

# 1 **A population-level invasion by transposable** 2 **elements in a fungal pathogen**

3  
4  
5 Ursula Oggenfuss<sup>1</sup>, Thomas Badet<sup>1</sup>, Thomas Wicker<sup>2</sup>, Fanny E. Hartmann<sup>3,4</sup>, Nikhil K. Singh<sup>1</sup>, Leen  
6 N. Abraham<sup>1</sup>, Petteri Karisto<sup>4,6</sup>, Tiziana Vonlanthen<sup>4</sup>, Christopher C. Mundt<sup>5</sup>, Bruce A. McDonald<sup>4</sup>,  
7 Daniel Croll<sup>1,\*</sup>

8  
9  
10 <sup>1</sup> Laboratory of Evolutionary Genetics, Institute of Biology, University of Neuchâtel, 2000 Neuchâtel,  
11 Switzerland

12 <sup>2</sup> Institute for Plant and Microbial Biology, University of Zurich, Zurich, Switzerland

13 <sup>3</sup> Ecologie Systématique Evolution, Bâtiment 360, Univ. Paris-Sud, AgroParisTech, CNRS,  
14 Université Paris-Saclay, 91400 Orsay, France

15 <sup>4</sup> Plant Pathology, Institute of Integrative Biology, ETH Zurich, Zurich, Switzerland

16 <sup>5</sup> Department of Botany and Plant Pathology, Oregon State University, Corvallis, OR 97331-2902,  
17 USA

18 <sup>6</sup> Department of Evolutionary Biology and Environmental Studies, University of Zurich, Zurich,  
19 Switzerland

20  
21 \* Author for correspondence: [daniel.croll@unine.ch](mailto:daniel.croll@unine.ch)

22  
23  
24 Data availability: Sequence data is deposited on the NCBI Short Read Archive under the accession  
25 numbers PRJNA327615, PRJNA596434 and PRJNA178194.

26  
27 Author contributions: UO and DC conceived the study, UO, TW and DC designed analyses, UO, TB  
28 and TV performed analyses, FEH, NKS, LNA, PK, CCM and BAM provided samples/datasets, BAM  
29 and DC provided funding, UO and DC wrote the manuscript with input from co-authors. All authors  
30 reviewed the manuscript and agreed on submission.

31 **ABSTRACT**

32 Transposable elements (TEs) are key drivers of adaptive evolution within species. Yet, the  
33 propagation of TEs across the genome can be highly deleterious and ultimately lead to genome  
34 expansions. Hence, TE activity is likely under complex selection regimes within species. To address  
35 this, we analyzed a large whole-genome sequencing dataset of the fungal wheat pathogen  
36 *Zymoseptoria tritici* harboring TE-mediated adaptations to overcome host defenses and fungicides.  
37 We built a robust map of genome-wide TE insertion and deletion loci for six populations and 284  
38 fungal individuals across the world. We identified a total of 2'456 unfixated TE loci within the species  
39 and a significant excess of rare insertions indicating strong purifying selection. A subset of TEs  
40 recently swept to near complete fixation with at least one locus likely contributing to higher levels of  
41 fungicide resistance. TE-driven adaptation was also supported by evidence for selective sweeps. In  
42 parallel, we identified a substantial genome-wide expansion of TE families from the pathogen's  
43 center of origin to more recently founded populations, suggesting that population bottlenecks played  
44 a major role in shaping TE content of the genome. The most dramatic expansion occurred among a  
45 pair of North American populations collected in the same field at an interval of 25 years. We show  
46 that both the activation of specific TEs and relaxed purifying selection likely underpin the  
47 expansion. Our study disentangles the effects of selection and TE bursts leading to intra-specific  
48 genome expansions, providing a model to recapitulate TE-driven genome evolution over deeper  
49 evolutionary timescales.

