# A population-level invasion by transposable elements triggers genome expansion in a fungal pathogen

Ursula Oggenfuss[1], Thomas Badet[1], Thomas Wicker[2], Fanny E. Hartmann[3,4], Nikhil K. Singh[1], Leen

N. Abraham[1], Petteri Karisto[4,6], Tiziana Vonlanthen[4], Christopher C. Mundt[5], Bruce A. McDonald[4],

Daniel Croll[1,*]

[1] Laboratory of Evolutionary Genetics, Institute of Biology, University of Neuchâtel, 2000 Neuchâtel, Switzerland
[2] Institute for Plant and Microbial Biology, University of Zurich, Zurich, Switzerland
[3] Ecologie Systématique Evolution, Bâtiment 360, Univ. Paris-Sud, AgroParisTech, CNRS, Université Paris-Saclay, 91400 Orsay, France
[4] Plant Pathology, Institute of Integrative Biology, ETH Zurich, Zurich, Switzerland
[5] Department of Botany and Plant Pathology, Oregon State University, Corvallis, OR 97331-2902, USA
[6] Department of Evolutionary Biology and Environmental Studies, University of Zurich, Zurich, Switzerland

[*] Author for correspondence: daniel.croll@unine.ch

Running title: Transposable element invasion triggers genome expansion

31 **ABSTRACT**

32 Genome evolution is driven by the activity of transposable elements (TEs). The spread of TEs can

33 have deleterious effects including the destabilization of genome integrity and expansions. However,

34 the precise triggers of genome expansions remain poorly understood because genome size evolution

35 is typically investigated only among deeply divergent lineages. Here, we use a large population

36 genomics dataset of 284 individuals from populations across the globe of *Zymoseptoria tritici*, a

37 major fungal wheat pathogen. We built a robust map of genome-wide TE insertions and deletions to

38 track a total of 2,456 polymorphic loci within the species. We show that purifying selection

39 substantially depressed TE frequencies in most populations but some rare TEs have recently risen in

40 frequency and likely confer benefits. We found that specific TE families have undergone a

41 substantial genome-wide expansion from the pathogen's center of origin to more recently founded

42 populations. The most dramatic increase in TE insertions occurred between a pair of North

43 American populations collected in the same field at an interval of 25 years. We find that both

44 genome-wide counts of TE insertions and genome size have increased with colonization bottlenecks.

45 Hence, the demographic history likely played a major role in shaping genome evolution within the

46 species. We show that both the activation of specific TEs and relaxed purifying selection underpin

47 this incipient expansion of the genome. Our study establishes a model to recapitulate TE-driven

48 genome evolution over deeper evolutionary timescales.

49

2

## INTRODUCTION

50

51    Transposable elements (TEs) are mobile repetitive DNA sequences with the ability to independently

52    insert into new regions of the genome. TEs are major drivers of genome instability and epigenetic

53    change (Eichler & Sankoff, 2003). Insertion of TEs can disrupt coding sequences, trigger

54    chromosomal rearrangements, or alter expression profiles of adjacent genes (Lim, 1988; Petrov *et*

55    *al.*, 2003; Slotkin & Martienssen, 2007; Hollister & Gaut, 2009; Oliver *et al.*, 2013). Hence, TE

56    activity can have phenotypic consequences and impact host fitness. While TE insertion dynamics are

57    driven by the selfish interest for proliferation, the impact on the host can range from beneficial to

58    highly deleterious. The most dramatic examples of TE insertions underpinned rapid adaptation of

59    populations or species (Feschotte, 2008; Chuong *et al.*, 2017), particularly following environmental

60    change or colonization events. Beneficial TE insertions are expected to experience strong positive

61    selection and rapid fixation in populations. However, most TE insertions have neutral or deleterious

62    effects upon insertions. Purifying selection is expected to rapidly eliminate deleterious insertions

63    from populations unless constrained by genetic drift (Walser *et al.*, 2006; Baucom *et al.*, 2008;

64    Cridland *et al.*, 2013; Stuart *et al.*, 2016; Lai *et al.*, 2017; Stritt *et al.*, 2017). Additionally, genomic

65    defense mechanisms can disable transposition activity. Across eukaryotes, epigenetic silencing is a

66    shared defense mechanism against TEs (Slotkin & Martienssen, 2007). Fungi evolved an additional

67    and highly specific defense system introducing repeat-induced point (RIP) mutations into any nearly

68    identical set of sequences. The relative importance of demography, selection and genomic defenses

69    determining the fate of TEs in populations remain poorly understood.

70

71    A crucial property predicting the invasion success of TEs in a genome is the transposition rate. TEs

72    tend to expand through family-specific bursts of transposition followed by prolonged phases of

73    transposition inactivity. Bursts of insertions of different retrotransposon families were observed

74    across eukaryotic lineages including *Homo sapiens*, *Zea mays*, *Oryza sativa* and *Blumeria graminis*

75    (Shen *et al.*, 1991; SanMiguel *et al.*, 1998; Eichler & Sankoff, 2003; Piegu *et al.*, 2006; Lu *et al.*,

76    2017; Frantzeskakis *et al.*, 2018). Prolonged bursts without effective counter-selection are thought to

77    underpin genome expansions. In the symbiotic fungus *Cenococcum geophilum*, the burst of TEs

78    resulted in a dramatically expanded genome compared to closely related species (Peter *et al.*, 2016).

79    Similarly, a burst of a TE family in brown hydras led to an approximately three-fold increase of the

80    genome size compared to related hydras (Wong *et al.*, 2019). Across the tree of life, genome sizes

81    vary by orders of magnitude and enlarged genomes invariably show hallmarks of historic TE

82    invasions (Kidwell, 2002). Population size variation is among the few correlates of genome size

83    across major groups, suggesting that the efficacy of selection plays an important role in controlling

84    TE activity (Lynch, 2007). Reduced selection efficacy against deleterious TE insertions is expected

85    to lead to a ratchet-like increase in genome size. In fungi, TE-rich genomes often show an isochore

86    structure alternating gene-rich and TE-rich compartments (Rouxel *et al.*, 2011). TE-rich

87    compartments often harbor rapidly evolving genes such as effector genes in pathogens or resistance

88    genes in plants (Raffaele & Kamoun, 2012; Jiao & Schneeberger, 2019). Taken together, incipient

89    genome expansions are likely driven by population-level TE insertion dynamics.

90

91    The fungal wheat pathogen *Zymoseptoria tritici* is one of the most important pathogens on crops,

92    causing high yield losses in many years (Torriani *et al.*, 2015). *Z. tritici* emerged during the

93    domestication of wheat in the Fertile Crescent where the species retained high levels of genetic

94    variation (Zhan *et al.*, 2005; Stukenbrock *et al.*, 2011). The pathogen migrated to all temperate zones

95    where wheat is currently grown and underwent multiple migration bottlenecks, in particular when

96    colonizing Oceania and North America (Zhan *et al.*, 2005; Estep *et al.*, 2015). The genome is

97    completely assembled and shows size variation between individuals sampled across the global

98    distribution range (Feurtey *et al.*, 2020; Badet *et al.*, 2020) (Goodwin *et al.*, 2011). The TE content

99    of the genome shows a striking variation of 17-24% variation among individuals (Badet *et al.*,

100   2020). *Z. tritici* recently gained major TE-mediated adaptations to colonize host plants and tolerate

101   environmental stress (Omrane *et al.*, 2015, 2017; Krishnan *et al.*, 2018; Meile *et al.*, 2018). Clusters

102   of TEs are often associated with genes encoding important pathogenicity functions (*i.e.* effectors),

103   recent gene gains or losses (Hartmann & Croll, 2017), and major chromosomal rearrangements

104   (Croll *et al.*, 2013; Plissonneau *et al.*, 2016). Transposition activity of TEs also had a genome-wide

4

105 impact on gene expression profiles during infection (Fouché *et al.*, 2019). The well-characterized

106 demographic history of the pathogen and evidence for recent TE-mediated adaptations make *Z.*

107 *tritici* an ideal model to recapitulate the process of TE insertion dynamics, adaptive evolution and

108 changes in genome size at the population level.

109

110 Here, we retrace the population-level context of TE insertion dynamics and genome size changes

111 across the species range by analyzing populations sampled on four continents for a total of 284

112 genomes. We developed a robust pipeline to detect newly inserted TEs using short read sequencing

113 datasets. Combining analyses of selection and knowledge of the colonization history of the

114 pathogen, we tested whether population bottlenecks were associated with substantial changes in the

115 TE content and the size of genomes.

116

117

118 **RESULTS**

119 A DYNAMIC TE LANDSCAPE SHAPED BY STRONG PURIFYING SELECTION

120 We detected 4,753 TE copies, grouped into 30 families with highly variable copy numbers in the

121 reference genome IPO323 (Figure 2 - figure supplement 1 and Figure 2 - figure supplement 2A). To

122 establish a comprehensive picture of within-species TE dynamics, we analyzed 295 genomes from a

123 worldwide set of six populations spanning the  distribution range of the wheat pathogen *Z. tritici*. To

124 ascertain the presence or absence of TEs across the genome, we developed a robust pipeline (Figure

125 1A). In summary, we called TE insertions by identifying reads mapping both to a TE sequence and a

126 specific location in the reference genome. Then, we assessed the minimum sequencing coverage to

127 reliably recover TE insertions and removed 11 genomes with an average read depth below 15X

128 (Figure 1B). We tested for evidence of TEs using read depth at target site duplications (Figure 1C)

129 and scanned the genome for mapped reads indicating gaps at TE loci (Figure 1D). We found robust

130 evidence for a total of 18,864 TE insertions grouping into 2,465 individual loci. Of these loci, 35.5%

131 ($n$ = 876) have singleton TEs (*i.e.*, this locus is only present in one isolate: Figure 2A, figure

5

132  supplement 3). An overwhelming proportion of loci (2,345 loci or 95.1%) have a TE frequency

133  below 1%. Singleton TE insertions in particular can be the product of spurious Illumina read

134  mapping errors (Nakamura *et al.*, 2011). To assess the reliability of the detected singletons, we

135  focused on seven isolates for which PacBio long-read data was available (Badet *et al.*, 2020).

136  Aligned PacBio reads confirmed the exact location of 71% (22 out of 31 singleton insertions among

137  seven isolates; see Methods for further details). We found no significant difference in read coverage

138  between confirmed and unconfirmed singleton insertions (Figure 2 - figure supplement 2C,-B and

139  Figure 2 - figure supplement 4).

140

141  The abundance of singleton TE insertions strongly supports the idea that TEs actively copy into new

142  locations but also indicates that strong purifying selection maintains nearly all TEs at low frequency

143  (Figure 2A). The density of TE loci on accessory chromosomes, which are not shared among all

144  isolates of the species, is almost twice the density found on core chromosomes (102 *versus* 58 TEs

145  per Mb; Figure 2B and Figure 2 - figure supplement 5A). This suggests relaxed selection against TE

146  insertion on the functionally dispensable and gene-poor accessory chromosomes. We found no

147  difference in TE allele frequency distribution between recombination hotspots and the rest of the

148  genome (Figure 2 - figure supplement 5B). Similarly, the TE density and the number of insertions

149  did not vary between recombination hotspots and the genomic background (Figure 2 - figure

150  supplement 5C).

151

152  TEs grouped into 23 families and 11 superfamilies, with 88.2% of all copies belonging to class

153  I/retrotransposons ($n = 2175$; Figure 2C and Figure 2 - figure supplements 6A-B). RLG/*Gypsy* ($n =$

154  1,483) and RLC/*Copia* ($n = 623$) elements constitute the largest long terminal repeats (LTR)

155  superfamilies. Class II/DNA transposons are dominated by DHH/*Helitron* ($n = 249$). As expected,

156  TE families shared among fewer isolates tend to show also lower global copy numbers (*i.e.*, all

157  isolates combined), while TE families that are present in all isolates generally have high global copy

158  numbers (Figure 2D).

159

160    We detected 153 loci with TEs inserted into genes with most of the insertions being singletons

161    (44.7%; $n$ = 68) or of very low frequency (Figure 2E). Overall, TE insertions into exonic sequences

162    were less frequent than expected compared to insertions into up- and downstream regions, which is

163    consistent with effective purifying selection (Figure 2F). Insertions into introns were also strongly

164    under-represented, likely due to the small size of most fungal introns (~ 50-100 bp) and the high

165    probability of disrupting splicing or adjacent coding sequences. We also found that insertions 800-

166    1000 bp away from coding sequences of a focal gene were under-represented. Given the high gene

167    density, with an average spacing between genes of 1.744 kb, TE insertions within 800-1,000 bp of a

168    coding gene tend to be near adjacent genes already. Taken together, TEs in the species show a high

169    degree of transposition activity and are subject to strong purifying selection.

170

171    DETECTION OF CANDIDATE TE LOCI UNDERLYING RECENT ADAPTATION

172    The TE transposition activity can generate adaptive genetic variation. To identify the most likely

173    candidate loci, we analyzed insertion frequency variation among populations as an indicator for

174    recent selection. Across all populations, the insertion frequencies differed only weakly with a strong

175    skew towards extremely low $F_{ST}$ values (mean = 0.0163; Figure 3A-B and Figure 3 - figure

176    supplement 1). To further analyze evidence for TE-mediated adaptive evolution, we screened a

177    genome-wide SNP dataset for evidence of selective sweeps using selection scans. We found 16.5 %

178    of all TE loci located in regions of selective sweep. Given our population sampling of two

179    population pairs, we tested for adaptive TE insertions in selective sweep regions either in the North

180    American or European population pairs. Hence, we selected loci having low TE insertion

181    frequencies (< 5%) in all populations except either the recent North American or European

182    population (> 20%) (Figure 3B). Based on these criteria, we obtained 7 candidate loci possibly

183    underlying local adaptation (6 in North America, one in Europe; Figure 4A and Figure 4 - figure

184    supplement 1). All loci carry inserted retrotransposons with 4 RLG_Luna, one RLG_Mercurius and

185    one RLG_Deimos.

