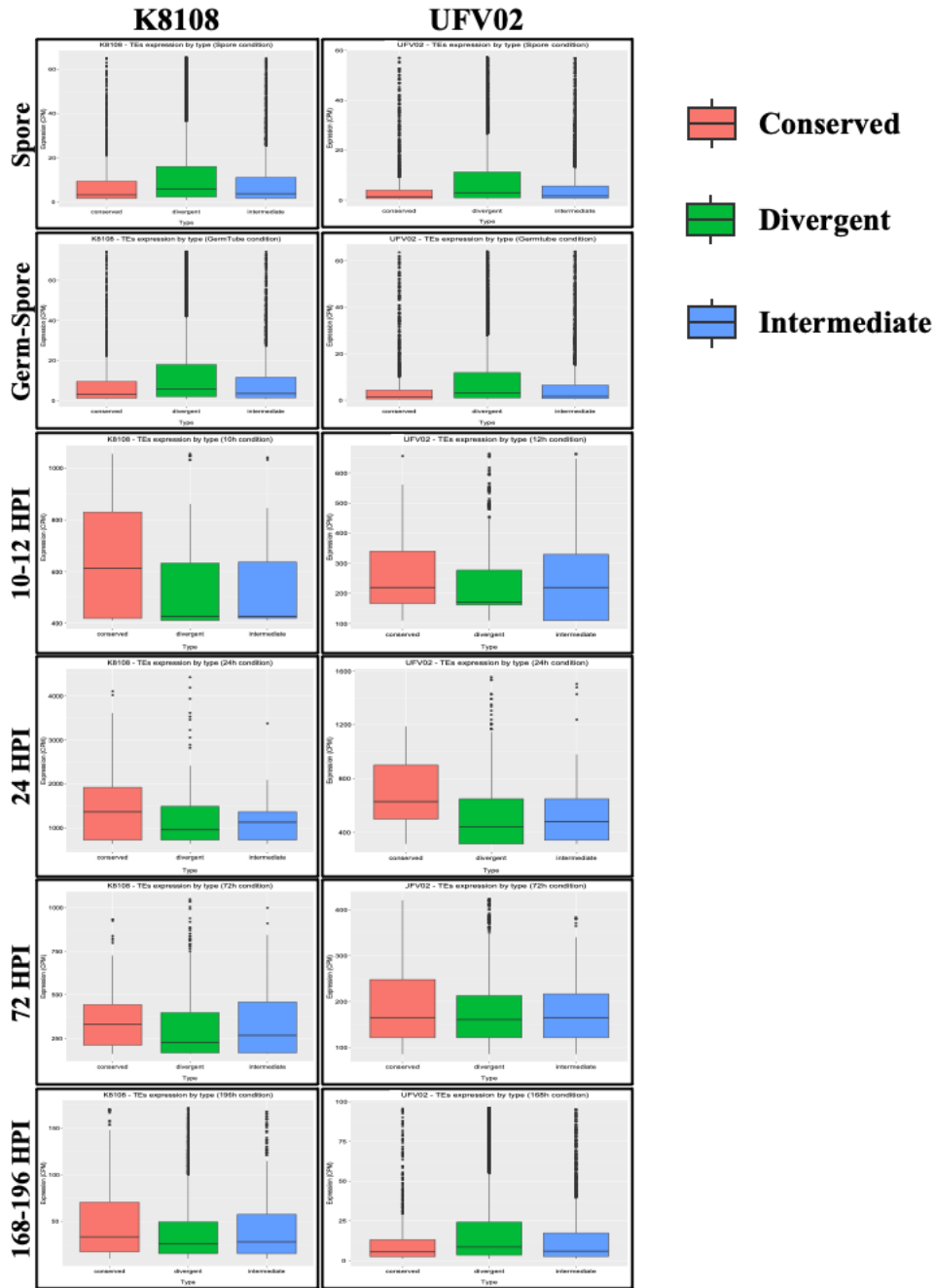
**Fig. S5. Expression profile of TEs on different condition in the *P. pachyrhizi* transcriptomes. Average of CPM (copies per million) of TEs. A, K8108; B, MT2006; C, UFV02.**

1228



1229  
1230  
1231  
1232  
1233

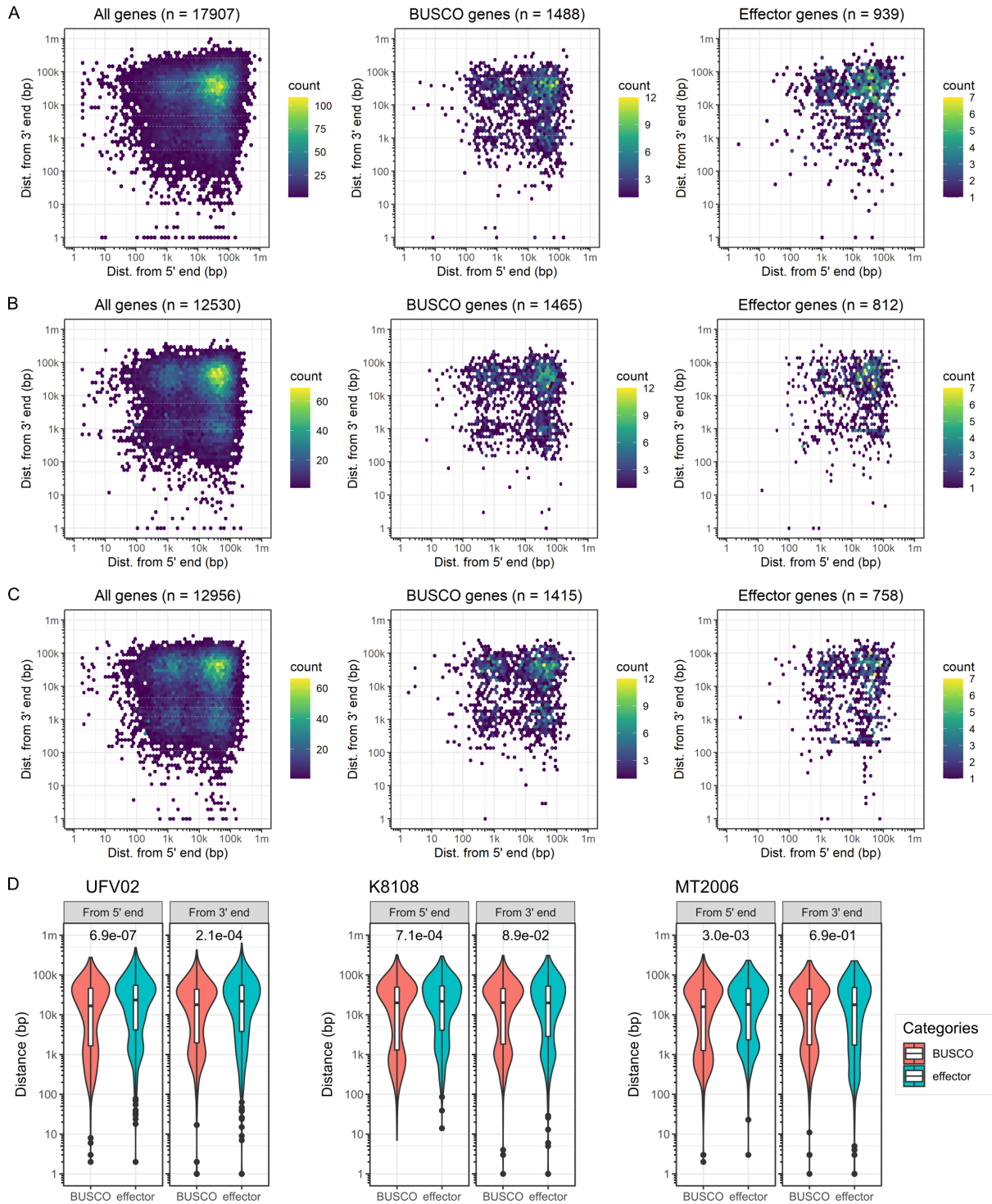
**Fig. S6. Expression profile of TEs in superfamilies different conditions (mentioned in Fig. 2) in the *P. pachyrhizi* transcriptomes. Average of CPM (copies per million) of TEs. A, K8108; B, MT2006; C, UFV02.**



1234  
1235  
1236  
1237  
1238

**Fig. S7. Expression profile of TEs based on conserved, divergent, and intermediate categories in the *P. pachyrhizi* genomes K8108 and UFV02. Average of CPM (copies per million) of TEs.**