50

51

## 52 INTRODUCTION

53 Transposable elements (TEs) are mobile repetitive DNA sequences with the ability to independently  
54 insert into new regions of the genome. TEs are major drivers of genome instability and epigenetic  
55 change (Eichler & Sankoff, 2003). Insertion of TEs can disrupt coding sequences, trigger  
56 chromosomal rearrangements, or alter expression profiles of adjacent genes (Lim, 1988; Petrov *et*  
57 *al.*, 2003; Slotkin & Martienssen, 2007; Hollister & Gaut, 2009; Oliver *et al.*, 2013). Hence, TE  
58 activity can have phenotypic consequences and impact host fitness. While TE insertion dynamics are  
59 driven by the selfish interest for proliferation, the impact on the host can range from beneficial to  
60 highly deleterious. For instance, TE insertions were shown to cause upregulation of genes  
61 influencing coloration and cold adaptation of fruits (Butelli *et al.*, 2012; Zhang *et al.*, 2019). The  
62 most dramatic examples of TE insertions underpinned rapid adaptation of populations or species  
63 (Feschotte, 2008; Chuong *et al.*, 2017), particularly following drastic environmental changes or  
64 colonization events. In the peppered moth, a TE insertion into an intron caused a darker phenotype,  
65 which provided better camouflage on tree bark darkened by pollution (van't Hof *et al.*, 2016). In  
66 *Drosophila melanogaster*, developmental timing adapted after migration to North America based on  
67 TE-mediated up-regulation of juvenile hormone production (González *et al.*, 2009). However, many  
68 studies indicate that the fate of TEs in populations is largely determined by purifying selection  
69 (Rizzon *et al.*, 2003; Walser *et al.*, 2006; Cridland *et al.*, 2013; Stuart *et al.*, 2016; Lai *et al.*, 2017;  
70 Stritt *et al.*, 2017). Thus, the broad range of fitness outcomes associated with individual TE  
71 insertions suggests that populations are likely to evolve complex selection regimes to constrain the  
72 transposition activity of TEs.

73

74 The fate of a new TE insertion in a population is dependent on the joint impact of selection and  
75 demography. TE insertions that have deleterious effects should be under strong purifying selection  
76 and are expected to be purged quickly from populations (Blumenstiel *et al.*, 2014). TE insertions  
77 with neutral effects on host fitness are governed by genetic drift alone. Hence, neutral TE insertion  
78 frequencies can fluctuate according to the strength of drift and historical bottlenecks. Effective

79 population size becomes crucial to determine these dynamics and TEs may reach fixation with a  
80 high probability in small populations (Jurka *et al.*, 2011). Low frequency TEs tend to be young or  
81 slightly deleterious insertions while high frequency TEs tend to be old insertions (Barron *et al.*,  
82 2014). Beneficial TEs are expected to experience strong positive selection and rapid fixation such as  
83 observed for the dark phenotype in the peppered moth (Barron *et al.*, 2014) (van't Hof *et al.*, 2016).  
84 In addition to negative selection against newly inserted TEs, genomic defense mechanisms can  
85 directly disable transposition activity. Across eukaryotes, epigenetic silencing is a shared defense  
86 mechanism against TEs (Slotkin & Martienssen, 2007). Fungi evolved an additional and highly  
87 specific defense system introducing repeat-induced point (RIP) mutations into any nearly identical  
88 set of sequences. TE control by RIP and RIP-like mechanisms shows a patchy distribution across the  
89 fungal tree of life and may carry significant costs, including occasional leakage of RIP into adjacent  
90 regions (Galagan & Selker, 2004; Rouxel *et al.*, 2011). Thus, the spread of TEs across the genome  
91 and the population-level frequencies at individual TE loci will be governed by a complex set of  
92 factors. However, the relative importance of demography, selection and genomic defenses on TE  
93 dynamics remains poorly understood.

94

95 A crucial property predicting the invasion success of TEs in a genome is the transposition rate. TEs  
96 mostly expand through family-specific bursts of transposition followed by prolonged phases of  
97 transposition inactivity. Bursts of insertions of different retrotransposon families were observed  
98 across eukaryotic lineages including *Homo sapiens*, *Zea mays*, *Oryza sativa* and *Blumeria graminis*  
99 (Shen *et al.*, 1991; SanMiguel *et al.*, 1998; Eichler & Sankoff, 2003; Lu *et al.*, 2017; Frantzeskakis  
100 *et al.*, 2018). Prolonged bursts without effective counter-selection are thought to underpin genome  
101 expansion. In the symbiotic fungus *Cenococcum geophilum*, the burst of TEs resulted in a  
102 dramatically expanded genome compared to closely related species (Peter *et al.*, 2016). Similarly, a  
103 burst of a TE family in brown hydras led to an approximately three-fold increase of the genome size  
104 compared to related hydras (Wong *et al.*, 2019). Across the tree of life, genome sizes vary by  
105 multiple orders of magnitude and enlarged genomes are invariably colonized by TEs (Kidwell,  
106 2002). Population size variation is among the few correlates of genome size across major groups,

107 suggesting that the efficacy of selection plays an important role in controlling TE activity (Lynch,  
108 2007). Reduced selection efficacy against deleterious TE insertions is expected to lead to a ratchet-  
109 like increase in genome size. TE-rich, enlarged genomes often show an isochores structure alternating  
110 gene-rich and TE-rich regions (Rouxel *et al.*, 2011). Genomic compartments rich in TEs often  
111 harbor genes showing high variability and may harbor rapidly evolving genes such as effectors in  
112 pathogens or resistance genes in plants (Raffaele & Kamoun, 2012; Jiao & Schneeberger, 2019).  
113 Hence, incipient genome expansions are likely driven by population-level processes such as TE  
114 insertion dynamics, strength of selection and demography.