186

187     One TE insertion is 3,815 bp downstream of a gene encoding an RTA1-like protein, which can

188     function as transporters with a transmembrane domain and have been associated with resistance

189     against several antifungal compounds (Soustre *et al.*, 1996). The insertion is also 5785 bp upstream

190     of a gene encoding a protein kinase domain (Figure 4B). The TE insertion was not detected in the

191     Middle East or the two European populations, and was at low frequencies in the Australian (3.7%)

192     and North American 1990 (1.7%) populations, but increased to 53% of all isolates in the North

193     American 2015 population (fixation index $F_{ST}$ = 0.42; Figure 4 - figure supplement 1). Isolates that

194     carry the insertion show a significantly higher resistance to azole antifungal compounds (Figure 4C).

195     The TE is in the subtelomeric region of chromosome 12, with a moderate GC content, a low TE and

196     a high gene density (Figure 4D). The TE belongs to the family RLG_Luna, which shows a

197     substantial burst across different chromosomes within the species (Figures 4E-F). We found no

198     association between the phylogenetic relationships among isolates based on the two closest genes

199     and the presence or absence of the TE insertion (Figure 4G). A second candidate adaptive TE

200     insertion belongs to the RLG_Mercurius family and is located between two genes of unknown

201     function (Figure 4 - figure supplement 2). A third potentially adaptive TE insertion of a

202     RLC_Deimos is 229 bp upstream of a gene encoding a SNARE domain protein and 286 bp upstream

203     of a gene encoding a flavin amine oxidoreductase. Furthermore, the TE is inserted in a selective

204     sweep region (Figure 4 - figure supplement 2). SNARE domains play a role in vesicular transport

205     and membrane fusion (Bonifacino & Glick, 2004). An additional four candidates for adaptive TE

206     insertions belong to RLG_Luna and were located distantly to genes (Figure 4 - figure supplement 2).

207     We experimentally tested whether the TE insertions in proximity to genes were associated with

208     higher levels of fungicide resistance. For this, we measured growth rates of the fungal isolates in the

209     presence or absence of an azole fungicide widely deployed against the pathogen. We found that the

210     insertion of TEs at two loci was positively associated with higher levels of fungicide resistance,

211     suggesting that the adaptation was mediated by the TE (Figure 4C and Figure 4 - figure supplement

212     2).

213

214 POPULATION-LEVEL EXPANSIONS IN TE CONTENT

215 If TE insertion dynamics are largely neutral across populations, TE frequencies across loci should

216 reflect neutral population structure. To test this, we performed a principal component analysis based

217 on a set of six populations on four continents that represent the global genetic diversity of the

218 pathogen (Figure 5A) and 900,193 genome-wide SNPs (Figure 5B). The population structure

219 reflected the demographic history of the pathogen with clear continental differentiation and only

220 minor within-site differentiation. To account for the lower number of TE loci, we performed an

221 additional principal component analysis using a random SNP set of similar size to the number of TE

222 loci. The reduced SNP set retained the geographic signal of the broader set of SNPs (Figure 5C). In

223 stark contrast, TE frequencies across loci showed only weak clustering by geographic origin with the

224 Australian population being the most distinct (Figure 5D). We found a surprisingly strong

225 differentiation of the two North American populations sampled at a 25-year interval in the same field

226 in Oregon.

227

228 Unusual patterns in population differentiation at TE loci suggests that TE activity may substantially

229 vary across populations (Figure 6, Figure 4 - figure supplement 1). To analyze this, we first

230 identified the total TE content across all loci per isolate. We found generally lower TE numbers in

231 the Middle Eastern population from Israel (Figure 6A-C, and Figure 6 - figure supplement 1), which

232 is close to the pathogen's center of origin (Stukenbrock *et al.*, 2007). Populations that underwent at

233 least one migration bottleneck showed a substantial burst of TEs across all major superfamilies.

234 These populations included the two populations from Europe collected in 1999 and 2016 and the

235 North American population from 1990, as well as the Australian population. We found a second

236 stark increase in TE content in the North American population sampled in 2015 at the same site as

237 the population from 1990. Strikingly, the isolate with the lowest number of analyzed TEs collected

238 in 2015 was comparable to the isolate with the highest number of TEs at the same site in 1990. We

239 tested whether sequencing coverage could explain variation in the detected TEs across isolates, but

240 we found no meaningful association (Figure 2 - figure supplement 6C). We analyzed whether the

241 population-specific expansions were correlated with shifts in the frequency spectrum of TEs in the

9

242 populations (Figure 6D). We found that the first step of expansions observed in Europe compared to

243 the Middle East (Israel) was associated with an upwards shift in allele frequencies. This is consistent

244 with transposition activity creating new copies in the genomes and stronger purifying selection in the

245 Middle East. Similarly, the North American populations showed also signatures consistent with

246 relaxation of selection against TEs (*i.e.*, fewer low frequency TEs). We found a significant

247 difference (Two-sample Kolmogorov-Smirnov test, two-sided) in the curve shapes between the

248 population from the Middle East and North America 2015 (Figure 6 - figure supplement 2). We

249 analyzed variation in TE copy numbers across families and found that the expansions were mostly

250 driven by RLG elements including the families Luna, Sol and Venus, the RLC family Deimos and

251 the LINE family Lucy (Figure 6E and Figure 6 - figure supplement 3A). We also found a North

252 American specific burst in DHH elements of the family Ada (increase from 4.6 to 6.1 copies on

253 average per isolate), an increase specific to Swiss populations in LINE elements, and an increase in

254 RLC elements in the Australian and the two North American populations. Analyses of complete *Z.*

255 *tritici* reference-quality genomes that include isolates from the Israel, Australia, Switzerland (1999)

256 and North American (1990) population revealed high TE contents in Australia and North America

257 (Oregon 1990) (Badet *et al.*, 2020). The reference-quality genomes confirmed also that the increase

258 in TEs was driven by LINE, RLG and RLC families in Australia and DHH, RLG and RLC families

259 in North America (Badet *et al.*, 2020).

260

261 TE-MEDIATED GENOME SIZE EXPANSIONS

262 The combined effects of actively copying TE families and relaxed purifying selection leads to an

263 accumulation of new TE insertions in populations. Consequently, mean genome sizes in populations

264 should increase over generations. We estimated the cumulative length of TE insertions based on the

265 length of the corresponding TE consensus sequences and found a strong increase in the total TE

266 length in populations outside the Middle East center of origin, and a second increase between the

267 two North American populations (Figure 1 - figure supplement 1). To test for incipient genome

268 expansions within the species, we first assembled genomes of all 284 isolates included in the study.

269 Given the limitations of short-read assemblies, we implemented corrective measures to compensate

270     for potential variation in assembly qualities. We corrected for variation in the GC content of

271     different sequencing datasets by downsampling reads to generate balanced sequencing read sets prior

272     to assembly (see Methods). We also excluded all reads mapping to accessory chromosomes because

273     different isolates are known to differ in the number of these chromosomes. Genome assemblies were

274     checked for completeness by retrieving the phylogenetically conserved BUSCO genes (Figure 7A).

275     Genome assemblies across different populations carry generally >99% complete BUSCO gene sets,

276     matching the completeness of reference-quality genomes of the same species (Badet *et al.*, 2020).

277     The completeness of the assemblies showed no correlation with either TE or GC content of the

278     genomes. GC content was inversely correlated with genome size consistent with the expansion of

279     repetitive regions having generally low GC content (Figure 7B). We found that the core genome size

280     varied substantially among populations with the Middle East, Australia as well as the two older

281     European and North American populations having the smallest core genome sizes (Figure 7C). We

282     found a notable increase in core genome size in both the more recent European and North American

283     populations. The increase in core genome size is positively correlated with the count and cumulative

284     length of all inserted TEs (Figure 7D, 7E and 7G) and negatively correlated with the genome-wide

285     GC content (Figure 7F and 7G). Hence, core genome size shows substantial variation within the

286     species matching the recent expansion in TEs across continents. We found the most variable genome

287     sizes in the more recent North American population (Figure 7 - figure supplement 1B). Finally, we

288     contrasted variation in genome size with the detected TE insertion dynamics. For this, we assessed

289     the variable genome segment as the difference between the smallest and largest analyzed core

290     genome. To reflect TE dynamics, we calculated the cumulative length of all detected TE insertions

291     in any given genome. We found that the cumulative length of inserted TEs represents between 4.8

292     and 184 % of the variable genome segment defined for the species or 0.2-2.6% of the estimated

293     genome size per isolate (Figure 7 - figure supplement 1C-D).

294

## DISCUSSION

296   TEs play a crucial role in generating adaptive genetic variation within species but are also drivers of

297   deleterious genome expansions. We analyzed the interplay of TEs with selective and neutral

298   processes including population differentiation and incipient genome expansions. TEs have

299   substantial transposition activity in the genome but are strongly counter-selected and are maintained

300   at low frequency. TE dynamics showed distinct trajectories across populations with more recently

301   established populations having higher TE content and a concurrent expansion of the genome.

302

303   RECENT SELECTION ACTING ON TE INSERTIONS

304   TE frequencies in the species show a strong skew towards singleton insertions across populations.

305   However, our short read based analyses are possibly skewed towards over-counting singletons as

306   indicated by independent long-read mapping evaluations. Nevertheless, the skew towards low

307   frequency TE insertions indicates both that TEs are undergoing transposition and that purifying

308   selection maintains frequencies at a low level. Similar effects of selection on active TEs were

309   observed across plants and animals, including *Drosophila melanogaster* and *Brachypodium*

310   *distachyon* (Cridland *et al.*, 2013; Stritt *et al.*, 2017; Luo *et al.*, 2020). TE insertions were under-

311   represented in or near coding regions, showing a stronger purifying selection against TEs inserting

312   into genes. Coding sequences in the *Z. tritici* genome are densely packed with an average distance of

313   only ~1 kb (Goodwin *et al.*, 2011). Consistent with this high gene density, TE insertions were most

314   frequent at a distance of 200-400 bp away from coding sequences. A rapid decay in linkage

315   disequilibrium in the *Z. tritici* populations (Croll *et al.*, 2015; Hartmann *et al.*, 2018) likely

316   contributed to the efficiency of removing deleterious insertions. Some TE superfamilies have

317   preferred insertion sites in coding regions and transcription start sites (Miyao *et al.*, 2003; Fu *et al.*,

318   2013; Gilly *et al.*, 2014; Quadrana *et al.*, 2016). Hence, some heterogeneity in the observed insertion

319   site distribution across the genome is likely due to insertion preferences of individual TEs. We also

320   found evidence for positive selection acting on TEs with the strongest candidate locus being a TE

321   insertion on chromosome 12. This locus showed a frequency increase only in the more recent North

322    American population, which experienced the first systematic fungicide applications and subsequent

323    emergence of fungicide resistance in the decade prior to the last sampling (Estep *et al.*, 2015). The

324    nearest gene encodes a RTA1-like protein, a transmembrane exporter which is associated with

325    resistance towards different stressors, including antifungal compounds, and shows strong copy

326    number variation in several fungi (Soustre *et al.*, 1996; Rogers & Barker, 2003; Sirisattha *et al.*,

327    2004; Ali *et al.*, 2014; Yew *et al.*, 2016; Liang *et al.*, 2018). Hence, the TE insertion may have

328    positively modulated RTA1 expression to resist antifungals.

329    Transposition activity in a genome and counter-acting purifying selection are expected to establish

330    an equilibrium over evolutionary time (Charlesworth & Charlesworth, 1983). However, temporal

331    bursts of TE families and changes in population size due to bottlenecks or founder events are likely

332    to shift the equilibrium. Despite purifying selection, we were able to detect signatures of positive

333    selection by scanning for short-term population frequency shifts. Population genomic datasets can be

334    used to identify the most likely candidate loci underlying recent adaptation. The shallow genome-

335    wide differentiation of *Z. tritici* populations provides a powerful background to test for outlier loci

336    (Hartmann *et al.*, 2018). We found the same TE families to have experienced genome-wide copy

337    number expansions, suggesting that the availability of adaptive TE insertions may be a by-product of

338    TE bursts in individual populations.

339

340    POPULATION-LEVEL TE INVASIONS AND RELAXED SELECTION

341    Across the surveyed populations from four continents, we identified substantial variation in TE

342    counts per genome. The increase in TEs matches the global colonization history of the pathogen

343    with an increase in TE copies in more recently established populations (Zhan *et al.*, 2003;

344    Stukenbrock *et al.*, 2007). Compared to the Israeli population located nearest the center of origin in

345    the Middle East, the European populations showed a three-fold increase in TE counts. The

346    Australian and North American populations established from European descendants retained high

347    TE counts. We identified a second increase at the North American site where TE counts nearly

348    doubled again over a 25-year period. Compared to the broader increase in TEs from the Middle East,

349    the second expansion at the North American site was driven by a small subset of TE families alone.