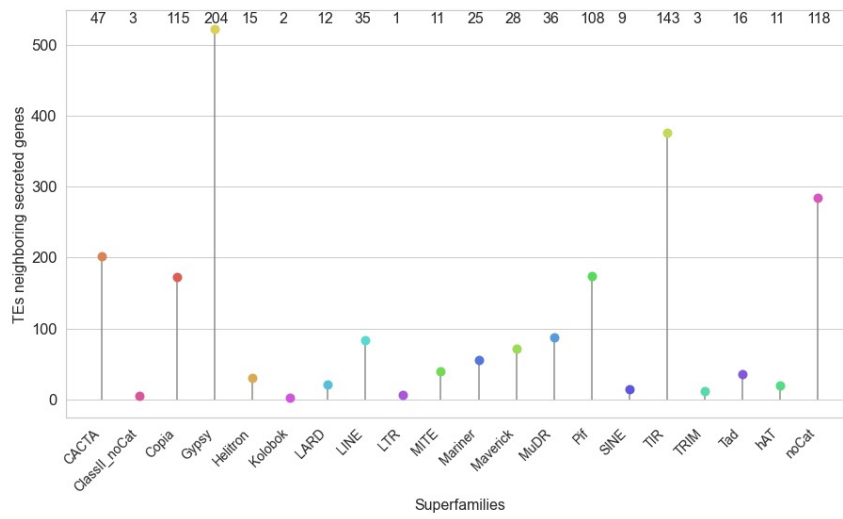


1239  
1240  
1241

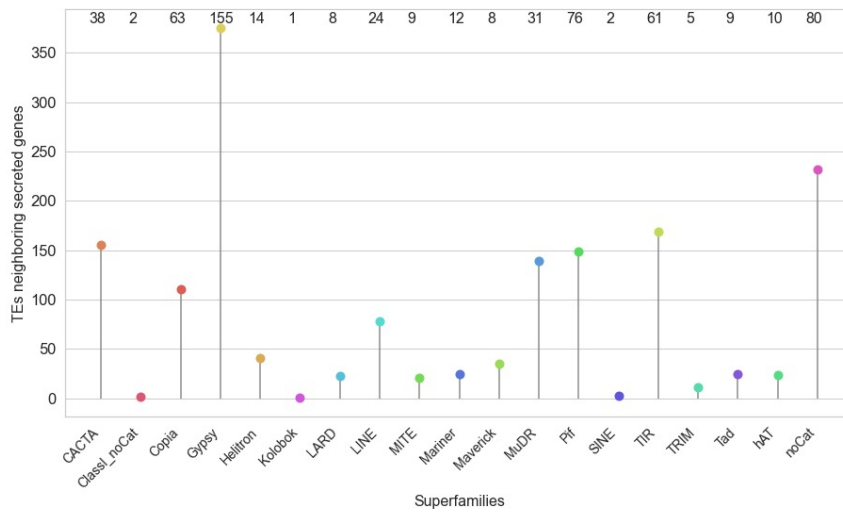
1242 **Fig. S8. Distribution of effector genes in comparison with gene catalogues and BUSCO genes**  
1243 **in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**

1244 **A-C.** Hexbin plots for 5'(x-axis) and 3'(y-axis) intergenic distances. The left-most column  
1245 represents profiles for all genes, the middle column for BUSCO genes, and the right-most  
1246 column for effector genes. The number of genes included in the analysis (genes with both  
1247 flanks within the same contig) is indicated in the parenthesis. **A.** UFV02; **B.** K8108; **C.** MT2006.  
1248 **D.** Violin plots for 5' and 3' intergenic distances of BUSCO and effector genes. *P* values from  
1249 Wilcox test are indicated in the plots. The basidiomycota\_odb10 dataset was used for the  
1250 BUSCO analysis.

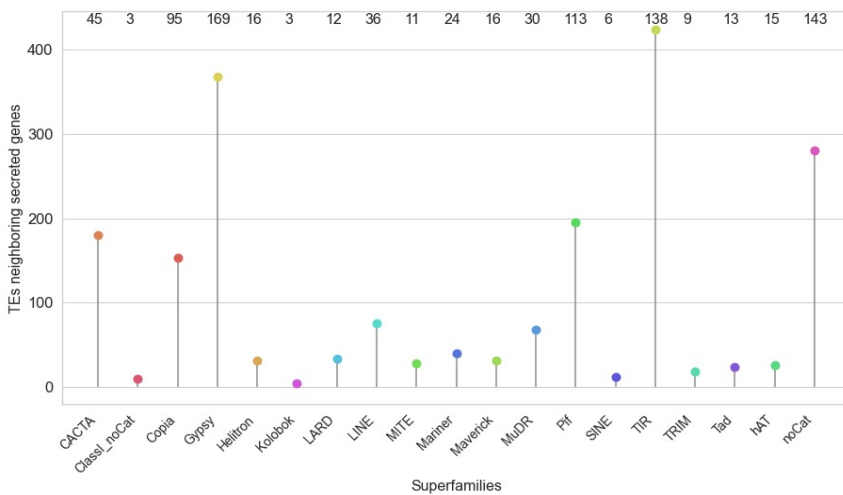
**A**



**B**



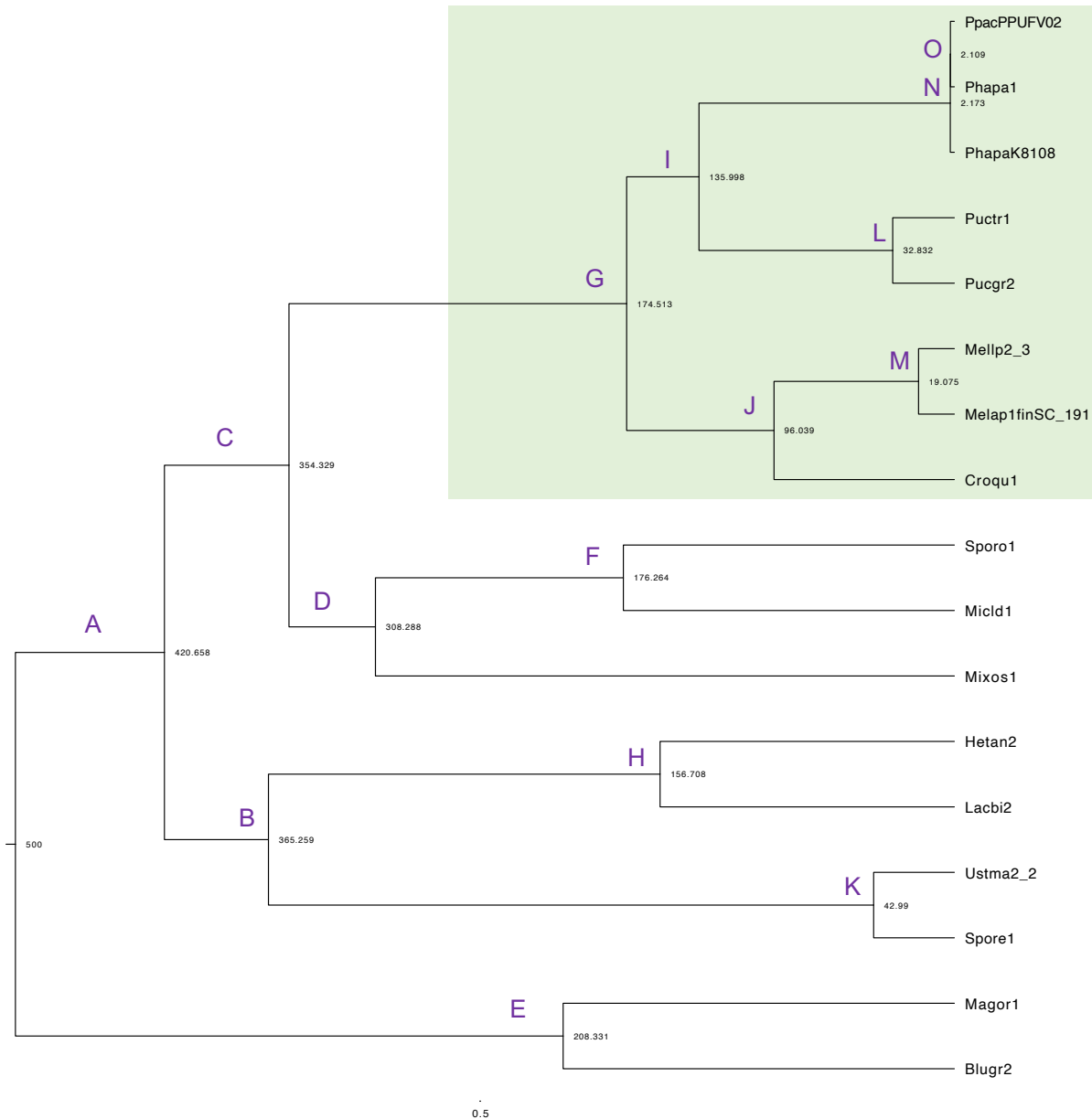
**C**



1251  
1252  
1253  
1254

**Fig. S9. Association of secreted genes to the neighboring TE families in the *P. pachyrhizi* genomes. A, K8108; B, UFV02; C, MT2006.**

The number of secreted genes correspond to the TE families shown at the top of the plot.