115

116 The fungal wheat pathogen *Zymoseptoria tritici* recently gained major TE-mediated adaptations to  
117 colonize host plants and tolerate environmental and agricultural stress (Omrane *et al.*, 2015, 2017;  
118 Krishnan *et al.*, 2018; Meile *et al.*, 2018). Clusters of TEs are often associated with genes encoding  
119 important pathogenicity functions (*i.e.* effectors), recent gene gains or losses (Hartmann & Croll,  
120 2017), and major chromosomal rearrangements (Croll *et al.*, 2013; Plissonneau *et al.*, 2016).  
121 Transposition activity of TEs also had a genome-wide impact on gene expression profiles during  
122 infection (Fouché *et al.*, 2019). The compact genome of ~39 Mb is completely assembled and  
123 contains ~17% TEs (Goodwin *et al.*, 2011; Dhillon *et al.*, 2014). The well-characterized  
124 demographic history of the pathogen and evidence for recent TE-mediated adaptations make *Z.*  
125 *tritici* an ideal model to recapitulate the process of TE insertion dynamics, adaptive evolution and  
126 changes in genome size at the population level.

127

128 We aimed to retrace the population-level context of TE insertion dynamics across the species range  
129 by analyzing six populations sampled on four continents for a total of 284 genomes. We first aimed  
130 to develop a robust pipeline to detect newly inserted TEs using short read sequencing datasets. Then,  
131 we tested for the strength of purifying selection against TE insertions within and across populations.  
132 In addition, we searched for signatures of adaptive evolution driven by TE activity across  
133 populations. Using knowledge of the colonization history of the pathogen, we analyzed whether  
134 population bottlenecks were associated with substantial changes in the TE content of individual

135 genomes. In particular, we tested whether geographically isolated populations experienced distinct  
136 bursts of TEs.

137

138

## 139 **METHODS**

### 140 FUNGAL ISOLATE COLLECTION AND SEQUENCING

141 We analyzed 295 *Z. tritici* isolates covering six populations originating from four geographic  
142 locations and four continents (Supplementary Table S1), including: Middle East 1992 ( $n = 30$   
143 isolates, Nahal Oz, Israel), Australia 2001 ( $n = 27$ , Wagga Wagga), Europe 1999 ( $n = 33$ , Berg am  
144 Irchel, Switzerland), Europe 2016 ( $n = 52$ , Eschikon, ca. 15km from Berg am Irchel, Switzerland),  
145 North America 1990 and 2015 ( $n = 56$  and  $n = 97$ , Willamette Valley, Oregon, United States)  
146 (McDonald *et al.*, 1996; Linde *et al.*, 2002; Zhan *et al.*, 2002, 2003, 2005). Illumina short read data  
147 from the Middle East, Australia, European 1999 and North American 1990 populations were  
148 obtained from the NCBI Short Read Archive under the BioProject PRJNA327615 (Hartmann *et al.*,  
149 2017). For, the Switzerland 2016 and Oregon 2015 populations, asexual spores were harvested from  
150 infected wheat leaves from naturally infected fields and grown in YSB liquid media including 50  
151  $\text{mgL}^{-1}$  kanamycin and stored in silica gel at  $-80^{\circ}\text{C}$ . High-quality genomic DNA was extracted from  
152 liquid cultures using the DNeasy Plant Mini Kit from Qiagen (Venlo, Netherlands). The isolates  
153 were sequenced on an Illumina HiSeq in paired-end mode and raw reads were deposited on the  
154 NCBI Short Read Archive under the BioProject PRJNA596434.

155

### 156 TE INSERTION DETECTION

157 The quality of Illumina short reads was determined with FastQC version 0.11.5 (Figure 1A)  
158 (Andrews *et al.*, 2013). To remove spuriously sequenced Illumina adaptors and low quality reads,  
159 we trimmed the sequences with Trimmomatic version 0.36, using the following filter parameters:  
160 illuminaclip:TruSeq3-PE-2.fa:2:30:10 leading:10 trailing:10 slidingwindow:5:10 minlen:50 (Bolger  
161 *et al.*, 2014). We created repeat sequence consensi for TE families (Supplementary File S1) in the