13

350    Analyses of completely assembled reference-quality genomes from the same populations confirmed

351    that genome expansions were primarily driven by the same TE families belonging to the RLG, RLC

352    and DHH superfamilies (Badet *et al.*, 2020). Consistent with the contributions from individual TEs,

353    we found that the first expansion in Europe led to an increase in low-frequency variants, suggesting

354    higher transposition activity of many TEs in conjunction with strong purifying selection. The second

355    expansion at the North American site shifted TE frequencies upwards, suggesting relaxed selection

356    against TEs. The population-level context of TEs in *Z. tritici* shows how heterogeneity in TE control

357    interacts with demography to determine extant levels of TE content and, ultimately, genome size.

358

359    TE INVASION DYNAMICS UNDERPINS GENOME SIZE EXPANSIONS

360    The number of detected TEs was closely correlated with core genome size, hence genome size

361    expansions were at least partly caused by the very recent proliferation of TEs. Genome assemblies of

362    large eukaryotic genomes based on short read sequencing are often fragmented and contain chimeric

363    sequences (Nagarajan & Pop, 2013). Focusing on the less repetitive core chromosomes in the

364    genome of *Z. tritici* reduces such artefacts substantially. Because genome assemblies are the least

365    complete in the most repetitive regions, any underrepresented sequences may rather underestimate

366    than overestimate within-species variation in genome size. Hence, we consider the assembly sizes to

367    be a robust correlate of total genome size. The core genome size differences observed across the

368    species range match genome size variation typically observed among closely related species. Among

369    primates, genome size varies by ~70% with ~10% between humans and chimpanzees (Rogers &

370    Gibbs, 2014; Miga *et al.*, 2020). In fungi, genome size varies by several orders of magnitude within

371    phyla but is often highly similar among closely related species (Raffaele & Kamoun, 2012).

372    Interestingly, drastic changes in genome size have been observed in the *Blumeria* and

373    *Pseudocercospora* genera where genome size changed by 35-130% between the closest known

374    species (González-Sayer *et al.*; Frantzeskakis *et al.*, 2018). Beyond analyses of TE content variation

375    correlating with genome size evolution, proximate mechanisms driving genome expansions are

376    poorly understood. By establishing large population genetic datasets, such as those possible for crop

377    pathogens, analyses of genome size evolution become tractable at the population level.

378    TEs might not only contribute to genome expansion directly by adding length through additional

379    copies, but also by increasing the rate of chromosomal rearrangements and ectopic recombination

380    (Bourque *et al.*, 2018; Blommaert, 2020). However, TEs are not the only repetitive elements that can

381    lead to a genome size expansion. In *Arabidopsis thaliana* genomes, the 45S rDNA has been shown

382    to have the strongest impact on genome size variation, followed by 5S rDNA variation, and

383    contributions by centromeric repeats and TEs (Long *et al.*, 2013). In conjunction, recent work

384    demonstrates how repetitive sequences are drivers of genome size evolution over short evolutionary

385    timescales.

386    The activity of TEs is controlled by complex selection regimes within species. Actively transposing

387    elements may accelerate genome evolution and underpin expansions. Hence, genomic defenses

388    should evolve to efficiently target recently active TEs. Here, we show that TE activity and

389    counteracting genomic defenses have established a tenuous equilibrium across the species range. We

390    show that population subdivisions are at the origin of highly differentiated TE content within a

391    species matching genome size changes emerging over the span of only decades and centuries. In

392    conclusion, population-level analyses of genome size can recapitulate genome expansions typically

393    observed across much deeper time scales providing fundamentally new insights into genome

394    evolution.

395

396    **METHODS**

397    FUNGAL ISOLATE COLLECTION AND SEQUENCING

398    We analyzed 295 *Z. tritici* isolates covering six populations originating from four geographic

399    locations and four continents (Figure 5 - figure supplement 1), including: Middle East 1992 ($n = 30$

400    isolates, Nahal Oz, Israel), Australia 2001 ($n = 27$, Wagga Wagga), Europe 1999 ($n = 33$, Berg am

401    Irchel, Switzerland), Europe 2016 ($n = 52$, Eschikon, ca. 15km from Berg am Irchel, Switzerland),

402    North America 1990 and 2015 ($n = 56$ and $n = 97$, Willamette Valley, Oregon, United States)

403    (McDonald *et al.*, 1996; Linde *et al.*, 2002; Zhan *et al.*, 2002, 2003, 2005). Illumina short read data

404    from the Middle Eastern, Australian, European 1999 and North American 1990 populations were

15

405 obtained from the NCBI Sequence Read Archive (SRA) under the BioProject PRJNA327615

406 (Hartmann *et al.*, 2017). For the Switzerland 2016 and Oregon 2015 populations, asexual spores

407 were harvested from infected wheat leaves from naturally infected fields and grown in YSB liquid

408 media including 50 mgL$^{-1}$ kanamycin and stored in silica gel at −80°C. High-quality genomic DNA

409 was extracted from liquid cultures using the DNeasy Plant Mini Kit from Qiagen (Venlo,

410 Netherlands). The isolates were sequenced on an Illumina HiSeq in paired-end mode and raw reads

411 were deposited at the NCBI SRA under the BioProject PRJNA596434.

412

413 TE INSERTION DETECTION

414 The quality of Illumina short reads was determined with FastQC version 0.11.5

415 (https://www.bioinformatics.babraham.ac.uk/projects/fastqc/) (Figure 1A). To remove spuriously

416 sequenced Illumina adaptors and low quality reads, we trimmed the sequences with Trimmomatic

417 version 0.36, using the following filter parameters: illuminaclip:TruSeq3-PE-2.fa:2:30:10 leading:10

418 trailing:10 slidingwindow:5:10 minlen:50 (Bolger *et al.*, 2014). We created repeat consensus

419 sequences for TE families (sequences are available on https://github.com/crolllab/datasets; Figure 1 -

420 figure supplement 5) in the complete reference genome IPO323 (Goodwin *et al.*, 2011) with

421 RepeatModeler version open-4.0.7 (http://www.repeatmasker.org/RepeatModeler/) based on the

422 RepBase Sequence Database and de novo (Bao *et al.*, 2015). TE classification into superfamilies and

423 families was based on an approach combining detection of conserved protein sequences and tools to

424 detect non-autonomous TEs (Badet *et al.*, 2020). To detect TE insertions, we used the R-based tool

425 ngs_te_mapper version 79ef861f1d52cdd08eb2d51f145223fad0b2363c integrated into the

426 McClintock pipeline version 20cb912497394fabddcdaa175402adacf5130bd1, using bwa version

427 0.7.4-r385 to map Illumina short reads, samtools version 0.1.19 to convert alignment file formats

428 and R version 3.2.3 (Li & Durbin, 2009; Li *et al.*, 2009; Linheiro & Bergman, 2012; R Core Team,

429 2017; Nelson *et al.*, 2017).

430

16

431 DOWN-SAMPLING ANALYSIS

432 We performed a down-sampling analysis to estimate the sensitivity of the TE detection with

433 ngs_te_mapper based on variation in read depth. We selected one isolate per population matching

434 the average coverage of the population. We extracted the per-base pair read depth with the

435 genomecov function of bedtools version 2.27.1 and calculated the genome-wide mean read depth

436 (Quinlan & Hall, 2010). The number of reads in the original fastq file was reduced in steps of 10%

437 to simulate the impact of reduced coverage. We analyzed each of the obtained reduced read subsets

438 with ngs_te_mapper using the same parameters as described above. The correlation between the

439 number of detected insertions and the read depth was visualized using the function nls with model

440 SSlogis in R and visualized with ggplot2 (Wickham, 2016). The number of detected TEs increased

441 with the number of reads until reaching a plateau indicating saturation (Figure 1B). Saturation was

442 reached at a coverage of approximately 15X, hence we retained only isolates with an average read

443 depth above 15X for further analyses. We thus excluded one isolate from the Oregon 2015

444 population and ten isolates from the Switzerland 2016 population.

445

446 VALIDATION PROCEDURE FOR PREDICTED TE INSERTIONS

447 ngs_te_mapper detects the presence but not the absence of a TE at any given locus. We devised

448 additional validation steps to ascertain both the presence as well as the absence of a TE across all

449 loci in all individuals. TEs absent in the reference genome were validated by re-analyzing mapped

450 Illumina reads. Reads spanning both parts of a TE sequence and an adjacent chromosomal sequence

451 should only map to the reference genome sequence and cover the target site duplication of the TE

452 (Figure 1C). We used bowtie2 version 2.3.0 with the parameter --very-sensitive-local to map

453 Illumina short reads of each isolate on the reference genome IPO323 (Langmead & Salzberg, 2012).

454 Mapped Illumina short reads were then sorted and indexed with samtools and the resulting bam file

455 was converted to a bed file with the function bamtobed in bedtools. We extracted all mapped reads

456 with an end point located within 100 bp of the target site duplication (Figure 1C). We tested whether

457 the number of reads with a mapped end around the target site duplication significantly deviated if the

458 mapping ended exactly at the boundary. A mapped read ending exactly at the target site duplication

17

459    boundary is indicative of a split read mapping to a TE sequence absent in the reference genome. To

460    test for the deviation in the number of read mappings around the target site duplication, we used a

461    Poisson distribution and the *ppois* function in R version 3.5.1 (Figure 1C). We identified a TE as

462    present in an isolate if tests on either side of the target site duplication had a *p*-value < 0.001 (Figure

463    5 - figure supplement 1; Figure 1 - figure supplement 1B and Figure 1 - figure supplement 2).

464

465    For TEs present in the reference genome, we analyzed evidence for spliced junction reads spanning

466    the region containing the TE. Spliced reads are indicative of a discontinuous sequence and, hence,

467    absence of the TE in a particular isolate (Figure 1D). We used STAR version 2.5.3a to detect spliced

468    junction reads with the following set of parameters: --runThreadN 1 --outFilterMultimapNmax 100 -

469    -winAnchorMultimapNmax 200 --outSAMmultNmax 100 --outSAMtype BAM Unsorted --

470    outFilterMismatchNmax 5 --alignIntronMin 150 --alignIntronMax 15000 (Dobin *et al.*, 2012). We

471    then sorted and indexed the resulting bam file with samtools and converted split junction reads with

472    the function bam2hints in bamtools version 2.5.1 (Barnett *et al.*, 2011). We selected loci without

473    overlapping spliced junction reads using the function intersect in bedtools with the parameter -loj -v.

474    We considered a TE as truly absent in an isolate if ngs_te_mapper did not detect a TE and evidence

475    for spliced junction reads were found, indicating that the isolate had no inserted TE in this region. If

476    the absence of a TE could not be confirmed by spliced junction reads, we labelled the genotype as

477    missing. Finally, we excluded TE loci with more than 20% missing data from further investigations

478    (Figure 1D and Figure 1 - figure supplement 1C).

479

480    CLUSTERING OF TE INSERTIONS INTO LOCI

481    We identified insertions across isolates as being the same locus if all detected TEs belonged to the

482    same TE family and insertion sites differed by ≤100 bp (Figure 1 - figure supplement 3). We used

483    the R package *GenomicRanges* version 1.28.6 with the functions makeGRangesFromDataFrame and

484    findOverlaps and the R package *devtools* version 1.13.4 (Lawrence *et al.*, 2013; Wickham & Chang,

485    2016). We used the R package *dplyr* version 0.7.4 to summarize datasets

486    (https://dplyr.tidyverse.org/). Population-specific frequencies of insertions were calculated with the

18

487 function allele.count in the R package *hierfstat* version 0.4.22 (Goudet, 2005). We conducted a

488 principal component analysis for TE insertion frequencies filtering for a minor allele frequency ≥

489 5%. We also performed a principal component analysis for genome-wide single nucleotide

490 polymorphism (SNP) data obtained from Hartmann et al (2017) and Singh et al (2020). As described

491 previously, SNPs were hard-filtered with VariantFiltration and SelectVariants tools integrated in the

492 Genome Analysis Toolkit (GATK) (McKenna *et al.*, 2010). SNPs were removed if any of the

493 following filter conditions applied: QUAL<250; QD<20.0; MQ<30.0; -2 > BaseQRankSum > 2; -2

494 > MQRankSum > 2; -2 > ReadPosRankSum > 2; FS>0.1. SNPs were excluded with vcftools version

495 0.1.17 and plink version 1.9 requiring a genotyping rate >90% and a minor allele frequency >5%

496 (https://www.cog-genomics.org/plink2, Chang et al., 2015). Finally, we converted tri-allelic SNPs to

497 bi-allelic SNPs by recoding the least frequent allele as a missing genotype. Principal component

498 analysis was performed using the *gdsfmt* and *SNPRelate* packages in R (Zheng *et al.*, 2012, 2017).

499 For a second principal component analysis with a reduced set of random markers, we randomly

500 selected SNPs with vcftools and the following set of parameters: --maf 0.05 –thin 200000 to obtain

501 an approximately equivalent number of SNPs as TE loci.

502

503 EVALUATION OF SINGLETON INSERTIONS

504 To evaluate the reliability of singleton TE insertion loci, we analyzed singleton loci in isolates for

505 which we had both Illumina datasets and complete reference-quality genomes (Badet *et al.*, 2020).