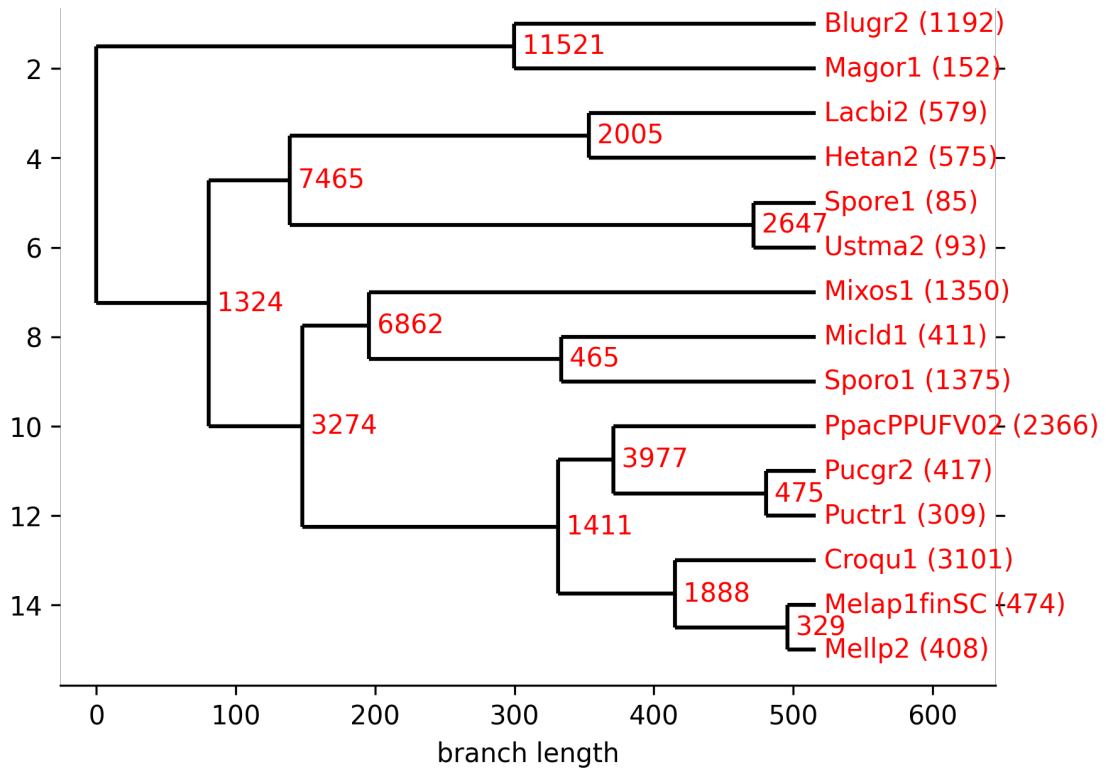
1255  
1256  
1257  
1258  
1259  
1260  
1261  
1262  
1263  
1264  
1265  
1266  
1267  
1268  
1269  
1270

**Fig. S10. Phylogenetic relationships and estimated divergence time of selected fungal species.**

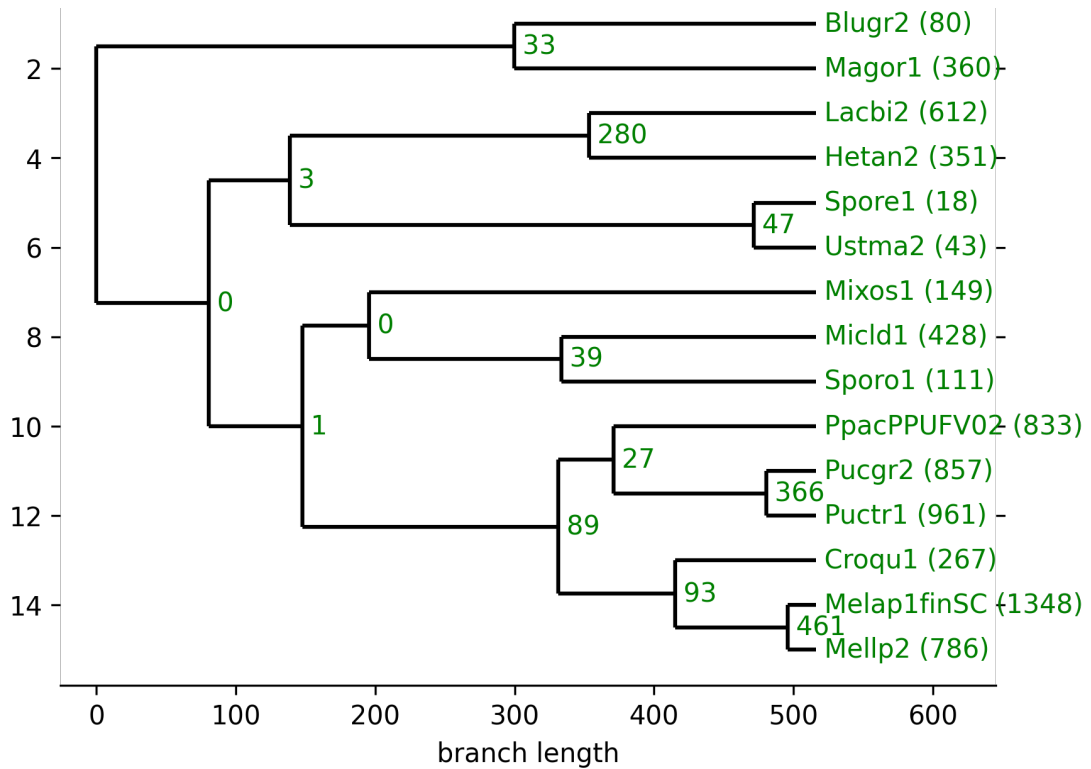
The phylogenetic tree was generated after alignment of 408 conserved orthologous markers identified from at least 13 genomes using PHYling (table S24a). The sequences were aligned and concatenated into a super-alignment with 408 partitions. Phylogenetic tree was built with RAxML-NG (v0.9.0) using a partitioned analysis and 200 bootstraps replicates.

**Abbreviations:** *P. pachyrhizi* MG2006 v1.0 (Phapa1), *P. pachyrhizi* K8108 v2.0 (PhapaK8108), *P. pachyrhizi* UFV02 v2.1 (PpacPPUFV02), *Cronartium quercuum* f. sp. *fusiforme* G11 (Croqu1), *Melampsora laricis-populina* v2.0 (Mellp2\_3), *M. allii-populina* 12AY07 v1.0 (Melap1finSC\_191), *P. graminis* f. sp. *tritici* v2.0 (Pucgr2), *P. triticina* 1-1 BBBD Race 1 (Puctr1), *Sporobolomyces roseus* v1 (Sporo1), *Mixia osmundae* (Mixos1), *Microbotryum lychnidisdioicae* p1A1 Lamole (Micld1), *Ustilago maydis* 521 v2.0 (Ustma2\_2), *Sporisorium reilianum* SRZ2 (Spore1), *Laccaria bicolor* v2 (Lacbi2), *Heterobasidion annosum* TC 32-1 (Hetan2), *Blumeria graminis* f. sp. *hordei* DH14 (Blugr2), *Magnaporthe oryzae* 70-15 v3.0 (Magor1).