162 complete reference genome IPO323 (Goodwin *et al.*, 2011) with RepeatModeler version open-4.0.7  
163 based on the RepBase Sequence Database (Smit & Hubley; Bao *et al.*, 2015). TE classification into  
164 superfamilies and families was based on an approach combining detection of conserved protein  
165 sequences and tools to detect non-autonomous TEs (Badet *et al.*, 2019). To detect TE insertions, we  
166 used the R-based tool `ngs_te_mapper` version 79ef861f1d52cdd08eb2d51f145223fad0b2363c  
167 integrated into the McClintock pipeline version 20cb912497394fabddcdad175402adacf5130bd1,  
168 using `bwa` version 0.7.4-r385 to map Illumina short reads, `samtools` version 0.1.19 to convert  
169 alignment file formats and R version 3.2.3 (Li & Durbin, 2009; Li *et al.*, 2009; Linheiro & Bergman,  
170 2012; Nelson *et al.*, 2017; R Core Team, 2017).

171

#### 172 DOWN-SAMPLING ANALYSIS

173 We performed a down-sampling analysis to estimate the sensitivity of the TE detection with  
174 `ngs_te_mapper` based on variation in read depth. We selected one isolate per population matching  
175 the average coverage of the population. We extracted the per-base pair read depth with the  
176 `genomecov` function of `bedtools` version 2.27.1 and calculated the genome-wide mean read depth  
177 (Quinlan & Hall, 2010). The number of reads in the original fastq file was reduced in steps of 10%  
178 to simulate the impact of reduced coverage. We analyzed each of the obtained reduced read subsets  
179 with `ngs_te_mapper` using the same parameters as described above. The correlation between the  
180 number of detected insertions and the read depth was visualized using the function `nls` with model  
181 `SSlogis` in R, (Wickham, 2016). The number of detected TEs increased with the number of reads  
182 until reaching a plateau indicating saturation (Figure 1B). Saturation was reached at a coverage of  
183 approximately 15X, hence we retained only isolates with an average read depth above 15X for  
184 further analyses. We thus excluded one isolate from the Oregon 2015 population and ten isolates  
185 from the Switzerland 2016 population.

186

187 VALIDATION PROCEDURE FOR PREDICTED TE INSERTIONS

188 ngs\_te\_mapper detects the presence but not the absence of a TE at any given locus. We devised  
189 additional validation steps to ascertain both the presence as well absence of a TE across all loci in all  
190 individuals. For TE loci with missing information about presence or absence, we conducted further  
191 analyses. TEs absent in the reference genome were validated by re-analyzing mapped Illumina reads.  
192 Reads spanning both parts of a TE sequence and an adjacent chromosomal sequence should only  
193 map to the reference genome sequence and cover the target site duplication (TSD) of the TE (Figure  
194 1C). We used bowtie2 version 2.3.0 with the parameter --very-sensitive-local to map Illumina short  
195 reads of each isolate on the reference genome IPO323 (Langmead & Salzberg, 2012). Mapped  
196 Illumina short reads were then sorted and indexed with samtools and the resulting bam file was  
197 converted to a bed file with the function bamtobed in bedtools. We extracted all mapped reads with  
198 an end point located within 100 bp of the TSD (Figure 1C). We tested whether the number of reads  
199 with a mapped end around the TSD significantly deviated if the mapping ended exactly at the  
200 boundary. A mapped read ending exactly at the TSD boundary is indicative of a split read mapping  
201 to a TE sequence not present in the reference genome. To test for the deviation in the number of read  
202 mappings around the TSD, we used a Poisson distribution and the *ppois* function in R version 3.5.1  
203 (Figure 1C). We identified a TE as present in an isolate if tests on either side of the TSD had a *p*-  
204 value < 0.001 (Supplementary Table S1, S2, Figure S1B).

205

206 For TEs present in the reference genome, we analyzed evidence for spliced junction reads spanning  
207 the region containing the TE. Spliced reads are indicative of a discontinuous sequence and, hence,  
208 absence of the TE in a particular isolate (Figure 1D). We used STAR version 2.5.3a to detect spliced  
209 junction reads with the following set of parameters: --runThreadN 1 --outFilterMultimapNmax 100 -  
210 --winAnchorMultimapNmax 200 --outSAMmultNmax 100 --outSAMtype BAM Unsorted --  
211 outFilterMismatchNmax 5 --alignIntronMin 150 --alignIntronMax 15000 (Dobin *et al.*, 2012). We  
212 then sorted and indexed the resulting bam file with samtools and converted split junction reads with  
213 the function bam2hints in bamtools version 2.5.1 (Barnett *et al.*, 2011). We selected loci without  
214 overlapping spliced junction reads using the function intersect in bedtools with the parameter -loj -v.