506 From a set of 19 long-read PacBio reference genomes spanning the global distribution of *Z. tritici,*

507 one isolate each from Australia, Israel, North America (1990) and four isolates from Europe (1999)

508 were also included in the TE insertion screening. To assess the reliability of singleton TE insertions,

509 we first investigated structural variation analyses among the reference genomes (Badet *et al.*, 2021,

510 Supplementary Data 1 and 2). The structural variation was called both based on split read mapping

511 of PacBio reads and pairwise whole-genome alignments. Using bedtools intersect, we recovered for

512 the 31 singleton TE loci in the 7 analyzed genomes a total of 17 loci showing either an indel,

513 translocation, copy number polymorphism, duplication, inverted duplication, inversion, or inverted

514    translocation at the same location. We visually inspected the PacBio read alignment bam files

515    against the IPO323 reference genome using IGV version 2.4.16 (Robinson *et al.*, 2011), and found a

516    typical coverage increase at the target site duplication, with most read mappings interrupted at the

517    target site duplication as expected for an inserted TE. For the 14 remaining TE loci, we extracted the

518    region of the predicted insertion and padded the sequence on both ends with an additional 500 bp

519    using samtools faidx. We used blast to identify a homologous region in the assembled reference-

520    quality genomes. Matching regions were inspected based on blastn for the presence of a TE

521    sequence matching the TE family originally detected at the locus. With this second approach, we

522    confirmed an additional five singletons to be true insertions. Both methods combined produced

523    supportive evidence for 22 out of 31 singleton insertions (71%). We calculated the read coverage

524    after mapping to the reference genome IPO323 with bedtools genomecov for each PacBio long-read

525    dataset and calculated mean coverage for 500 bp regions around singleton TE insertions.

526

527    POPULATION DIFFERENTIATION IN TE FREQUENCIES

528    We calculated Nei's fixation index ($F_{ST}$) between pairs of populations using the R packages *hierfstat*

529    and *adegenet* version 2.1.0 (Jombart, 2008; Jombart & Ahmed, 2011). To understand the

530    chromosomal context of TE insertion loci across isolates, we analyzed draft genome assemblies. We

531    generated *de novo* genome assemblies for all isolates using SPAdes version 3.5.0 with the parameter

532    --careful and a kmer range of "21, 29, 37, 45, 53, 61, 79, 87" (Bankevich *et al.*, 2012). We used

533    blastn to locate genes adjacent to TE insertion loci on genomic scaffolds of each isolate. We then

534    extracted scaffold sequences surrounding 10 kb up- and downstream of the localized gene with the

535    function faidx in samtools and reverse complemented the sequence if needed. Then, we performed

536    multiple sequence alignments for each locus across all isolates with MAFFT version 7.407 with

537    parameter --maxiterate 1000 (Katoh & Standley, 2013). We performed visual inspections to ensure

538    correct alignments across isolates using Jalview version 2.10.5 (Waterhouse *et al.*, 2009). To

539    generate phylogenetic trees of individual gene or TE loci, we extracted specific sections of the

540    alignment using the function extractalign in EMBOSS version 6.6.0 (Rice *et al.*, 2000) and

541    converted the multiple sequence alignment into PHYLIP format with jmodeltest version 2.1.10 using

20

542   the -getPhylip parameter. We then estimated maximum likelihood phylogenetic trees with the

543   software PhyML version 3.0, the K80 substitution model and 100 bootstraps on the ATGC South of

544   France bioinformatics platform (Guindon & Gascuel, 2003; Guindon *et al.*, 2010; Darriba *et al.*,

545   2012). Bifurcations with a supporting value lower than 10% were collapsed in TreeGraph version

546   2.15.0-887 beta and trees were visualized as circular phylograms in Dendroscope version 2.7.4

547   (Huson *et al.*, 2007; Stöver & Müller, 2010). For loci showing complex rearrangements, we

548   generated synteny plots using 19 completely sequenced genomes from the same species using the R

549   package *genoplotR* version 0.8.9 (Guy *et al.*, 2010; Badet *et al.*, 2020). We calculated the

550   population-specific allele frequency for TE loci and estimated the exponential decay curve with a

551   self-starting Nls asymptomatic regression model nls(p_loci ~ SSasymp(p_round, Asym, R0, lrc) in

552   R.

553   We analyzed signatures of selective sweeps based on genome-wide SNPs using the extended

554   haplotype homozygosity (EHH) tests implemented in the R package *REHH* (Sabeti *et al.*, 2007;

555   Gautier & Vitalis, 2012). We analyzed within-population signatures based on the iHS statistic and

556   chose a maximum gap distance of 20 kb. We also analyzed cross-population signatures based on the

557   XP-EHH statistic for the following two population pairs: North America 1990 versus North America

558   2015, Europe 1999 versus Europe 2016. We defined significant selective sweeps as being among the

559   99.9th percentile outliers of the iHS and XP-EHH statistics. Significant SNPs at less than 5 kb were

560   clustered into a single selective sweep region adding +/- 2.5 kb. Finally, we analyzed whether TE

561   loci in the population pairs were within 10 kb of a region identified as a selective sweep by XP-EHH

562   using the function intersect from bedtools.

563

564   GENOMIC LOCATION OF TE INSERTIONS

565   To characterize the genomic environment of TE insertion loci, we split the reference genome into

566   non-overlapping windows of 10 kb using the function splitter from EMBOSS. TEs were located in

567   the reference genome using RepeatMasker providing consensus sequences from RepeatModeler

568   (http://www.repeatmasker.org/). To analyze coding sequence, we retrieved the gene annotation for

21

569    the reference genome (Grandaubert *et al.*, 2015). We estimated the percentage covered by genes or

570    TEs per window using the function intersect in bedtools. Additionally, we calculated the GC content

571    using the tool get_gc_content (https://github.com/spundhir/RNA-

572    Seq/blob/master/get_gc_content.pl). We extracted the number of TEs present in 1 kb windows

573    around each annotated core gene in the reference genome IPO323, using the function window in

574    bedtools. We calculated the relative distances between each gene and the closest TE with the

575    function bedtools closest. For the TEs inserted into genes, we used the function intersect in bedtools

576    to distinguish intron and exon insertions with the parameters -wo and -v, respectively. TEs that

577    overlap more than one exon were only counted once. For each 100 bp segment in the 1 kb windows

578    as well as for introns and exons, we calculated the mean number of observed TE insertions per base

579    pair. We calculated the mean number of TEs per window and calculated the log2 of the observed

580    number of TE insertions divided by the expected value. We extracted information about

581    recombination hotspots from Croll *et al.* (2015). This dataset is based on two experimental crosses

582    initiated from isolates included in our analyses (1A5x1E4, 3D1x3D7). The recombination rates were

583    assessed based on the reference genome IPO323 and analyzed with the *R/qtl* package in R. We used

584    bedtools intersect to compare both TE density in IPO323 and TE insertion polymorphism with

585    predicted recombination hotspots.

586

587    CORE GENOME SIZE ESTIMATION

588    Accessory chromosomes show presence/absence variation within the species and length

589    polymorphism (Goodwin *et al.*, 2011; Croll *et al.*, 2013) and thus impact genome size. We

590    controlled for this effect by first mapping sequencing reads to the reference genome IPO323 using

591    bowtie2 with --very-sensitive-local settings and retained only reads mapping to any of the 13 core

592    chromosomes using seqtk subseq v1.3-r106 (https://github.com/lh3/seqtk/). Furthermore, we found

593    that different sequencing runs showed minor variation in the distribution of the per read GC content.

594    In particular, reads of a GC content lower than 30 % were underrepresented in the Australian (mean

595    reads < 30 % of the total readset: 0.05 %), North American 1990 (0.07 %) and Middle East (0.1 %)

596    populations, and higher in the Europe 1999 (1.3 %), North American 2015 (3.0 %) and Europe 2016

22

597 (4.02 %) populations (Figure 1 - figure supplement 4). Library preparation protocols and Illumina

598 sequencer generations are known factors influencing the recovery of reads of varying GC content

599 (Benjamini & Speed, 2012).

600

601 To control a potential bias stemming from this, we subsampled reads based on GC content to create

602 homogeneous datasets. For this, we first retrieved the mean GC content for each read pair using

603 geecee in EMBOSS and binned reads according to GC content. For the bins with a GC content

604 <30%, we calculated the mean proportion of reads from the genome over all samples. We then used

605 seqtk subseq to subsample reads of <30% to adjust the mean GC content among readsets. We

606 generated *de novo* genome assemblies using the SPAdes assembler version with the parameters --

607 careful and a kmer range of "21, 29, 37, 45, 53, 61, 79, 87". The SPAdes assembler is optimized for

608 the assembly of relatively small eukaryotic genomes. We evaluated the completeness of the

609 assemblies using BUSCO v4.1.1 with the fungi_odb10 gene test set (Simão *et al.*, 2015). We finally

610 ran Quast v5.0.2 to retrieve assembly metrics including scaffolds of at least 1 kb (Mikheenko *et al.*,

611 2018).

612

613 FUNGICIDE RESISTANCE ASSAY

614 To quantify susceptibility towards propiconazole we used a previously published microtiter plate

615 assay dataset with 3 replicates performed for each isolate and concentration. Optical density was

616 used to estimate growth rates under different fungicide concentrations (0, 0.00006, 0.00017, 0.0051,

617 0.0086, 0.015, 0.025, 0.042, 0.072, 0.20, 0.55, 1.5 mgL$^{-1}$) (Hartmann *et al.*, 2020). We calculated

618 dose-response curves and estimated the half-maximal lethal concentration EC$_{50}$ with a 4-parameter

619 logistics curve in the R package *drc* (Ritz & Streibig, 2005).

620

621 **Data availability**

23

622 Sequence data is deposited at the NCBI SRA under the accession numbers PRJNA327615,

623 PRJNA596434 and PRJNA178194. Transposable element consensus sequences are available from

624 https://github.com/crolllab/datasets.

625

626 **Author contributions**

627 UO and DC conceived the study, UO, TW and DC designed analyses, UO, TB, TV and FEH

628 performed analyses, FEH, NKS, LNA, PK, CCM and BAM provided samples/datasets, BAM and

629 DC provided funding, UO and DC wrote the manuscript with input from co-authors. All authors

630 reviewed the manuscript and agreed on submission.

631

632 **Acknowledgments**

638

639 **Competing interests**

640 We declare to have no competing interests.

641

642

643 **REFERENCES**

644 **Ali SS, Khan M, Mullins E, Doohan FM**. **2014**. Identification of Fusarium oxysporum Genes
645 Associated with Lignocellulose Bioconversion Competency. *Bioenergy Research* **7**: 110–119.

646 **Badet T, Fouché S, Hartmann FE, Zala M, Croll D**. **2021**. Machine-learning predicts genomic
647 determinants of meiosis-driven structural variation in a eukaryotic pathogen. *Nature*
648 *Communications* **12**.

649 **Badet T, Oggenfuss U, Abraham L, McDonald BA, Croll D**. **2020**. A 19-isolate reference-quality
650 global pangenome for the fungal wheat pathogen Zymoseptoria tritici. *BMC Biology* **18**: 12.

651 **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,**
652 **Nikolenko SI, Pham S, Prjibelski AD,** *et al.* **2012**. SPAdes: a new genome assembly algorithm and

653    its applications to single-cell sequencing. *Journal of computational biology: a journal of*
654    *computational molecular cell biology* **19**: 455–77.

655    **Bao W, Kojima KK, Kohany O**. **2015**. Repbase Update, a database of repetitive elements in
656    eukaryotic genomes. *Mobile DNA* **6**: 4–9.

657    **Barnett DW, Garrison EK, Quinlan AR, Strmberg MP, Marth GT**. **2011**. Bamtools: A C++
658    API and toolkit for analyzing and managing BAM files. *Bioinformatics* **27**: 1691–1692.

659    **Baucom RS, Estill JC, Leebens-Mack J, Bennetzen JL**. **2008**. Natural selection on gene function
660    drives the evolution of LTR retrotransposon families in the rice genome. *Genome Research* **19**: 243–
661    254.

662    **Benjamini Y, Speed TP**. **2012**. Summarizing and correcting the GC content bias in high-throughput
663    sequencing. *Nucleic Acids Research* **40**: 1–14.

664    **Blommaert J**. **2020**. Genome size evolution: towards new model systems for old questions.
665    *Proceedings. Biological sciences* **287**: 20201441.

666    **Bolger AM, Lohse M, Usadel B**. **2014**. Trimmomatic: a flexible trimmer for Illumina sequence
667    data. *Bioinformatics* **30**: 2114–2120.

668    **Bonifacino JS, Glick BS**. **2004**. The Mechanisms of Vesicle Budding and Fusion. *Cell* **116**: 153–
669    166.

670    **Bourque G, Burns KH, Gehring M, Gorbunova V, Seluanov A, Hammell M, Imbeault M,**
671    **Izsvák Z, Levin HL, Macfarlan TS, *et al.* 2018**. Ten things you should know about transposable
672    elements. *Genome Biology* **19**: 199.

673    **Charlesworth B, Charlesworth D**. **1983**. The population dynamics of transposable elements.
674    *Genetical Research* **42**: 1–27.

675    **Chuong EB, Elde NC, Feschotte C**. **2017**. Regulatory activities of transposable elements: from
676    conflicts to benefits. *Nature Reviews Genetics* **18**: 71–86.

677    **Cridland JM, Macdonald SJ, Long AD, Thornton KR**. **2013**. Abundance and distribution of
678    transposable elements in two drosophila QTL mapping resources. *Molecular Biology and Evolution*
679    **30**: 2311–2327.