**A**

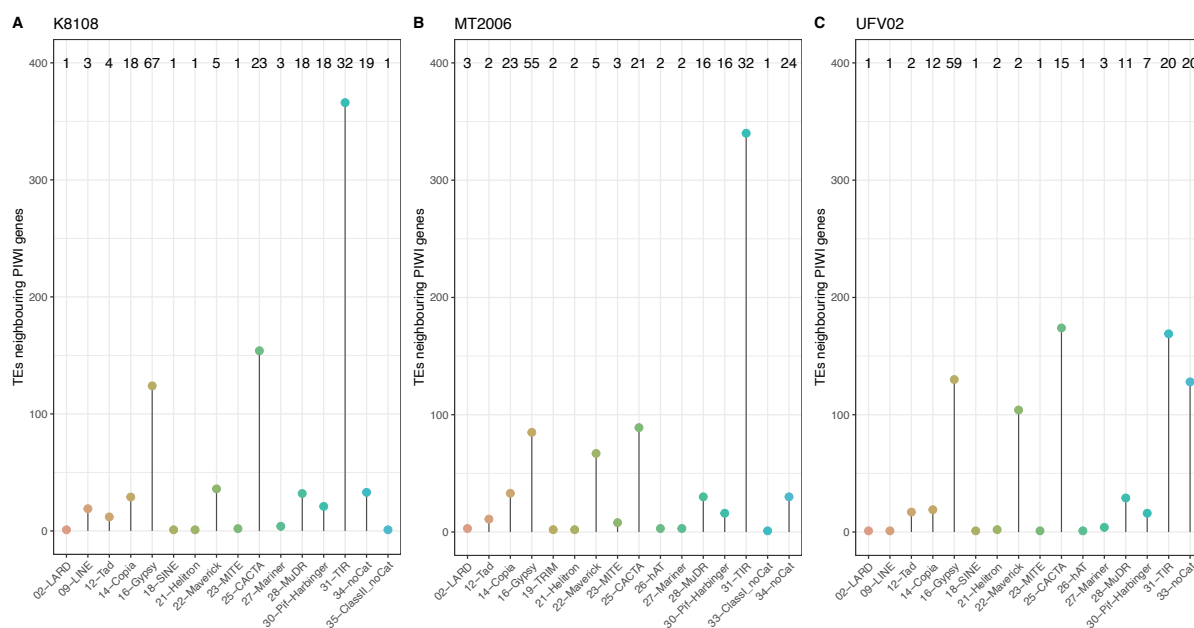


**B**



1271  
1272

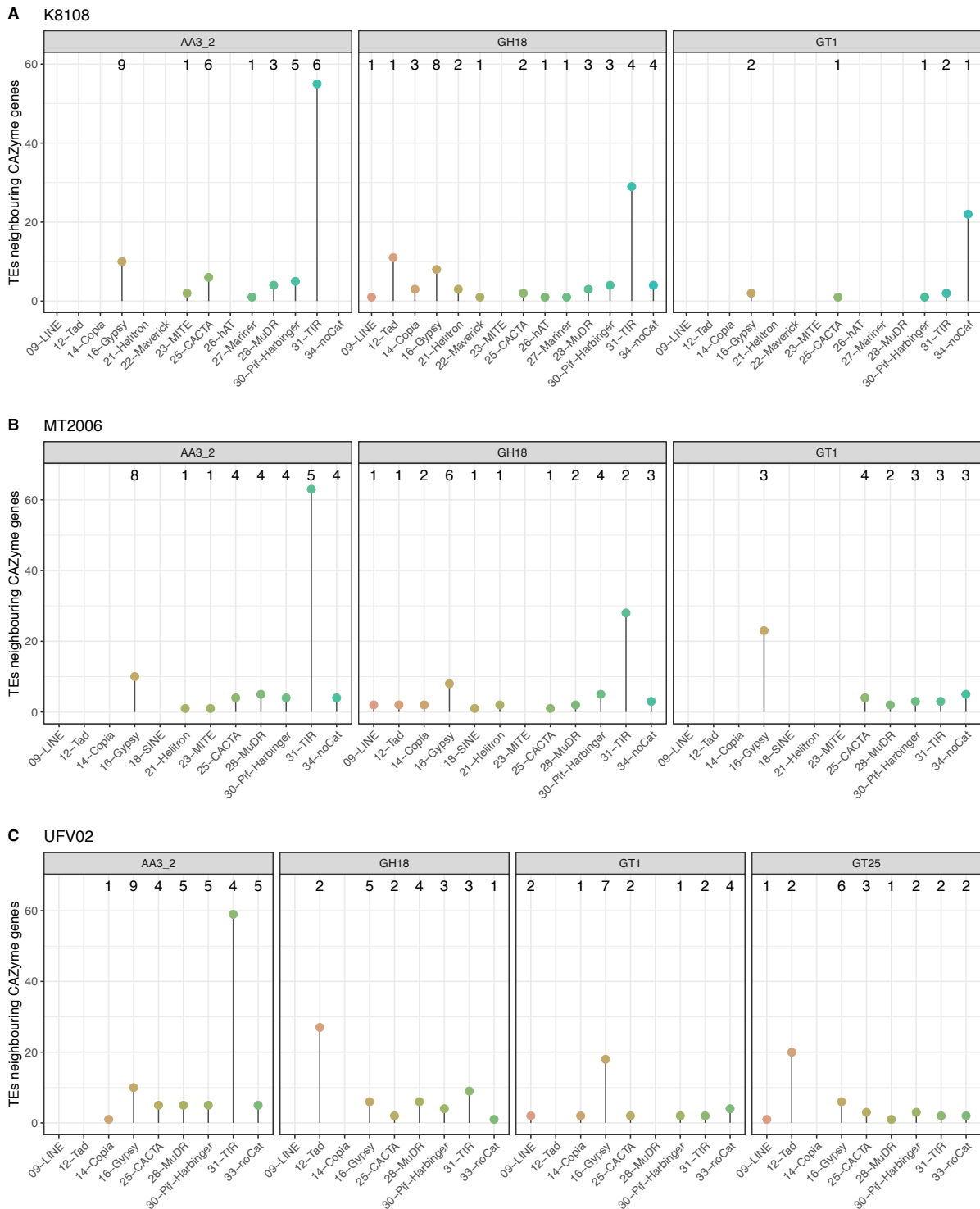
1273 **Fig. S11. Gene families contraction (A) and expansion (B) in 15 different fungal pathogens.**  
1274 The branch length represents differentiation time. Number of expanded and contracted gene  
1275 families are shown after the species name. The numbers on the nodes correspond to the  
1276 ancestral protein families.  
1277 **Abbreviations:** *P. pachyrhizi* UFV02 v2.1 (PpacPPUFV02), *Cronartium quercuum* f. sp.  
1278 *fusiforme* G11 (Croqu1), *Melampsora laricis-populina* v2.0 (Mellp2\_3), *M. allii-populina*  
1279 12AY07 v1.0 (Melap1finSC\_191), *P. graminis* f. sp. *tritici* v2.0 (Pucgr2), *P. triticina* 1-1 BBBD  
1280 Race 1 (Puctr1), *Sporobolomyces roseus* v1 (Sporo1), *Mixia osmundae* (Mixos1),  
1281 *Microbotryum lychnidis-dioicae* p1A1 Lamole (Micld1), *Ustilago maydis* 521 v2.0 (Ustma2\_2),  
1282 *Sporisorium reilianum* SRZ2 (Spore1), *Laccaria bicolor* v2 (Lacbi2), *Heterobasidion annosum*  
1283 TC 32-1 (Hetan2), *Blumeria graminis* f. sp. *hordei* DH14 (Blugr2), *Magnaporthe oryzae* 70-15  
1284 v3.0 (Magor1)  
1285



1286  
1287  
1288  
1289  
1290

**Fig. S12. Association of Piwi genes to the neighboring TE families in the *P. pachyrhizi* genomes. A, K8108; B, MT2006; C, UFV02.**

The number of Piwi genes correspond to the TE families shown at the top of the plot.