215 We considered a TE as truly absent in an isolate if `ngs_te_mapper` did not detect a TE and no  
216 evidence for spliced junction reads were found. If the absence of a TE could not be confirmed by  
217 spliced junction reads, we labelled the genotype as missing. Finally, we excluded TE loci with more  
218 than 20% missing data from further investigations (Figure 1D and Supplementary Figure S1C).

219

## 220 CLUSTERING OF TE INSERTIONS INTO LOCI

221 We identified insertions across isolates as being the same locus if all detected TEs belonged to the  
222 same TE family and insertion sites differed by  $\leq 100$  bp (Supplementary Figure S2). We used the R  
223 package *GenomicRanges* version 1.28.6 with the functions `makeGRangesFromDataFrame` and  
224 `findOverlaps` and the R package *devtools* version 1.13.4 (Lawrence *et al.*, 2013; Wickham & Chang,  
225 2016). We used the R package *dplyr* version 0.7.4 to summarize datasets (Wickham *et al.*, 2017).  
226 Population-specific frequencies of insertions were calculated with the function `allele.count` in the R  
227 package *hierfstat* version 0.4.22 (Goudet & Jombart, 2015). We conducted a principal component  
228 analysis for TE insertion frequencies filtering for a minor allele frequency  $\geq 5\%$ . We also performed  
229 a principal component analysis for genome-wide single nucleotide polymorphism (SNP) data  
230 obtained from Hartmann *et al.* (2017). As described previously, SNPs were hard-filtered with  
231 `VariantFiltration` and `SelectVariants` tools integrated in the Genome Analysis Toolkit (GATK)  
232 (McKenna *et al.*, 2010). SNPs were removed if any of the following filter conditions applied:  
233 `QUAL<250; QD<20.0; MQ<30.0; -2 > BaseQRankSum > 2; -2 > MQRankSum > 2; -2 >`  
234 `ReadPosRankSum > 2; FS>0.1`. SNPs were excluded with `vcftools` version 0.1.17 and `plink` version  
235 1.9 requiring a genotyping rate  $>90\%$  and a minor allele frequency  $>5\%$  (<https://www.cog->  
236 [genomics.org/plink2](https://www.cog-genomics.org/plink2), Chang *et al.*, 2015). Finally, we converted tri-allelic SNPs to bi-allelic SNPs  
237 by recoding the least frequent allele as a missing genotype. Principal component analysis was  
238 performed using the *gdsfmt* and *SNPRelate* packages in R (Zheng *et al.*, 2012, 2017). For a second  
239 principal component analysis with a reduced set of random markers, we randomly selected SNPs  
240 with `vcftools` and the following set of parameters: `--maf 0.05 --thin 200'000` to obtain an  
241 approximately equivalent number of SNPs as TE loci.

242

#### 243 GENOMIC LOCATION OF TE INSERTIONS

244 To characterize the genomic environment of TE insertion loci, we split the reference genome into  
245 non-overlapping windows of 10 kb using the function splitter from EMBOSS version 6.6.0 (Rice *et*  
246 *al.*, 2000). TEs were located in the reference genome using RepeatMasker providing consensi from  
247 *RepeatModeler Open-1.0* (Smit & Hubley; <http://www.repeatmasker.org>). To analyze coding  
248 sequence, we retrieved the gene annotation for the reference genome (Grandaubert *et al.*, 2015). We  
249 estimated the percentage covered by genes or TEs per window using the function intersect in  
250 bedtools. Additionally, we calculated the GC content using the tool `get_gc_content`  
251 ([https://github.com/spundhir/RNA-Seq/blob/master/get\\_gc\\_content.pl](https://github.com/spundhir/RNA-Seq/blob/master/get_gc_content.pl)). We also extracted the  
252 number of TEs present in 1 kb windows up- and downstream of each annotated gene with the  
253 function window in bedtools with the parameters `-l 1000 -r 1000` and calculated the relative  
254 distances with the closest function in bedtools. For the TEs inserted into genes, we used the intersect  
255 function in bedtools to distinguish intron and exon insertions with the parameters `-wo` and `-v`,  
256 respectively. For each 100 bp segment in the 1kb windows as well as for introns and exons, we  
257 calculated the mean number of observed TE insertions per base pair.