680    **Croll D, Lendenmann MH, Stewart E, McDonald BA**. **2015**. The Impact of Recombination
681    Hotspots on Genome Evolution of a Fungal Plant Pathogen. *Genetics* **201**: 1213-U787.

682    **Croll D, Zala M, McDonald BA**. **2013**. Breakage-fusion-bridge Cycles and Large Insertions
683    Contribute to the Rapid Evolution of Accessory Chromosomes in a Fungal Pathogen (J Heitman,
684    Ed.). *PLOS Genetics* **9**: e1003567.

685    **Darriba D, Taboada GL, Doallo R, Posada D**. **2012**. jModelTest 2: more models, new heuristics
686    and parallel computing. *Nature Methods* **9**: 772.

687    **Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Gingeras TR, Batut P,**
688    **Chaisson M**. **2012**. STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15–21.

689    **Eichler EE, Sankoff D**. **2003**. Structural dynamics of eukaryotic chromosome evolution. *Science*
690    **301**: 793–797.

691    **Estep LK, Torriani SFF, Zala M, Anderson NP, Flowers MD, Mcdonald BA, Mundt CC,**
692    **Brunner PC**. **2015**. Emergence and early evolution of fungicide resistance in North American
693    populations of Zymoseptoria tritici. *Plant Pathology* **64**: 961–971.

694    **Feschotte C**. **2008**. Transposable elements and the evolution of regulatory networks. *Nature*
695    *Reviews Genetics* **9**: 397–405.

696    **Feurtey A, Lorrain C, Croll D, Eschenbrenner C, Freitag M, Habig M, Haueisen J, Möller M,**
697    **Schotanus K, Stukenbrock EH**. **2020**. Genome compartmentalization predates species divergence
698    in the plant pathogen genus Zymoseptoria. *BMC genomics* **21**: 588.

699    **Fouché S, Badet T, Oggenfuss U, Plissonneau C, Francisco CS, Croll D**. **2019**. Stress-driven
700    transposable element de-repression dynamics in a fungal pathogen. *Molecular Biology and*
701    *Evolution*.

702    **Frantzeskakis L, Kracher B, Kusch S, Yoshikawa-Maekawa M, Bauer S, Pedersen C, Spanu**

25

703 **PD, Maekawa T, Schulze-Lefert P, Panstruga R**. **2018**. Signatures of host specialization and a
704 recent transposable element burst in the dynamic one-speed genome of the fungal barley powdery
705 mildew pathogen. *BMC Genomics* **19**: 1–23.

706 **Fu Y, Kawabe A, Etcheverry M, Ito T, Toyoda A, Fujiyama A, Colot V, Tarutani Y, Kakutani**
707 **T**. **2013**. Mobilization of a plant transposon by expression of the transposon-encoded anti-silencing
708 factor. *EMBO Journal* **32**: 2407–2417.

709 **Gautier M, Vitalis R**. **2012**. Rehh An R package to detect footprints of selection in genome-wide
710 SNP data from haplotype structure. *Bioinformatics* **28**: 1176–1177.

711 **Gilly A, Etcheverry M, Madoui MA, Guy J, Quadrana L, Alberti A, Martin A, Heitkam T,**
712 **Engelen S, Labadie K,** *et al.* **2014**. TE-Tracker: systematic identification of transposition events
713 through whole-genome resequencing. *Bmc Bioinformatics* **15**.

714 **González-Sayer S, Oggenfuss U, García I, Aristizabal F**. High-quality genome assembly of
715 Pseudocercospora ulei the main threat to natural rubber trees. : 0–1.

716 **Goodwin SB, Ben M'Barek S, Dhillon B, Wittenberg AHJ, Crane CF, Hane JK, Foster AJ,**
717 **Van der Lee TAJ, Grimwood J, Aerts A,** *et al.* **2011**. Finished Genome of the Fungal Wheat
718 Pathogen Mycosphaerella graminicola Reveals Dispensome Structure, Chromosome Plasticity, and
719 Stealth Pathogenesis (HS Malik, Ed.). *PLOS Genetics* **7**: e1002070.

720 **Goudet J**. **2005**. Hierstat, a package for R to compute and test heirarchical F-statistics. *Molecular*
721 *Ecology Notes* **5**: 184–186.

722 **Grandaubert J, Bhattacharyya A, Stukenbrock EH**. **2015**. RNA-seq-Based Gene Annotation and
723 Comparative Genomics of Four Fungal Grass Pathogens in the Genus Zymoseptoria Identify Novel
724 Orphan Genes and Species-Specific Invasions of Transposable Elements. *G3-Genes Genomes*
725 *Genetics* **5**: 1323–1333.

726 **Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O**. **2010**. New
727 Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance
728 of PhyML 3.0. *Systematic Biology* **59**: 307–321.

729 **Guindon S, Gascuel O**. **2003**. A simple, fast, and accurate algorithm to estimate large phylogenies
730 by maximum likelihood. *Systematic Biology* **52**: 696–704.

731 **Guy L, Kultima JR, Andersson SGE**. **2010**. GenoPlotR: comparative gene and genome
732 visualization in R. *Bioinformatics* **26**: 2334–2335.

733 **Hartmann F, Croll D**. **2017**. Distinct Trajectories of Massive Recent Gene Gains and Losses in
734 Populations of a Microbial Eukaryotic Pathogen. *Molecular Biology and Evolution*.

735 **Hartmann F, McDonald M, Croll D**. **2018**. Genome-wide evidence for divergent selection
736 between populations of a major agricultural pathogen. *Molecular Ecology* **27**: 2725–2741.

737 **Hartmann FE, Sánchez-Vallet A, McDonald BA, Croll D**. **2017**. A fungal wheat pathogen
738 evolved host specialization by extensive chromosomal rearrangements. *The ISME Journal* **11**: 1189–
739 1204.

740 **Hartmann FE, Vonlanthen T, Singh NK, McDonald MC, Milgate A, Croll D**. **2020**. The
741 complex genomic basis of rapid convergent adaptation to pesticides across continents in a fungal
742 plant pathogen. *Molecular Ecology*.

743 **Hollister JD, Gaut BS**. **2009**. Epigenetic silencing of transposable elements: A trade-off between
744 reduced transposition and deleterious effects on neighboring gene expression. *Genome Research* **19**:
745 1419–1428.

746 **Huson DH, Richter DC, Rausch C, Dezulian T, Franz M, Rupp R**. **2007**. Dendroscope: An
747 interactive viewer for large phylogenetic trees. *BMC Bioinformatics* **8**: 1–6.

748 **Jiao W-B, Schneeberger K**. **2019**. Chromosome-level assemblies of multiple Arabidopsis thaliana
749 accessions reveal hotspots of genomic rearrangements. *bioRxiv*: 738880.

750 **Jombart T**. **2008**. Adegenet: A R package for the multivariate analysis of genetic markers.
751 *Bioinformatics* **24**: 1403–1405.

752 **Jombart T, Ahmed I**. **2011**. adegenet 1.3-1: New tools for the analysis of genome-wide SNP data.

26

753     *Bioinformatics* **27**: 3070–3071.

754     **Katoh K, Standley DM**. **2013**. MAFFT multiple sequence alignment software version 7:
755     Improvements in performance and usability. *Molecular Biology and Evolution* **30**: 772–780.

756     **Kidwell MG**. **2002**. Transposable elements and the evolution of genome size in eukaryotes.
757     *Genetica* **115**: 49–63.

758     **Krishnan P, Meile L, Plissonneau C, Ma X, Hartmann FE, Croll D, McDonald BA, Sánchez-**
759     **Vallet A**. **2018**. Transposable element insertions shape gene regulation and melanin production in a
760     fungal pathogen of wheat. *BMC Biology* **16**: 1–18.

761     **Lai X, Schnable JC, Liao Z, Xu J, Zhang G, Li C, Hu E, Rong T, Xu Y, Lu Y**. **2017**. Genome-
762     wide characterization of non-reference transposable element insertion polymorphisms reveals
763     genetic diversity in tropical and temperate maize. *BMC Genomics* **18**: 1–13.

764     **Langmead B, Salzberg SL**. **2012**. Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**:
765     357–359.

766     **Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, Morgan MT, Carey**
767     **VJ**. **2013**. Software for Computing and Annotating Genomic Ranges (A Prlic, Ed.). *PLOS*
768     *Computational Biology* **9**: e1003118.

769     **Li H, Durbin R**. **2009**. Fast and accurate short read alignment with Burrows-Wheeler transform.
770     *Bioinformatics* **25**: 1754–1760.

771     **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R**.
772     **2009**. The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079.

773     **Liang X, Wang B, Dong Q, Li L, Rollins JA, Zhang R, Sun G**. **2018**. Pathogenic adaptations of
774     Colletotrichum fungi revealed by genome wide gene family evolutionary analyses. *PLoS ONE* **13**:
775     1–25.

776     **Lim JK**. **1988**. Intrachromosomal rearrangements mediated by hobo transposons in Drosophila
777     melanogaster. *PNAS* **85**: 9153–9157.

778     **Linde CC, Zhan J, McDonald BA**. **2002**. Population Structure of *Mycosphaerella graminicola*:
779     From Lesions to Continents. *Phytopathology* **92**: 946–955.

780     **Linheiro RS, Bergman CM**. **2012**. Whole Genome Resequencing Reveals Natural Target Site
781     Preferences of Transposable Elements in Drosophila melanogaster (JE Stajich, Ed.). *PLOS ONE* **7**:
782     e30008.

783     **Long Q, Rabanal FA, Meng D, Huber CD, Farlow A, Platzer A, Zhang Q, Vilhjálmsson BJ,**
784     **Korte A, Nizhynska V, *et al.* 2013**. Massive genomic variation and strong selection in Arabidopsis
785     thaliana lines from Sweden. *Nature Genetics* **45**: 884–890.

786     **Lu L, Chen J, Robb SMC, Okumoto Y, Stajich JE, Wessler SR**. **2017**. Tracking the genome-
787     wide outcomes of a transposable element burst over decades of amplification. *Proceedings of the*
788     *National Academy of Sciences*: 201716459.

789     **Luo S, Zhang H, Duan Y, Yao X, Clark AG, Lu J**. **2020**. The evolutionary arms race between
790     transposable elements and piRNAs in Drosophila melanogaster. *BMC Evolutionary Biology* **20**: 14.

791     **Lynch M**. **2007**. *The Origins of Genome Architecture*. Sunderland MA: Sinauer Associates.

792     **McDonald BA, Mundt CC, Chen R**. **1996**. The role of selection on the genetic structure of
793     pathogen populations: Evidence from field experiments with Mycosphaerella graminicola on
794     wheat. *Euphytica* **92**: 73–80.

795     **McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,**
796     **Altshuler D, Gabriel S, Daly M, *et al.* 2010**. The Genome Analysis Toolkit: A MapReduce
797     framework for analyzing next-generation DNA sequencing data. *Genome Research* **20**: 1297–1303.

798     **Meile L, Croll D, Brunner PC, Plissonneau C, Hartmann FE, McDonald BA, Sánchez-Vallet**
799     **A**. **2018**. A fungal avirulence factor encoded in a highly plastic genomic region triggers partial
800     resistance to septoria tritici blotch. *New Phytologist* **219**: 1048–1061.

801     **Miga KH, Koren S, Rhie A, Vollger MR, Gershman A, Bzikadze A, Brooks S, Howe E,**
802     **Porubsky D, Logsdon GA, *et al.* 2020**. Telomere-to-telomere assembly of a complete human X

803     chromosome. *Nature* **585**: 79–84.

804     **Mikheenko A, Prjibelski A, Saveliev V, Antipov D, Gurevich A**. **2018**. Versatile genome
805     assembly evaluation with QUAST-LG. *Bioinformatics* **34**: i142–i150.

806     **Miyao A, Tanaka K, Murata K, Sawaki H, Takeda S, Abe K, Shinozuka Y, Onosato K,**
807     **Hirochika H**. **2003**. Target site specificity of the Tos17 retrotransposon shows a preference for
808     insertion within genes and against insertion in retrotransposon-rich regions of the genome. *Plant*
809     *Cell* **15**: 1771–1780.

810     **Nagarajan N, Pop M**. **2013**. Sequence assembly demystified. *Nature Reviews Genetics* **14**: 157–
811     167.

812     **Nakamura K, Oshima T, Morimoto T, Ikeda S, Yoshikawa H, Shiwa Y, Ishikawa S, Linak**
813     **MC, Hirai A, Takahashi H, *et al.* 2011**. Sequence-specific error profile of Illumina sequencers.
814     *Nucleic Acids Research* **39**.

815     **Nelson MG, Linheiro RS, Bergman CM**. **2017**. McClintock: An Integrated Pipeline for Detecting
816     Transposable Element Insertions in Whole-Genome Shotgun Sequencing Data. *G3&amp;#58;*
817     *Genes|Genomes|Genetics* **7**: 2763–2778.

818     **Oliver KR, McComb JA, Greene WK**. **2013**. Transposable elements: Powerful contributors to
819     angiosperm evolution and diversity. *Genome Biology and Evolution* **5**: 1886–1901.

820     **Omrane S, Audéon C, Ignace A, Duplaix C, Aouini L, Kema G, Walker A-S, Fillinger S**. **2017**.
821     Plasticity of the MFS1 promoter leads to multi drug resistance in the wheat pathogen Zymoseptoria
822     tritici. *mSphere*: 1–42.