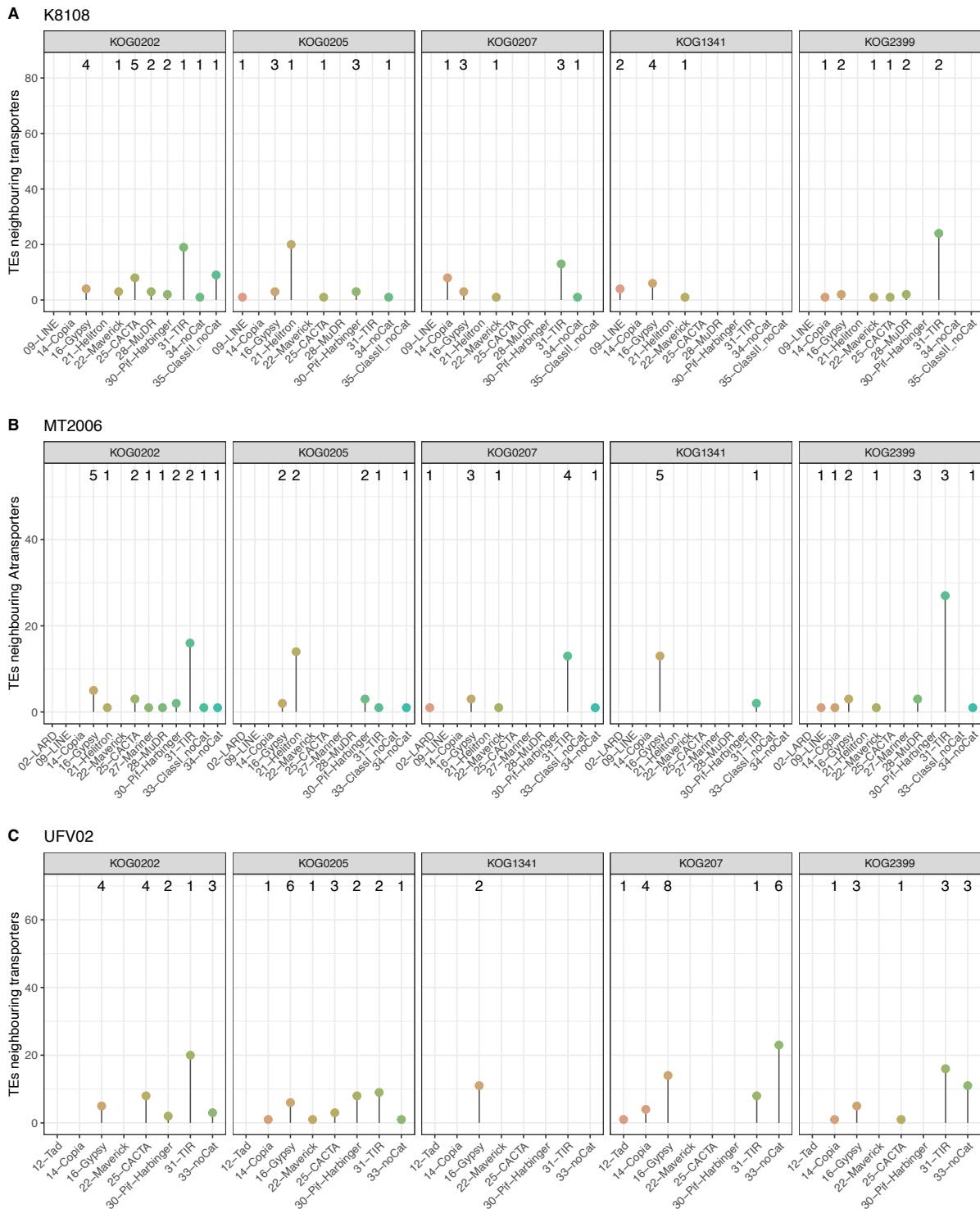


1291  
 1292 **Fig. S13. Association of CAZyme related genes to the neighboring TE families in the *P.***  
 1293 ***pachyrhizi* genomes. A, K8108; B, MT2006; C, UFV02.**

1294 The number of CAZyme correspond to the TE families shown at the top of the plot.

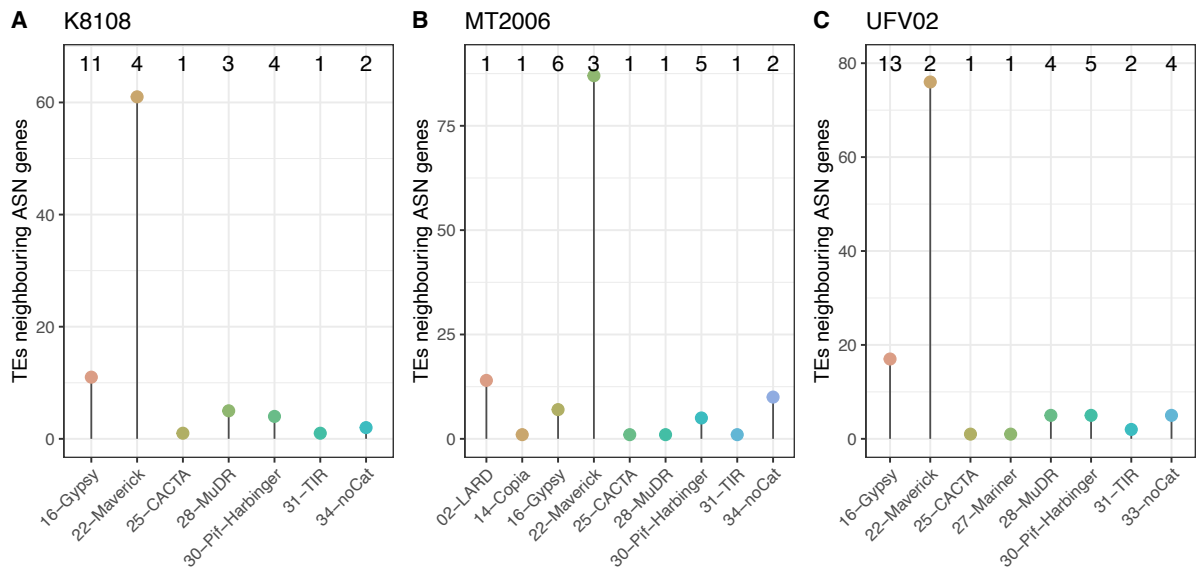
1295





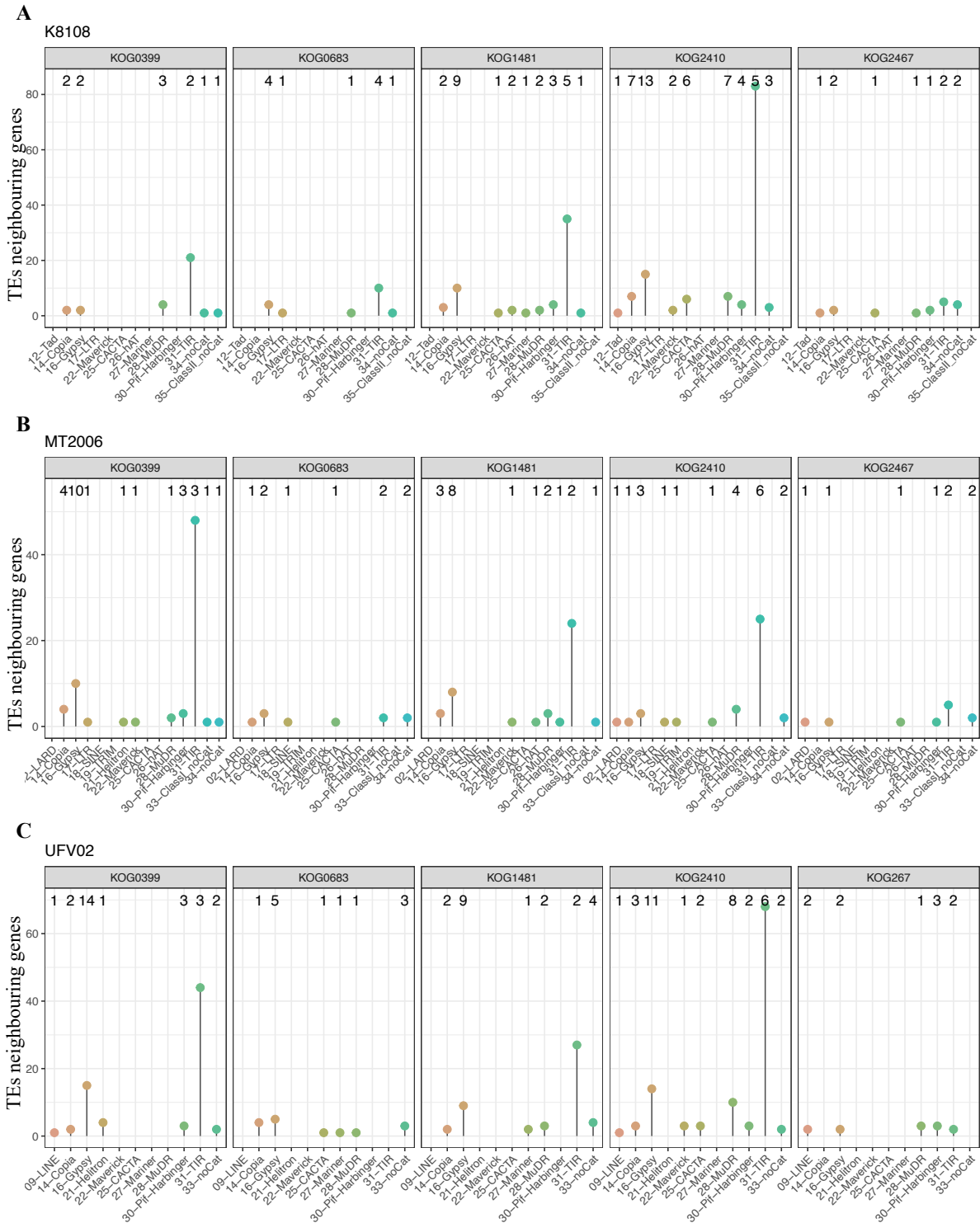
1296 **Fig. S14. Association of transporter related genes to the neighboring TE families in the *P.***  
 1297 ***pachyrhizi* genomes. A, K8108; B, MT2006; C, UFV02.**