258

#### 259 POPULATION DIFFERENTIATION IN TE FREQUENCIES

260 We calculated Nei's fixation index ( $F_{ST}$ ) between pairs of populations using the R packages *hierfstat*  
261 and *adegenet* version 2.1.0 (Jombart, 2008; Jombart & Ahmed, 2011). To understand the  
262 chromosomal context of TE insertion loci across isolates, we analyzed draft genome assemblies. We  
263 generated *de novo* genome assemblies for all isolates using SPAdes version 3.5.0 with the parameter  
264 `--careful` and a kmer range of "21, 29, 37, 45, 53, 61, 79, 87" (Bankevich *et al.*, 2012). We used  
265 `blastn` to locate genes adjacent to TE insertion loci on genomic scaffolds of each isolate. We then  
266 extracted scaffold sequences surrounding 10 kb up- and downstream of the localized gene with the  
267 function `faidx` in samtools and reverse complemented the sequence if needed. Then, we performed  
268 multiple sequence alignments for each locus across all isolates with MAFFT version 7.407 with

269 parameter --maxiterate 1000 (Kato & Standley, 2013). We performed visual inspections to ensure  
270 correct alignments across isolates using Jalview version 2.10.5 (Waterhouse *et al.*, 2009). To  
271 generate phylogenetic trees of individual gene or TE loci, we extracted specific sections of the  
272 alignment using the function extractalign in EMBOSS and converted the multiple sequence  
273 alignment into PHYLIP format with jmodeltest version 2.1.10 using the -getPhylip parameter. We  
274 then estimated maximum likelihood phylogenetic trees with the software PhyML version 3.0, the  
275 K80 substitution model and 100 bootstraps on the ATGC South of France bioinformatics platform  
276 (Guindon & Gascuel, 2003; Guindon *et al.*, 2010; Darriba *et al.*, 2012). Bifurcations with a  
277 supporting value lower than 10% were collapsed in TreeGraph version 2.15.0-887 beta and trees  
278 were visualized as circular phylograms in Dendroscope version 2.7.4 (Huson *et al.*, 2007; Stöver &  
279 Müller, 2010). For loci showing complex rearrangements, we generated synteny plots using 19  
280 completely sequenced genomes from the same species using the R package *genoplots* version 0.8.9  
281 (Guy *et al.*, 2010; Badet *et al.*, 2019).

282 We analyzed signatures of selective sweeps using the extended haplotype homozygosity (EHH) tests  
283 (Sabeti *et al.*, 2007) implemented in the R package REHH (Gautier & Vitalis, 2012). We analyzed  
284 within-population signatures based on the iHS statistic and chose a maximum gap distance of 20 kb.  
285 We also analyzed cross-population EHH (XP-EHH) signatures testing the following two population  
286 pairs: North America 1990 versus North America 2015, Europe 1999 versus Europe 2016. We  
287 defined significant selective sweeps as being among the 99.9th percentile outliers of the iHS and  
288 XP-EHH statistics. Significant SNPs at less than 5 kb were clustered into a single selective sweep  
289 region adding +/- 2.5 kb. Finally, we analyzed whether TE loci were within 10 kb of a region  
290 identified as a selective sweep using the function intersect from bedtools.

291

## 292 FUNGICIDE RESISTANCE ASSAY

293 To quantify susceptibility towards propiconazole we performed a microtiter plate assay. Isolates  
294 were grown on yeast malt sucrose agar for five days and spores were harvested. We then tested for  
295 growth inhibition by growing spores ( $2.5 \times 10^4$  spores/ml) in Sabouraud-dextrose liquid medium

296 with differing concentrations of propiconazole (0.00006, 0.00017, 0.0051, 0.0086, 0.015, 0.025,  
297 0.042, 0.072, 0.20, 0.55, 1.5 mg/L). We incubated the plates stationary in the dark at 21°C and 80%  
298 relative humidity for four days and measured optical density at 605 nm. We calculated EC<sub>50</sub> with the  
299 R package *drc* (Ritz & Streibig, 2005).