823     **Omrane S, Sghyer H, Audeon C, Lanen C, Duplaix C, Walker AS, Fillinger S**. **2015**. Fungicide
824     efflux and the MgMFS1 transporter contribute to the multidrug resistance phenotype in
825     Zymoseptoria tritici field isolates. *Environmental Microbiology* **17**: 2805–2823.

826     **Peter M, Kohler A, Ohm RA, Kuo A, Krützmann J, Morin E, Arend M, Barry KW, Binder M,**
827     **Choi C, *et al.* 2016**. Ectomycorrhizal ecology is imprinted in the genome of the dominant symbiotic
828     fungus Cenococcum geophilum. *Nature Communications* **7**: 1–15.

829     **Petrov DA, Aminetzach YT, Davis JC, Bensasson D, Hirsh AE**. **2003**. Size matters: Non-LTR
830     retrotransposable elements and ectopic recombination in Drosophila. *Molecular Biology and*
831     *Evolution* **20**: 880–892.

832     **Piegu B, Guyot R, Picault N, Roulin A, Sanyal A, Kim H, Collura K, Brar DS, Jackson S,**
833     **Wing RA, *et al.* 2006**. Doubling genome size without polyploidization: Dynamics of
834     retrotransposition-driven genomic expansions in Oryza australiensis, a wild relative of rice. *Genome*
835     *Research* **21**: 1201.

836     **Plissonneau C, Stürchler A, Croll D**. **2016**. The Evolution of Orphan Regions in Genomes of a
837     Fungal Pathogen of Wheat. *mBio* **7**: 1–13.

838     **Quadrana L, Silveira AB, Mayhew GF, LeBlanc C, Martienssen RA, Jeddeloh JA, Colot V**.
839     **2016**. The Arabidopsis thaliana mobilome and its impact at the species level. *Elife* **5**.

840     **Quinlan AR, Hall IM**. **2010**. BEDTools: A flexible suite of utilities for comparing genomic
841     features. *Bioinformatics* **26**: 841–842.

842     **R Core Team**. **2017**. R: A language and environment for statistical computing. R Foundation for
843     Statistical Computing, Vienna, Austria.

844     **Raffaele S, Kamoun S**. **2012**. Genome evolution in filamentous plant pathogens: why bigger can be
845     better. *Nature Reviews Microbiology* **10**: 417–430.

846     **Rice P, Longden L, Bleasby A**. **2000**. EMBOSS: The European Molecular Biology Open Software
847     Suite. *Trends in Genetics* **16**: 276–277.

848     **Ritz C, Streibig JC**. **2005**. Bioassay analysis using R. *Journal of Statistical Software* **12**: 1–22.

849     **Robinson JT, Thorvaldsdóttir H, Winckler W, Guttman M, Lander ES, Getz G, Mesirov JP**.
850     **2011**. Integrative Genome Viewer. *Nature Biotechnology* **29**: 24–6.

851     **Rogers PD, Barker KS**. **2003**. Genome-wide expression profile analysis reveals coordinately
852     regulated genes associated with stepwise acquisition of azole resistance in Candida albicans clinical

853    isolates. *Antimicrobial Agents and Chemotherapy* **47**: 1220–1227.

854    **Rogers J, Gibbs RA**. **2014**. Content and Dynamics. *Nature Reviews Genetics* **15**: 347–359.

855    **Rouxel T, Grandaubert J, Hane JK, Hoede C, van de Wouw AP, Couloux A, Dominguez V,**
856    **Anthouard V, Bally P, Bourras S,** *et al.* **2011**. Effector diversification within compartments of the
857    Leptosphaeria maculans genome affected by Repeat-Induced Point mutations. *Nature*
858    *communications* **2**: 202.

859    **Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH,**
860    **McCarroll SA, Gaudet R,** *et al.* **2007**. Genome-wide detection and characterization of positive
861    selection in human populations. *Nature* **449**: 913–918.

862    **SanMiguel P, Gaut BS, Tikhonov A, Nakajima Y, Bennetzen JL**. **1998**. The paleontology of
863    intergene retrotransposons of maize. *Nature Genetics* **20**: 43–45.

864    **Shen RM, Batzer MA, Deininger PL**. **1991**. Evolution of the master Alu gene(s). *Journal of*
865    *Molecular Evolution* **33**: 311–320.

866    **Simão FA, Waterhouse RM, Ioannidis P, Kriventseva E V., Zdobnov EM**. **2015**. BUSCO:
867    Assessing genome assembly and annotation completeness with single-copy orthologs.
868    *Bioinformatics* **31**: 3210–3212.

869    **Singh NK, Chanclud E, Croll D**. **2020**. Population-level deep sequencing reveals the interplay of
870    clonal and sexual reproduction in the fungal wheat pathogen Zymoseptoria tritici.

871    **Sirisattha S, Momose Y, Kitagawa E, Iwahashi H**. **2004**. Toxicity of anionic detergents
872    determined by Saccharomyces cerevisiae microarray analysis. *Water Research* **38**: 61–70.

873    **Slotkin RK, Martienssen R**. **2007**. Transposable elements and the epigenetic regulation of the
874    genome. *Nature Reviews Genetics* **8**: 272–285.

875    **Soustre I, Letourneux Y, Karst F**. **1996**. Characterization of the Saccharomyces cerevisiae RTA1
876    gene involved in 7-aminocholesterol resistance. *Current Genetics* **30**: 121–125.

877    **Stöver BC, Müller KF**. **2010**. TreeGraph 2: Combining and visualizing evidence from different
878    phylogenetic analyses. *BMC Bioinformatics* **11**: 1–9.

879    **Stritt C, Gordon SP, Wicker T, Vogel JP, Roulin AC**. **2017**. Recent activity in expanding
880    populations and purifying selection have shaped transposable element landscapes across natural
881    accessions of the Mediterranean grass Brachypodium distachyon. *Genome Biology and Evolution*
882    **10**: 1–38.

883    **Stuart T, Eichten SR, Cahn J, Karpievitch Y V, Borevitz JO, Lister R**. **2016**. Population scale
884    mapping of transposable element diversity reveals links to gene regulation and epigenomic variation.
885    *eLife* **5**: 1–27.

886    **Stukenbrock EH, Banke S, Javan-Nikkhah M, McDonald BA**. **2007**. Origin and domestication of
887    the fungal wheat pathogen Mycosphaerella graminicola via sympatric speciation. *Molecular Biology*
888    *and Evolution* **24**: 398–411.

889    **Stukenbrock EH, Bataillon T, Dutheil JY, Hansen TT, Li RQ, Zala M, McDonald BA, Wang J,**
890    **Schierup MH**. **2011**. The making of a new pathogen: Insights from comparative population
891    genomics of the domesticated wheat pathogen Mycosphaerella graminicola and its wild sister
892    species. *Genome Research* **21**: 2157–2166.

893    **Torriani SFF, Melichar JPE, Mills C, Pain N, Sierotzki H, Courbot M**. **2015**. Zymoseptoria
894    tritici: A major threat to wheat production, integrated approaches to control. *Fungal Genetics and*
895    *Biology* **79**: 8–12.

896    **Walser J-C, Chen B, Feder ME**. **2006**. Heat-Shock Promoters: Targets for Evolution by P
897    Transposable Elements in Drosophila. *PLOS Genetics* **2**: e165.

898    **Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ**. **2009**. Jalview Version 2-A
899    multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**: 1189–1191.

900    **Wickham H**. **2016**. *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer-Verlag.

901    **Wickham H, Chang W**. **2016**. devtools: Tools to Make Developing R Packages Easier.

902    **Wong WY, Simakov O, Bridge DM, Cartwright P, Bellantuono AJ, Kuhn A, Holstein TW,**

903  **David CN, Steele RE, Martínez DE**. **2019**. Expansion of a single transposable element family is
904  associated with genome-size increase and radiation in the genus Hydra. *Proceedings of the National*
905  *Academy of Sciences* **116**: 22915–22917.

906  **Yew SM, Chan CL, Kuan CS, Toh YF, Ngeow YF, Na SL, Lee KW, Hoh CC, Yee WY, Ng KP**.
907  **2016**. The genome of newly classified Ochroconis mirabilis: Insights into fungal adaptation to
908  different living conditions. *BMC Genomics* **17**: 1–17.

909  **Zhan J, Kema GHJ, Waalwijk C, McDonald BA**. **2002**. Distribution of mating type alleles in the
910  wheat pathogen Mycosphaerella graminicola over spatial scales from lesions to continents. *Fungal*
911  *Genetics and Biology* **36**: 128–136.

912  **Zhan J, Linde CC, Jurgens T, Merz U, Steinebrunner F, McDonald BA**. **2005**. Variation for
913  neutral markers is correlated with variation for quantitative traits in the plant pathogenic fungus
914  Mycosphaerella graminicola. *Mol Ecol* **14**: 2683–2693.

915  **Zhan J, Pettway RE, McDonald BA**. **2003**. The global genetic structure of the wheat pathogen
916  Mycosphaerella graminicola is characterized by high nuclear diversity, low mitochondrial diversity,
917  regular recombination, and gene flow. *Fungal Genetics and Biology* **38**: 286–297.

918  **Zheng X, Gogarten SM, Lawrence M, Stilp A, Conomos MP, Weir BS, Laurie C, Levine D**.
919  **2017**. SeqArray-a storage-efficient high-performance data format for WGS variant calls.
920  *Bioinformatics* **33**: 2251–2257.

921  **Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS**. **2012**. A high-performance
922  computing toolset for relatedness and principal component analysis of SNP data. *Bioinformatics* **28**:
923  3326–3328.

924

925

926     **Figure Legends**

927

928     **Figure 1: Robust discovery and validation of transposable element (TE) insertions**: (A) General

929     analysis pipeline. (B) Read depth down-sampling analysis for one isolate per population with an

930     average coverage of the population. The vertical black line indicates the coverage at which on average

931     90% of the maximally detectable variants were recovered. Dashed black lines indicate the standard

932     error. The threshold for a minimal mean coverage was set at 15X (red line). (C) Validation of

933     insertions absent in the reference genome. (i) TE insertions that are not present in the reference

934     genome show a duplication of the target site and the part of the reads that covers the TE will not be

935     mapped against the reference genome. We thus expect reads to map to the TE surrounding region and

936     the target site duplication but not the TE itself. At the target site, a local duplication of read depth is

937     expected. (ii) We selected all reads in an interval of 100 bp up- and downstream including the target

938     site duplication to detect deviations in the number of reads terminating near the target site duplication.

939     (D) Validation of insertions present in the reference genome. (i) Analyses read coverage at target site

940     duplications. (ii) Decision map if a TE should be kept as a true insertion or rejected as a false positive.

941     Only predicted TE insertions that overlap evidence of split reads were kept as TE insertions in

942     downstream analyses. (E) Singleton validation using long-read PacBio sequencing. (i) Analysis if TE

943     insertions overlap with a detected insertion/deletion locus (Badet *et al*, 2021). (ii) Homology search of

944     the TE insertion flanking sequences based on the reference genome against PacBio reads. In addition,

945     the consensus sequence of the inserted TE was used for matches between the flanks.

946     **Figure supplement 1.** Validation of transposable element (TE) insertion predictions. (A) TEs not

947     present in the reference genome: distribution of additional TE hits found per locus after the outlier

948     test. Color indicates superfamilies. (B) TEs not present in the reference genome: distribution of

949     additional TE hits found per population after the outlier test. Colors indicate populations. (C) TEs

950     present in the reference genome: distribution of missing data per locus after the validation with

951     spliced junction reads. Missing data indicates that the TE was not predicted with ngs_te_mapper and

952     that there was no indication of spliced reads. The red line (=20 %) indicates the threshold for missing

953     data. TE loci with an amount of missing data > 20 % were completely excluded from further analyses.

954     Color indicates superfamily. (D) TEs present in the reference genome: detection of strong outlier

955     isolates with a high number of split reads. Color indicates the population.

956     **Figure supplement 2.** TE insertion validations for non-reference copies (Table).

957     **Figure supplement 3.** Establishment of transposable element (TE) loci with differing start and end

958     positions in the isolates. Distribution of length of distance for start position, end position and both

959     start and end combined after the correction.

960     **Figure supplement 4.** Bias for reads with a GC content lower than 30 % per population. Red lines

961     indicate the mean.

962     **Figure supplement 5.** TE consensus sequences (Table).

963

964 **Figure 2: Transposable element (TE) landscape across populations**. (A) Allele frequencies of the

965 TE insertions across all isolates. (B) TE insertions per Mb on core chromosomes (dark) and accessory

966 chromosomes (light). Dashed lines represent mean values. Blue: global mean of 75.65 insertions/Mb,

967 dark: core chromosome mean of 58 TEs/Mb, light: accessory chromosome mean of 102.24

968 insertions/Mb). (C) Number of TE insertions per family. (D) TE frequencies among isolates and copy

969 numbers across the genome. The blue line indicates the maximum number of isolates ($n = 284$). (E)

970 Allele frequency distribution of TE insertions into introns and exons. (F) Number of TE insertions

971 within 1 kb up- and downstream of genes on core chromosomes including introns and exons (100 bp

972 windows). The blue arrow indicates a gene schematic with exons and an intron, the green triangles

973 indicate TE insertions. The dotted blue line indicates no deviation from the expected value (*i.e.,* mean

974 number of TEs per window).