1298 The number of transporter related genes correspond to the TE families shown at the top of  
 1299 the plot.  
 1300



1301  
1302  
1303  
1304  
1305

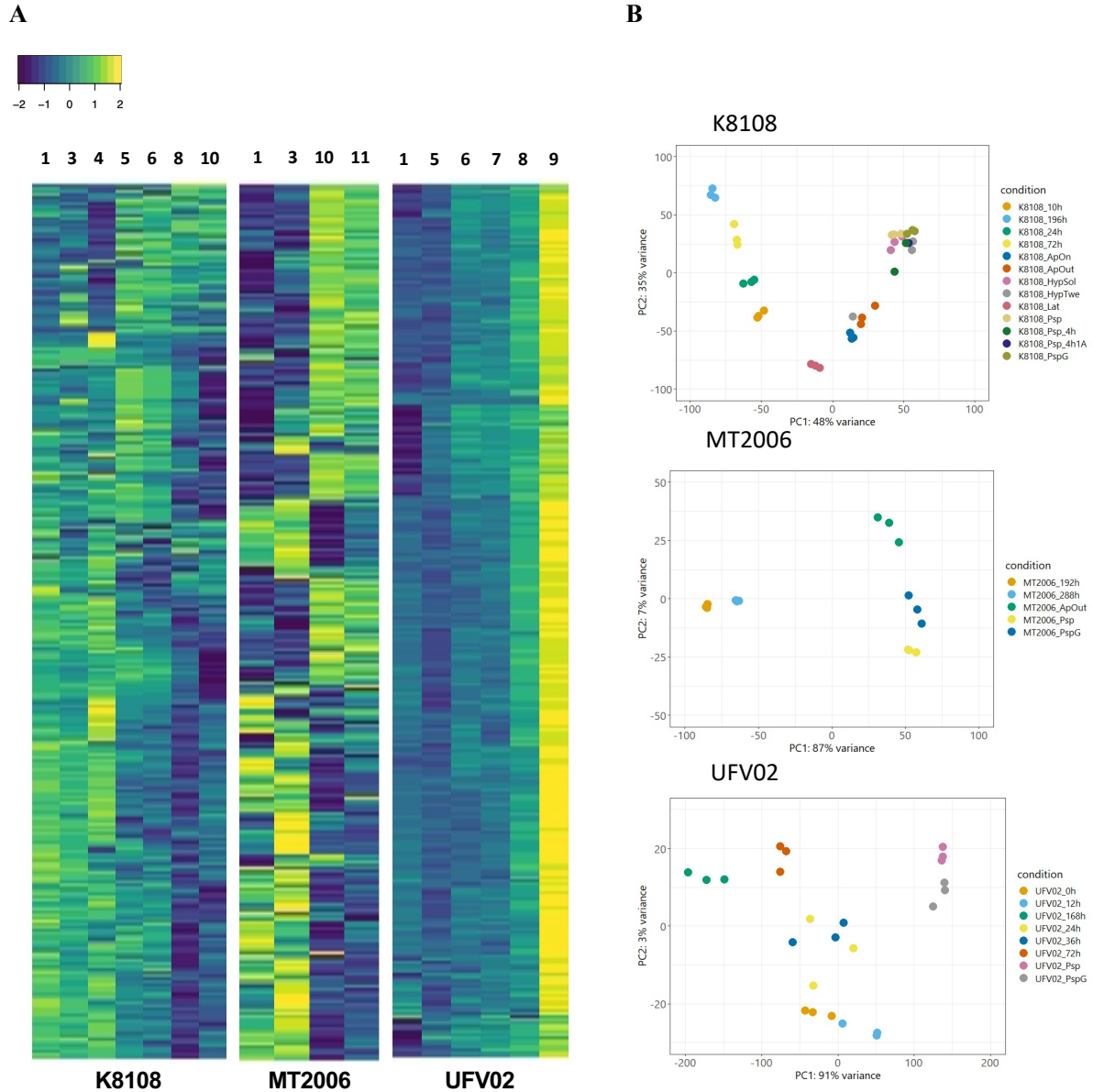
**Fig. S15. Association of Asparagine synthase (KOG0573) metabolism genes to the neighboring TE families in the *P. pachyrhizi* genomes. A, K8108; B, MT2006; C, UFV02.**  
The number of asparagine synthase genes correspond to the TE families shown at the top of the plot.



1306  
1307 **Fig. S16. Association of amino acid metabolism genes to the neighboring TE families in the**  
1308 ***P. pachyrhizi* genomes. A, K8108; B, MT2006; C, UFV02.**