300

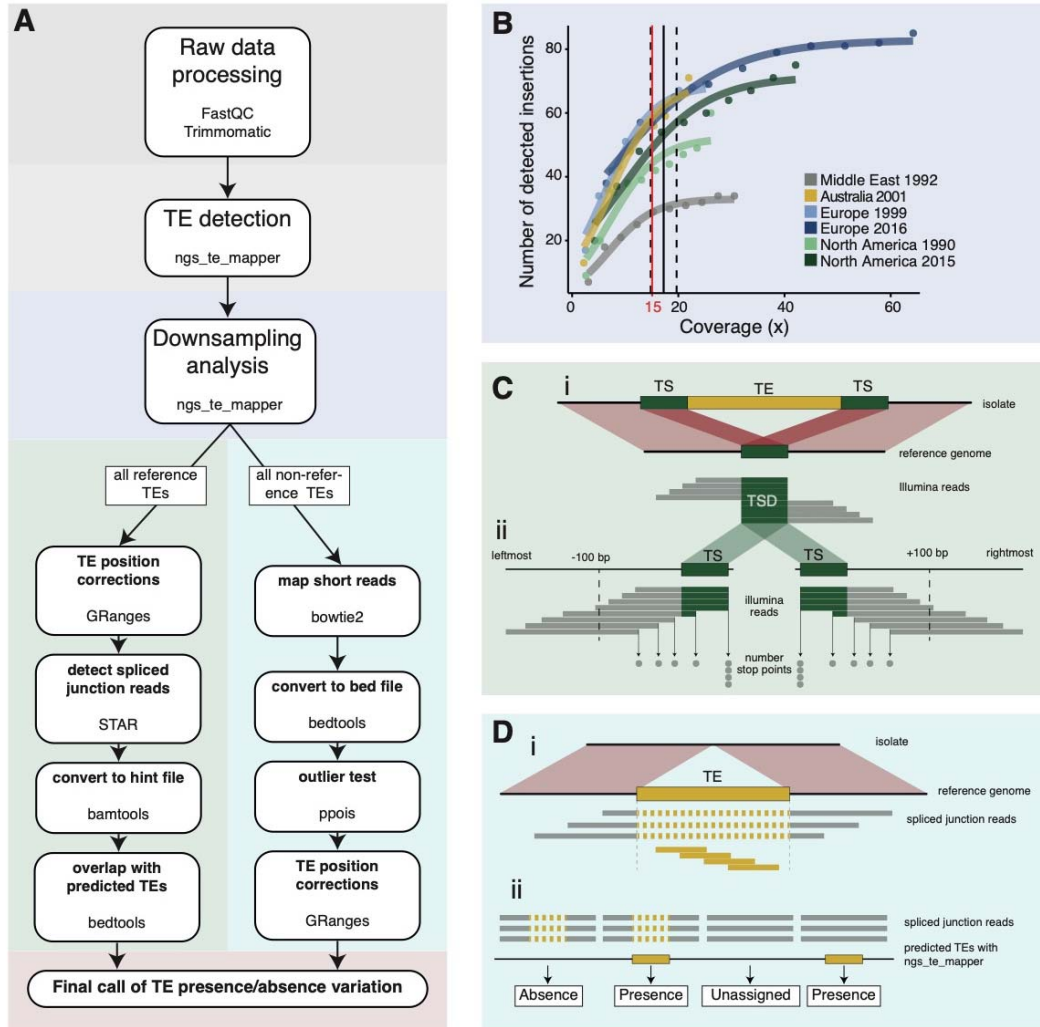
301

## 302 **RESULTS**

### 303 A DYNAMIC TE LANDSCAPE SHAPED BY STRONG PURIFYING SELECTION

304 We analyzed 284 genomes from a worldwide set of six populations spanning the distribution range  
305 of the wheat pathogen *Z. tritici*. To ascertain the presence or absence of TEs across the genome, we  
306 developed a robust pipeline (Figure 1A) to address the fact that *ngs\_te\_mapper* is only testing for TE  
307 presence and not TE absence. In summary, we called TE insertions by identifying reads mapping  
308 both to a TE sequence and a specific location in the reference genome. Then, we assessed the  
309 minimum sequencing coverage to reliably recover TE insertions, tested for evidence of TEs using  
310 read depth at target site duplications, and scanned the genome for mapped reads indicating gaps at  
311 TE loci. We found robust evidence for a total of 18'864 TE insertions grouping into 2'465 individual  
312 loci. More than 30% of these loci have singleton TEs (Figure 2B, Supplementary Table S3). An  
313 overwhelming proportion of loci have a TE frequency below 1%. This pattern strongly supports the  
314 hypothesis that TEs actively copy into new locations but also indicates that strong purifying  
315 selection maintains nearly all TEs at low frequency (Figure 2B). We found a higher density of TE  
316 loci on accessory chromosomes, which are not shared among all isolates of the species, compared to  
317 core chromosomes (Figure 2C). This suggests relaxed selection against TE insertion on the  
318 functionally dispensable accessory chromosomes.