975 **Figure supplement 1.** TEs in reference (Table).

976 **Figure supplement 2.** Validation of singleton insertions detected by mapped Illumina reads using

977 PacBio read alignments for confirmation. (A) Comparison of TE family copy numbers per isolate to

978 the number of copies found in the reference genome (IPO323). The color is indicating superfamilies.

979 This figure includes only TE families that were detected in any of the isolates used for validation. (B)

980 Confirmation of singleton TE insertions detected in the isolates CH99_SW5, CH99_SW39,

981 CH99_3D7, CH99_3D1, ISR92_Ar_4f, AUS01_1H8 and ORE90Ste_4A10 using aligned PacBio

982 reads. Confirmed/not confirmed TE insertions are shown by TE family. (C) PacBio read coverage (in

983 500 bp window) at singleton loci.

984 **Figure supplement 3.** Presence absence matrix TE loci (Table).

985 **Figure supplement 4.** Singletons (Table).

986 **Figure supplement 5.** TE insertion loci characteristics. (A) Number of TE insertions and density

987 (insertions per Mb) in accessory and core genes. (B) Allele frequencies of TEs genome-wide and

988 restricted to recombination hotspots. (C) TE insertion density and TE copy numbers within and

989 outside of recombination hotspots.

990 **Figure supplement 6.** Hierarchy superfamilies. (A) Number of transposable element (TE) insertions

991 per superfamily. Colors indicate the superfamily. (B) Number of TE loci and classification hierarchy.

992 (C) Comparison of mean genome sequencing coverage and the number of detected TEs with

993 ngs_te_mapper in isolates of the Middle East population. Dots indicate the coverage and colors

994 indicate the superfamily.

995

996 **Figure 3: Differentiation in transposable element insertions frequencies across the genome**. (A)

997 Global pairwise $F_{ST}$ distributions shown across the 21 chromosomes. The red horizontal line indicates

998 the mean $F_{ST}$ (= 0.0163). TEs with a strong local short-term frequency difference among populations

999 are highlighted (blue: increase in Europe; green: increase in North America). (B) Allele frequency

1000    changes between the populations. The same TE loci as in panel A are highlighted. (C) Circos plot

1001    describing from the outside to the inside: The black line indicates chromosomal position in Mb. Blue

1002    bars indicate the gene density in windows of 100 kb with darker blue representing higher gene

1003    density. Red bars indicate the TE density in windows of 100 kb with a darker red representing higher

1004    TE density. Green triangles indicate positions of TE insertions with among population $F_{ST}$ value

1005    shown on the y-axis.

1006    **Figure supplement 1.** Global pairwise FST distributions shown separately for the 21 chromosomes.

1007    The red horizontal line indicates the mean FST = 0.0163. Colors are according to the three main

1008    superfamilies (RLG, RLC, DHH).

1009

1010    **Figure 4: Candidate adaptive transposable element (TE) insertions**. (A) Distribution of all

1011    extremely differentiated TEs and their distance to the closest gene. Color indicates the superfamily.

1012    The stars indicate TE insertions not found in the reference genome. (B) Location of the RLG_Luna

1013    TE insertion on chromosome 12 corresponding to its two closest genes. (C) Resistance against azole

1014    fungicides among isolates as a function of TE presence or absence. (D) Genomic niche of the

1015    RLG_Luna TE insertion on chromosome 12: $F_{ST}$ values for each TE insertion, gene content (blue), TE

1016    content (green) and GC content (yellow). The grey section highlights the insertion site. (E) Number of

1017    RLG_Luna copies per isolate and population. (F) Frequency changes of RLG_Luna between the two

1018    North American populations compared to the other populations. Colors indicate the number of copies

1019    per chromosome. (G) Phylogenetic trees of the coding sequences of either the gene encoding the

1020    RTA1-like protein or the protein kinase domain. Isolates of the two North American populations and

1021    an additional 11 isolates from other populations not carrying the insertion are shown. Blue color

1022    indicates TE presence, yellow indicates TE absence.

1023    **Figure supplement 1.** Top loci information (Table).

1024    **Figure supplement 2.** Additional top loci. Six additional candidate adaptive transposable element

1025    (TE) insertions. Each row corresponds to a candidate, with the first five being candidates detected in

1026    the North American populations and the last one in the European populations. For each candidate, the

1027    direction of the TE and the direction, function and distance of the closest two genes are indicated. The

1028    middle column indicates the location of the TE in the genomic niche, with TE content, gene content

1029    and GC content for the surrounding windows. The third column indicates resistance levels towards

1030    azole antifungals for isolates with and without the TE insertion.

1031

1032    **Figure 5: Population differentiation at transposable element (TE) and genome-wide SNP loci.**

1033    (A) Sampling locations of the six populations. Middle East represents the region of origin of the

1034    pathogen. In North America, the two populations were collected at an interval of 25 years in the same

1035    field in Oregon. In Europe, two populations were collected at an interval of 17 years from two fields

1036    in Switzerland <20 km apart. Dark arrows indicate the historic colonization routes of the pathogen.

33

1037   (B) Principal component analysis (PCA) of 284 *Zymoseptoria tritici* isolates, based on 900,193

1038   genome-wide SNPs. (C) PCA of a reduced SNP data set with randomly selected 203 SNPs matching

1039   approximately the number of analyzed TE loci. (D) PCA based on 193 TE insertion loci. Loci with

1040   allele frequency < 5% are excluded.

1041   **Figure supplement 1.** Isolates (Table).

1042

1043   **Figure 6: Global population structure of transposable element (TE) insertion polymorphism**.

1044   (A) Total TE copies per isolate. Colors identify TE superfamilies. (B) TE copies per family and (C)

1045   superfamily. (D) TE insertion frequency spectrum per population. The curve fitting was performed

1046   with a self-starting Nls asymptomatic regression model (E). TE family copy numbers per isolate.

1047   **Figure supplement 1.** Population changes additional. Variation in transposable element (TE) content

1048   per isolate across populations. (A) Total TE copies per superfamily (colored) and per isolate only

1049   including LTR (long terminal repeat) TEs *Copia* and *Gypsy*. Color indicates the family. (B) Total TE

1050   copies per superfamily (colored) and per isolate only on the core chromosomes. (C) Total TE copies

1051   per superfamily (colored) and per isolate only on the accessory chromosomes.

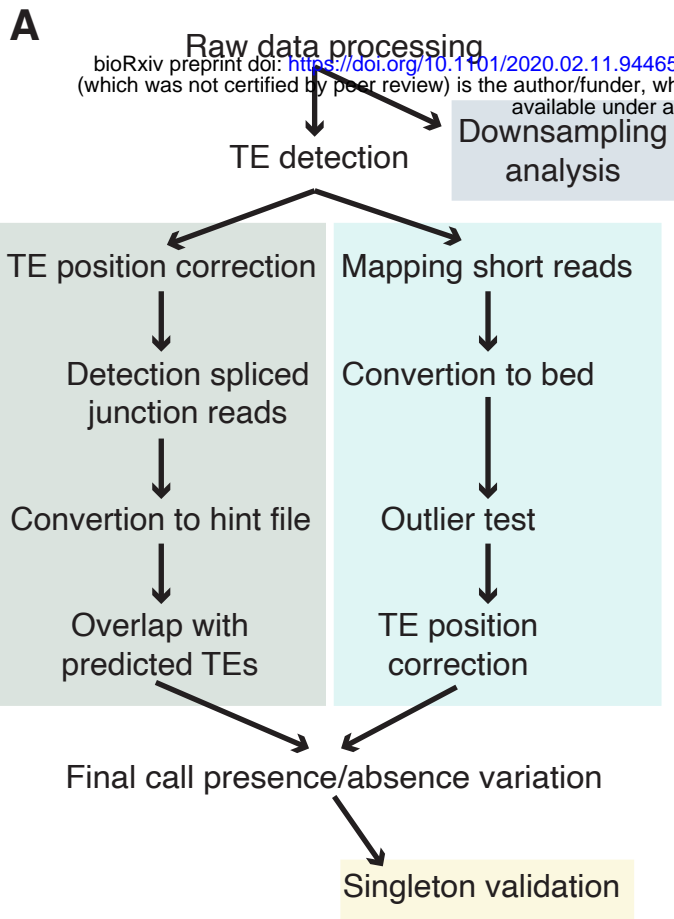1052   **Figure supplement 2.** Kolmogorof-Smirnov (Table).

1053   **Figure supplement 3.** Heatmap loci. (A) Presence (blue) and absence (yellow) matrix for all

1054   transposable element (TE) loci in all isolates per population. Colors on the left side indicate the