1309 The number of amino acid metabolism genes correspond to the TE families shown at the top  
1310 of the plot.

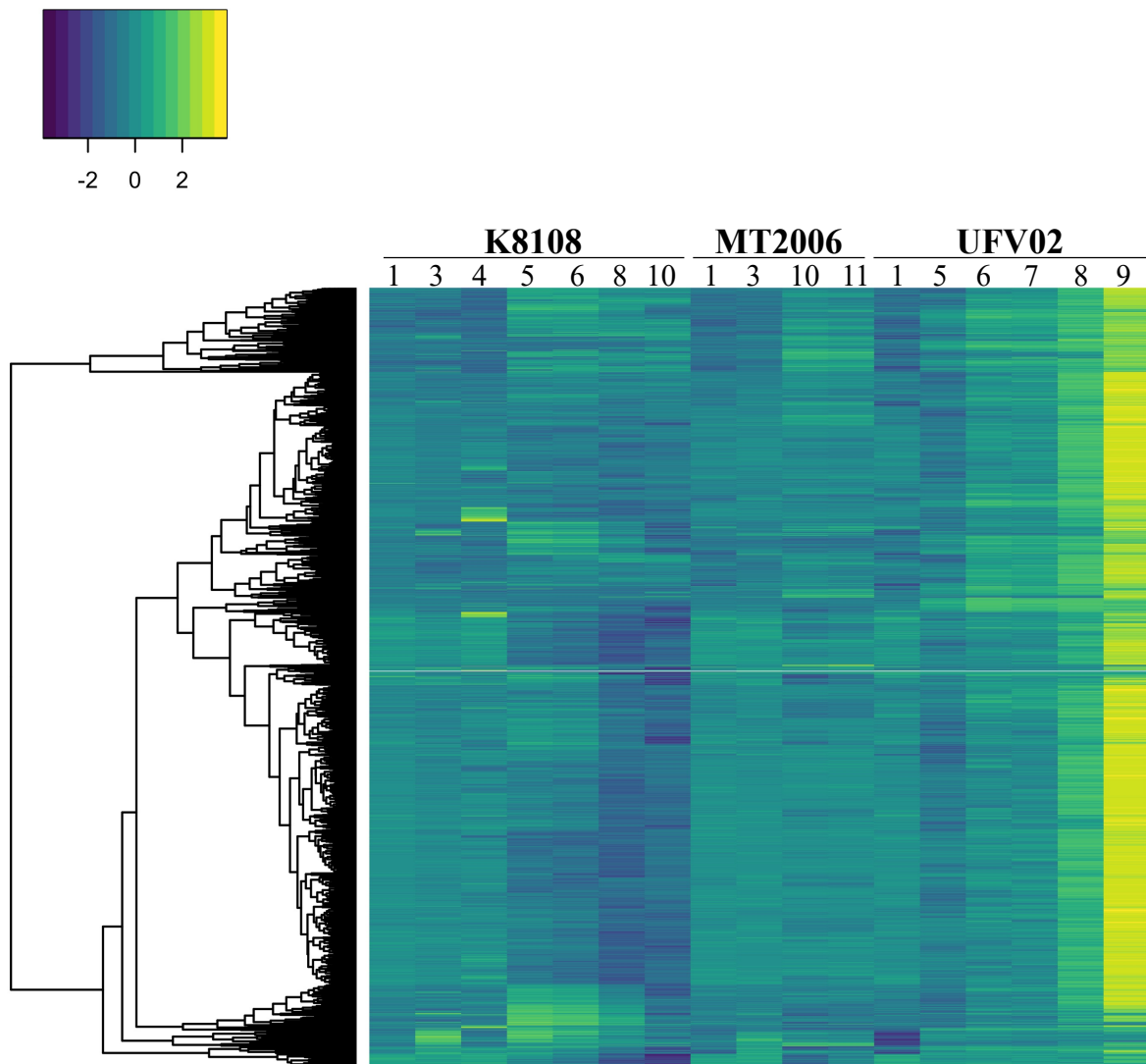
1311



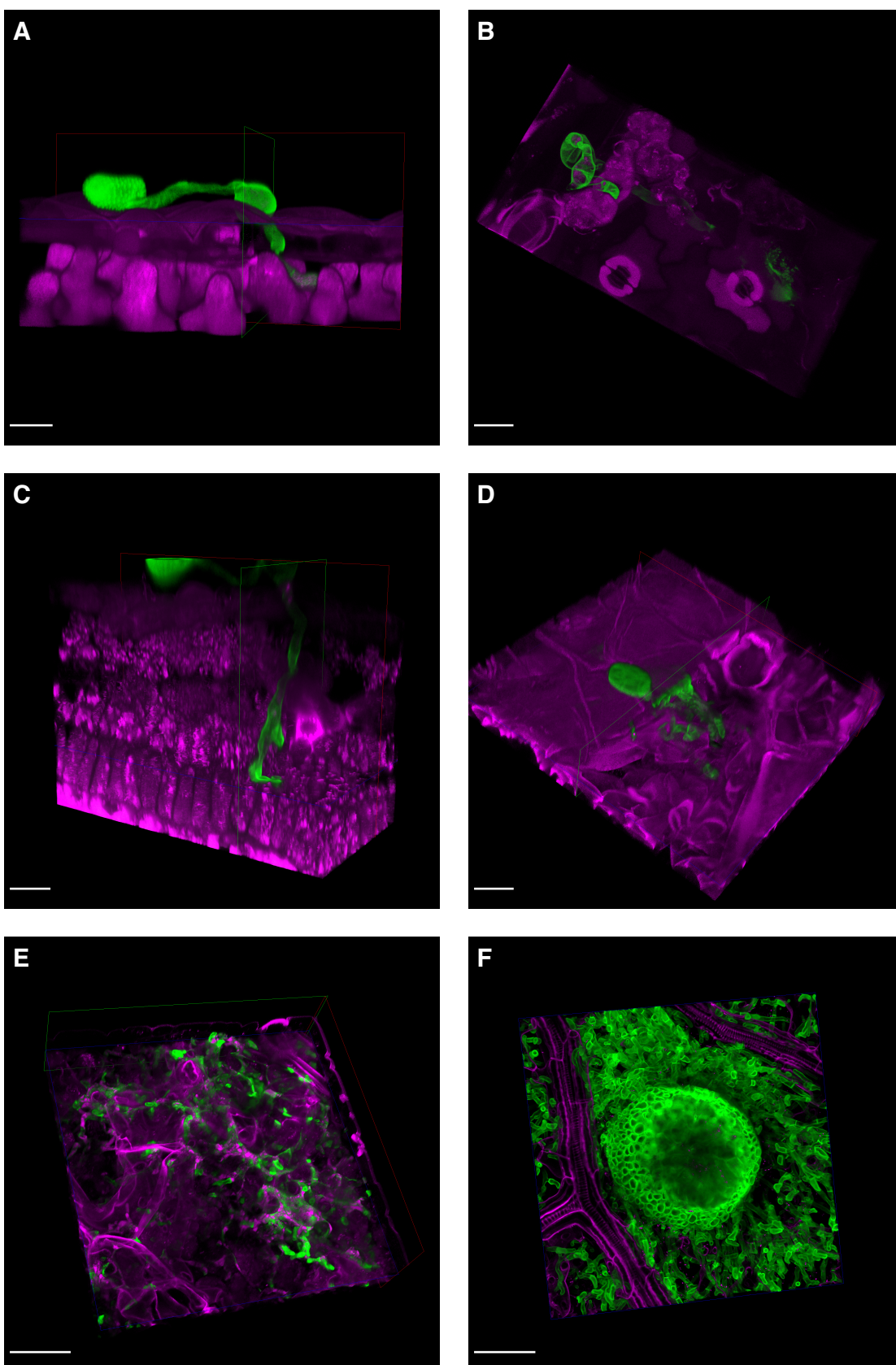
1312  
1313  
1314  
1315  
1316  
1317  
1318  
1319  
1320

**Fig. S17. Heatmap of differentially expressed genes (DEGs) in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**

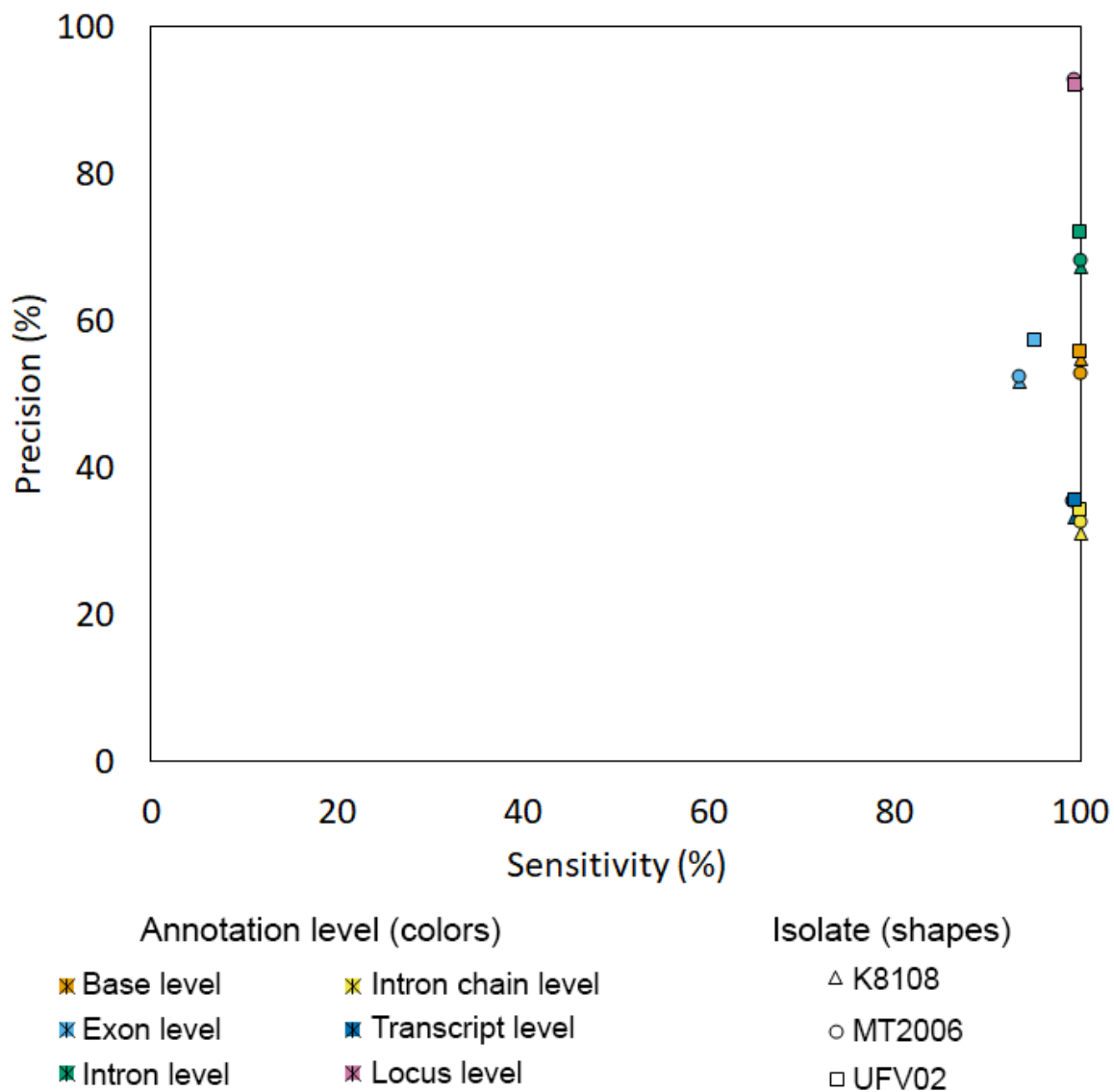
**A.** DEGs were hierarchical clustered by treatment, applying hclust method using R package (150). The germinated-spore (2) was used as calibrator. The other conditions are (1) Spore; (3) appressorium *in vitro*; (4) appressorium *in planta*; (5) 10–12 HPI; (6) 24h HPI; (7) 36 HPI; (8) 72 HPI; (9) 168 HPI; (10) 192–196 HPI; (11) 288 HPI. **B.** Principal component analysis of transcriptomic data.



1321  
1322 **Fig. S18. Heatmap of common DEGs between the *P. pachyrhizi* genomes K8108, MT2006**  
1323 **and UFV02.**  
1324 DEGs were hierarchical clustered by treatment, applying hclust method using R package  
1325 (142). The germinated-spore (2) was used as calibrator. The other conditions are (1) Spore;  
1326 (3) appressorium *in vitro*; (4) appressorium *in planta*; (5) 10-12 HPI; (6) 24 HPI; (7) 36 HPI; (8)  
1327 72 HPI; (9) 168 HPI; (10) 192-196 HPI; (11) 288 HPI.