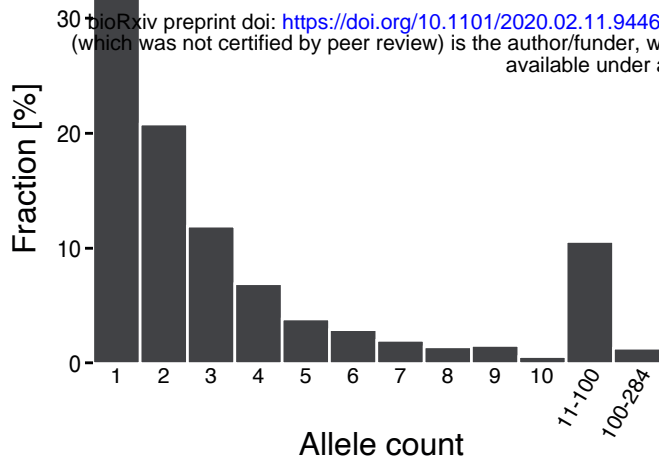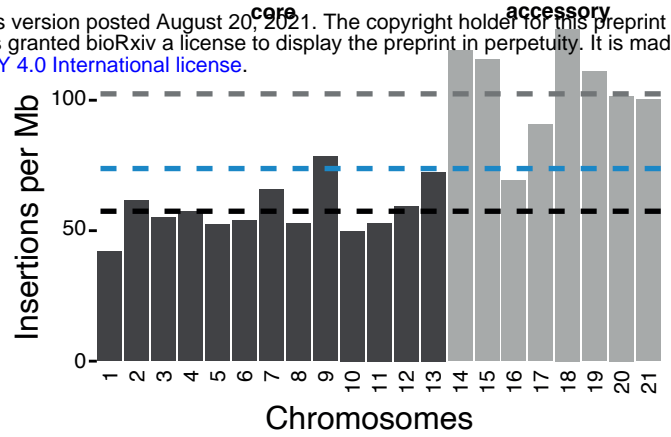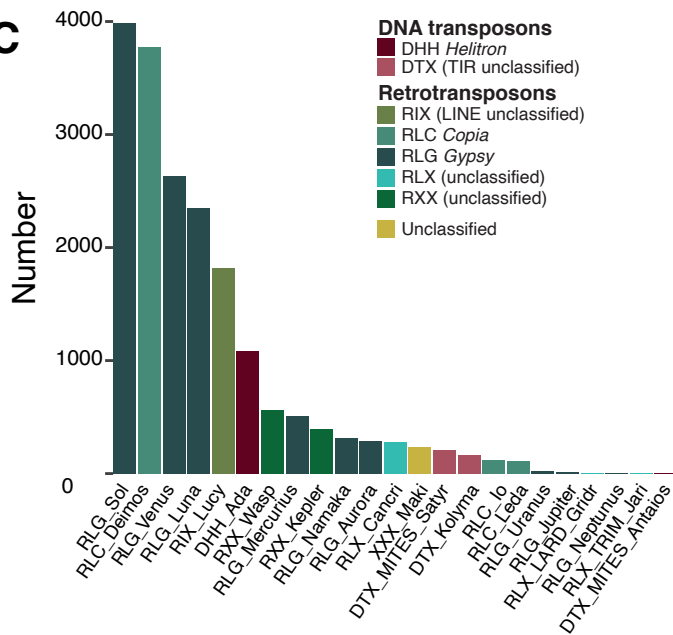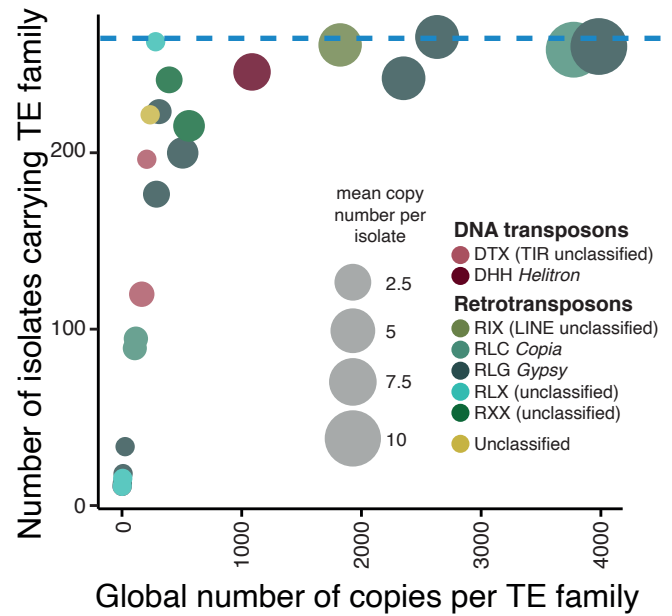1055   superfamily. (B) Comparison of different genomic regions with and without TE insertions in IPO323.
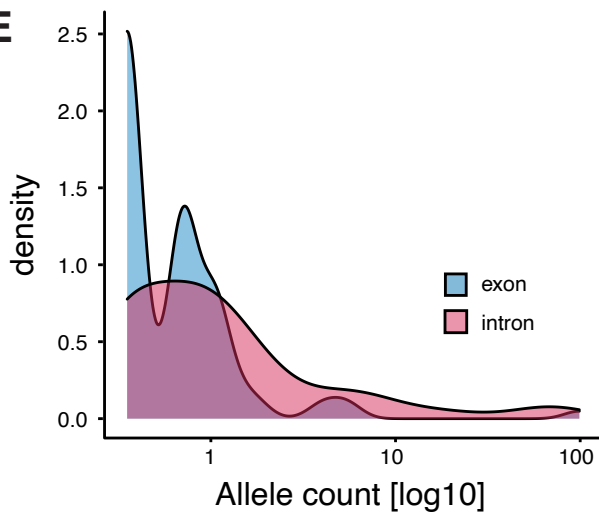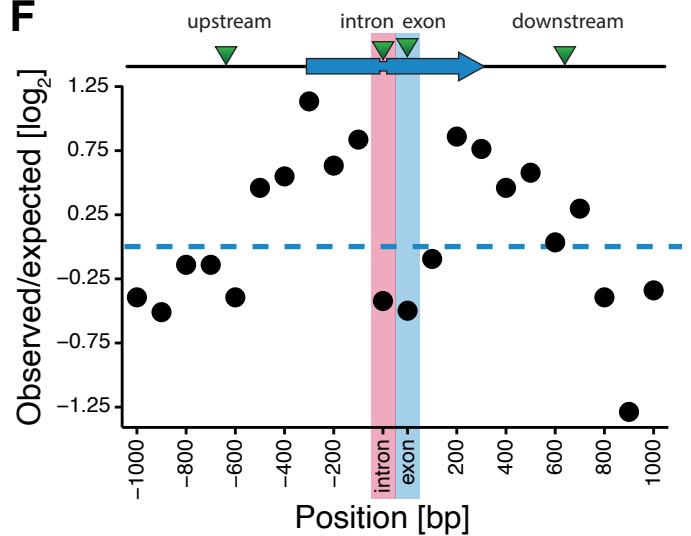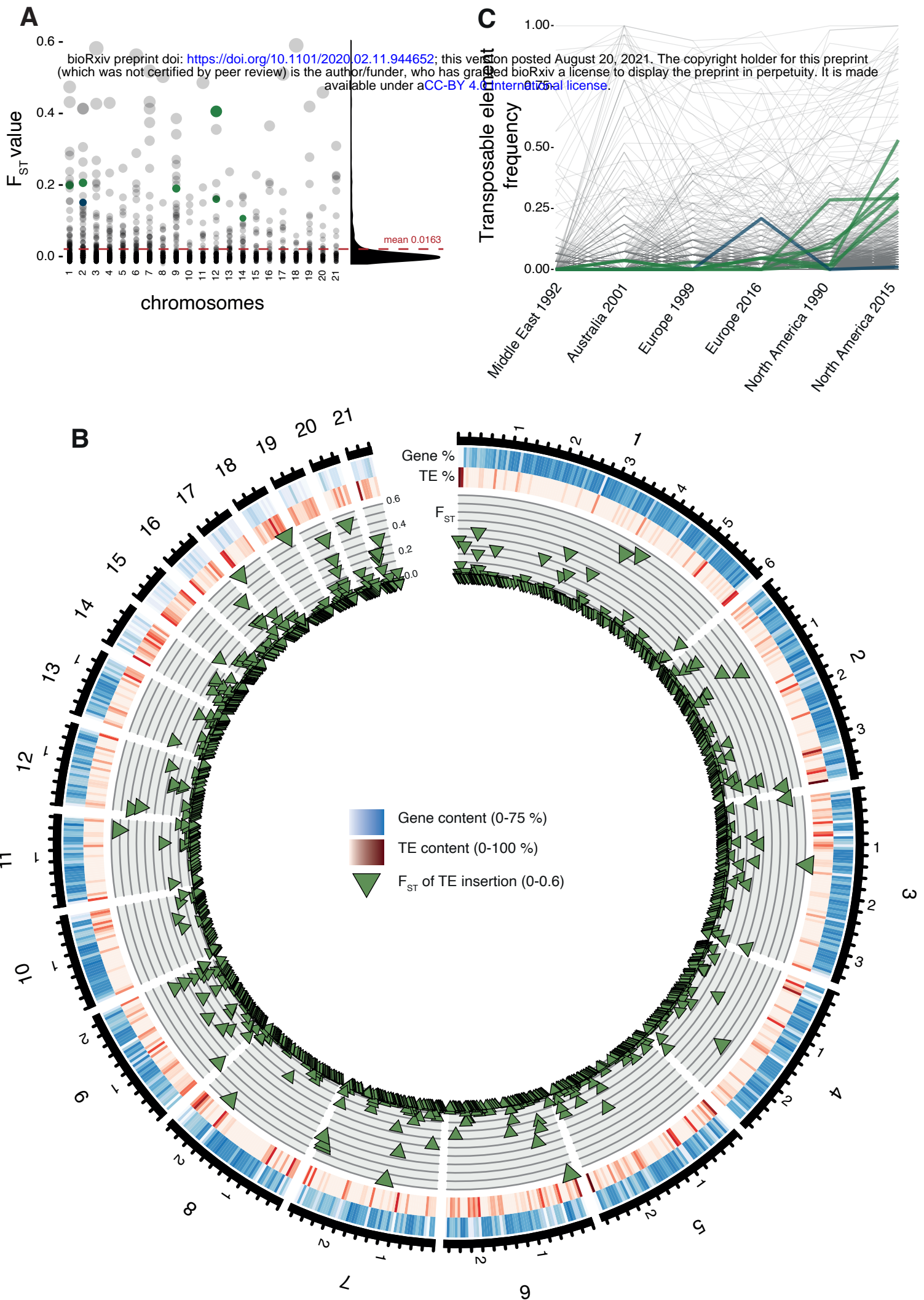
1056

1057   **Figure 7: Core genome size and transposable element (TE) evolution across populations**. (A)

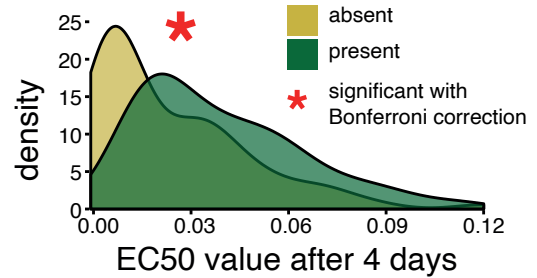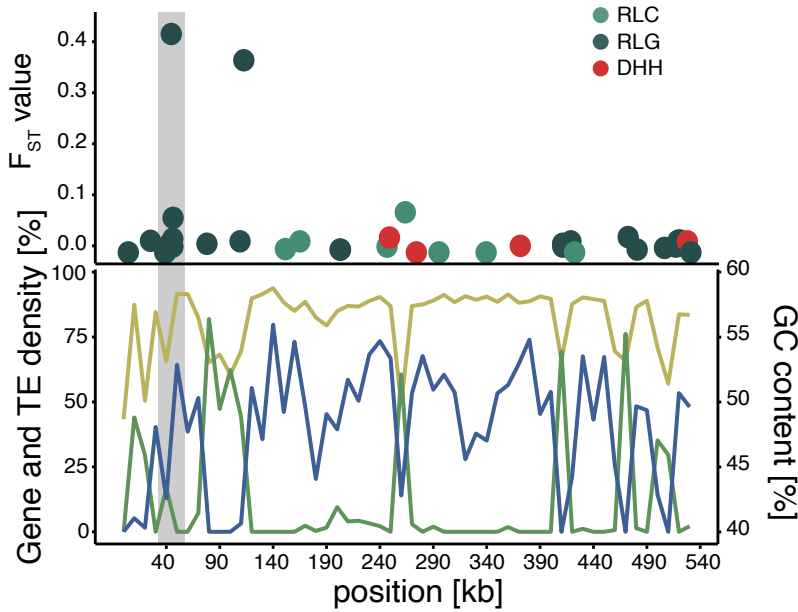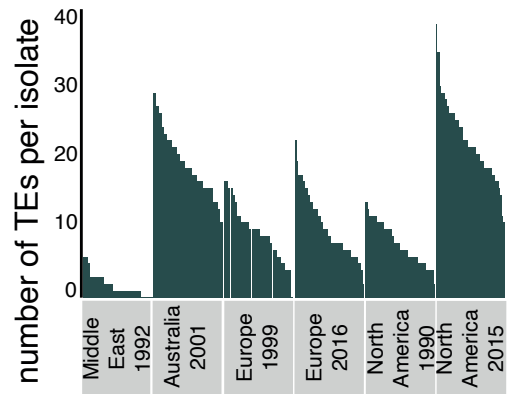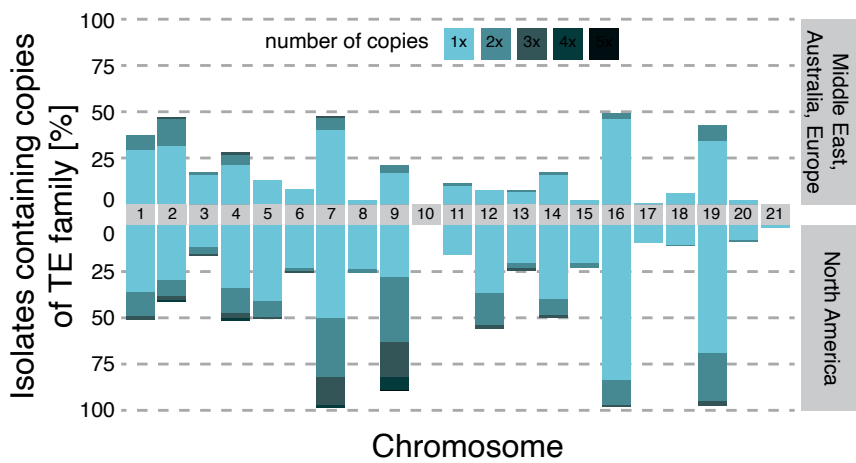1058   BUSCO completeness variation among genome assemblies. Black lines indicate the mean genome

1059   size per population. (B) Genome-wide GC content variation. (C) Core genome size variation among

1060   the isolates of the populations (excluding accessory chromosomes). (D) Correlation of core genome

1061   size and number of detected TEs. (E) Correlation of core genome size and the cumulative length of all

1062   TEs detected as inserted. (F) Correlation of core genome size and genome-wide GC content. (G)

1063   Spearman correlation matrix of BUSCO completeness, core genome size, number of detected TEs and

1064   genome-wide GC content.

1065   **Figure supplement 1.** Genome size expansion. (A) Estimated length of TE insertions per isolate and

1066   population. (B) Genome size variation per population. (C) Percentage of TEs content variation

1067   compared to the variation in genome size. (D) TE contributions to genome size variation compared to
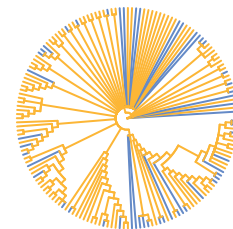
1068   full genome size.

34

**A**

Raw data processing

TE detection

Downsampling analysis

TE position correction

Mapping short reads

Detection spliced junction reads

Convertion to bed

Convertion to hint file

Outlier test

Overlap with predicted TEs

TE position correction

Final call presence/absence variation

Singleton validation

**B**



Number of detected insertions

Middle East 1992
Australia 2001
Europe 1999
Europe 2016
North America 1990
North America 2015

Coverage (x)

**C**

**i**

TS    TE    TS    isolate

reference genome

Illumina reads

TSD

**ii**

leftmost    -100 bp    TS    TS    +100 bp    rightmost

illumina reads

number stop points

**D**

**i**

isolate

TE    reference genome

spliced junction reads

**ii**

spliced junction reads

isolate

accept as absence | confirmed presence | Unassigned | reject presence

predicted TE in ngs_te_mapper

**E**

**i**

TS    TE    TS    isolate

structural variation locus    PacBio sequence

reference genome

**ii**

500bp    TS    500bp    reference genome

PacBio sequence

TE consensus sequence

**A**

n=56  n=97  n=33  n=53  n=30  n=27

- Middle East 1992
- Australia 2001
- Europe 1999
- Europe 2016
- North America 1990
- North America 2015

**B** Genome-wide SNPs

PC 2 (5.20 %)

PC 1 (8.00 %)

**C** Genome-wide SNPs (subset)

PC 2 (6.10 %)

PC 1 (10.40 %)

**D** Transposable element loci

PC 2 (1.57 %)

PC 1 (2.39 %)

Figure legend:
- Middle East 1992
- Australia 2001
- Europe 1999
- Europe 2016
- North America 1990
- North America 2015

**A** BUSCO completeness [%]

**B** GC content [%]

**C** Core genome size [Mb] — Middle East 1992, Australia 2001, Europe 1999, Europe 2016, North America 1990, North America 2015

**D** Core genome size [Mb] vs TEs per isolate — R = 0.76, p < 2.2e⁻¹⁶

**E** Core genome size [Mb] vs cumulative TE length [Mb] — R = 0.78, p < 2.2e⁻¹⁶

**F** Core genome size [Mb] vs Genome-wide GC [%] — R = 0.95, p < 2.2e⁻¹⁶

**G** Spearman Correlation