1329 **Fig. S19. Microscopic images of soybean leaf tissue (magenta) infected with *P. pachyrhizi***  
1330 **(green) at different time points of the infection.**  
1331 **(A) 12 HPI (B) 24 HPI (C) 32 HPI (D) 72 HPI (E) 168 HPI (F) 192 HPI.** Shown are 3D images  
1332 obtained from z-stacks. Red, green and blue frames indicate sites of clipping to reveal areas  
1333 inside the leaf. Scale bars represent 20  $\mu\text{m}$  (A, B, C, D) and 50  $\mu\text{m}$  (E, F).  
1334

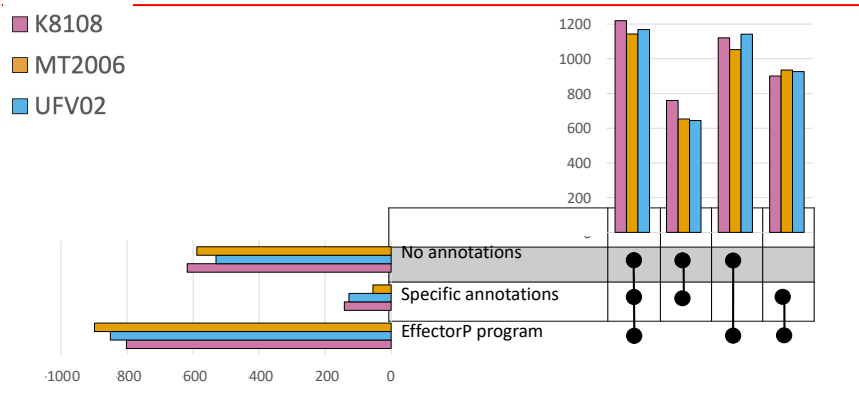


1335  
1336  
1337  
1338

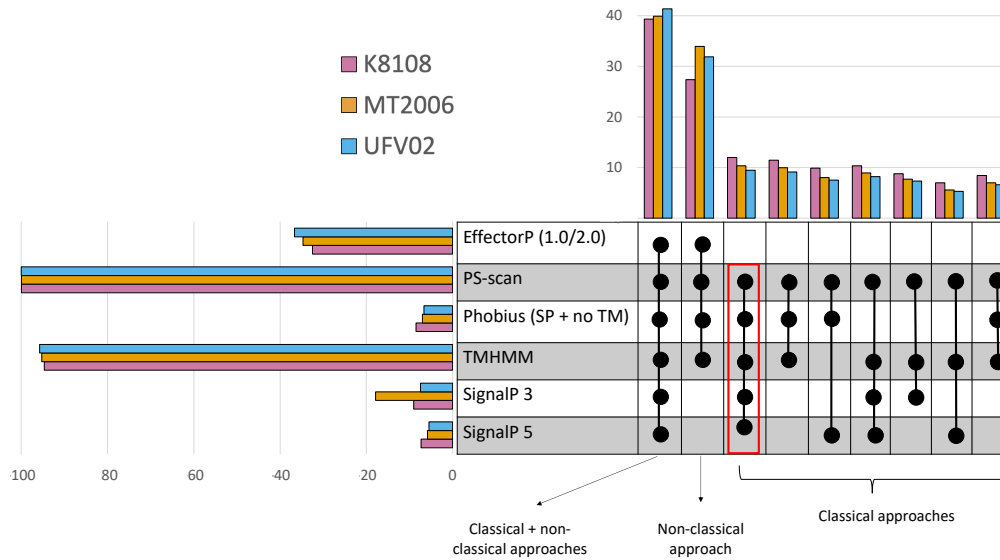
**Fig. S20. Validation of gene annotation based on expression data from the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**



A

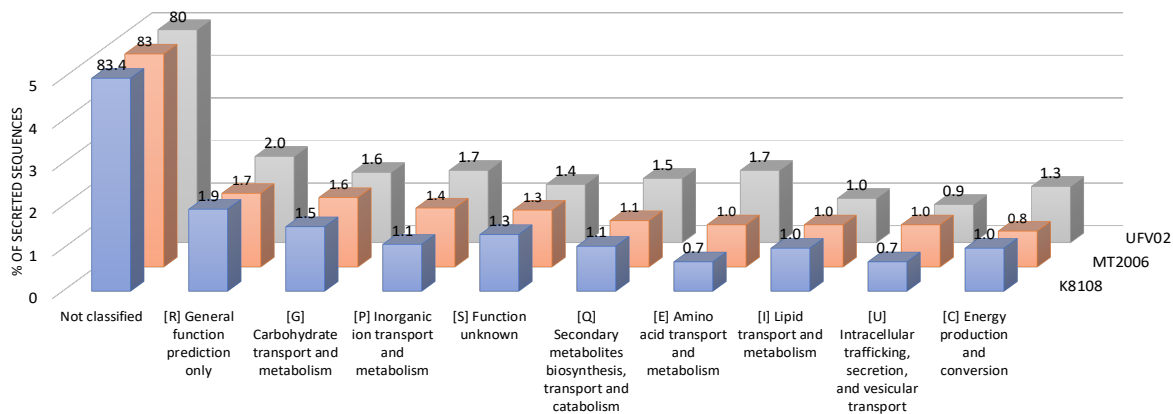


B



1339  
1340  
1341  
1342

**Fig. S21. Prediction of secreted proteins from the *P. pachyrhizi* genomes K8108, MT2006 and UFV02 using different effector prediction tools.**



1343  
 1344 **Fig. S22. Gene categories of the secreted proteins from the *P. pachyrhizi* genomes K8108,**  
 1345 **MT2006 and UFV02.**

- 1346 **Supplementary tables:**
- 1347 **Table S1. Summary metrics of TE annotation in the *P. pachyrhizi* genomes K8108, MT2006**
- 1348 **and UFV02.**
- 1349 **Table S2. Complete TE annotation in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02**
- 1350 **Table S3. Conserved TEs in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**
- 1351 **Table S4. Intermediate TEs in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**
- 1352 **Table S5. Divergent TEs in the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**
- 1353 **Table S6. List of candidate effectors from K8108 isolate.**
- 1354 **Table S7. List of candidate effectors from MT2006 isolate.**
- 1355 **Table S8: List of candidate effectors from UFV02 isolate.**
- 1356 **Table S9. Expression profile of common secreted genes in the 3 Phapa transcriptomes.**
- 1357 **Table S10. Number of expressed TEs in the *P. pachyrhizi* genomes K8108, MT2006 and**
- 1358 **UFV02 per order and superfamily.**
- 1359 **Table S11. Number of expressed TEs per conditions in the *P. pachyrhizi* genomes K8108,**
- 1360 **MT2006 and UFV02.**
- 1361 **Table S12. Expression of TEs under different conditions in K8108 isolate.**
- 1362 **Table S13. Expression of TEs under different conditions in MT2006 isolate.**
- 1363 **Table S14. Expression of TEs under different conditions in UFV02 isolate.**
- 1364 **Table S15a. Summary and functional impact of variants and in the *P. pachyrhizi* genomes**
- 1365 **K8108, MT2006 and UFV02.**
- 1366 **Table S15b. Predication of the SNP impact in the *P. pachyrhizi* genome.**
- 1367 **Table S16. Haplotype phasing of the *P. pachyrhizi* genomes K8108, MT2006 and UFV02.**
- 1368 **Table S17. Differentially expressed genes in K8108 transcriptome.**
- 1369 **Table S18. Differentially expressed genes in MT2006 transcriptome.**
- 1370 **Table S19. Differentially expressed genes in UFV02 transcriptome.**
- 1371 **Table S20a. Summary of the fungal species used for the MCL and CAFÉ analysis.**
- 1372 **Table S20b. Dated tree with time to Most Recent Common Ancestor (tMRCA).**
- 1373 **Table S21. Distribution of contracted gene families in 15 different fungal species.**
- 1374 **Table S22. Distribution of expanded gene families in 15 different fungal species.**
- 1375 **Table S23a. Summary of fungal species used for the CAZyme comparisons.**
- 1376 **Table S23b. Summary of the CAZyme profile in fungal species.**
- 1377 **Table S23C. Summary of the CAZyme families in fungal species.**
- 1378 **Table S24. Summary metrics of *P. pachyrhizi* genome annotations and gene models.**
- 1379 **Table S25. Expression profile of common genes shared by the *P. pachyrhizi* transcriptomes**
- 1380 **K8108, MT2006 and UFV02.**
- 1381 **Table S26. Allelic Correspondence among the *P. pachyrhizi* genomes K8108, MT2006 and**
- 1382 **UFV02 gene catalogues.**
- 1383 **Table S27. Precision and sensitivity of the gene annotations in the *P. pachyrhizi* genomes**
- 1384 **K8108, MT2006 and UFV02.**
- 1385 **Table S28. List of authors and their contributions.**