319

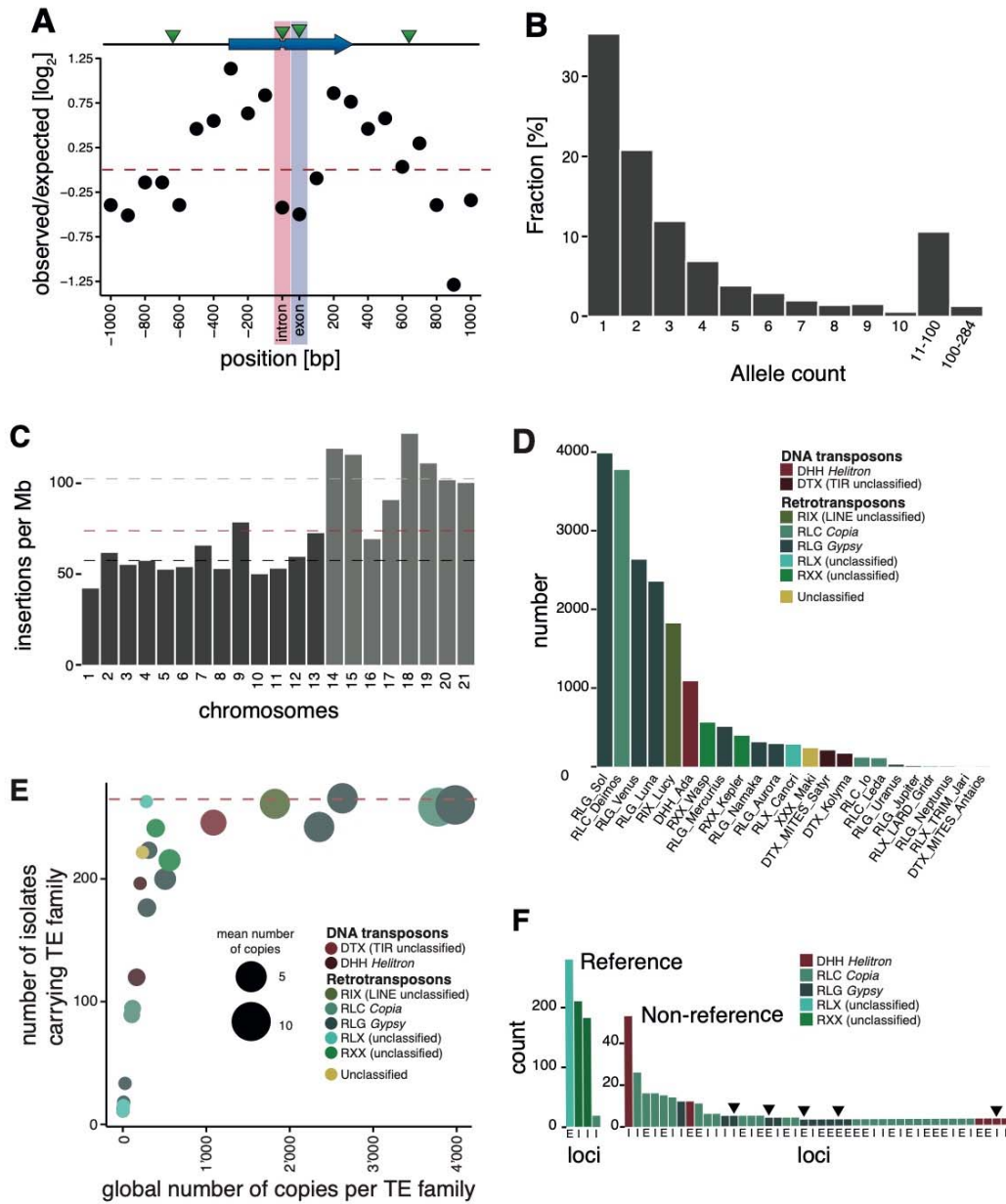


320

321 **Figure 1: Overview methods and validations of transposable element (TE) insertions:** (A)  
 322 Bioinformatic pipeline. (B) Read depth down-sampling analysis for one isolate per population with  
 323 an average coverage of the population. The vertical black line indicates the coverage at which on  
 324 average 90% of the maximally detectable variants were recovered. Dashed black lines indicate the  
 325 standard error. The threshold for a minimal mean coverage was set at 15X (red line). (C) Validation  
 326 of insertions not present in the reference genome. (i) TE insertions that are not present in the  
 327 reference genome show a duplication of the target site and the part of the reads that covers the TE  
 328 will not be mapped against the reference genome. We thus expect reads to map to the TE  
 329 surrounding region and the target site duplication but not the TE itself. At the target site, a local  
 330 duplication of read depth is expected. (ii) We selected all reads in an interval of 100 bp up- and  
 331 downstream including the target site duplication to detect deviations in the number of reads  
 332 terminating near the target site duplication. (D) Validation of insertions present in the reference  
 333 genome. (i) Analyses read coverage at target site duplications. (ii) Synthesis of evidence from  
 334 ngs\_te\_mapper and split read mapping to determine TE presence or absence.

335

336 Inserted TEs grouped into 11 superfamilies and 23 families with most TEs belonging to class  
337 I/retrotransposons ( $n = 2175$ ; Supplementary Figure S3A; Figure 2D). Most class I TEs are long  
338 terminal repeats (LTR) with 1'483 belonging to *Gypsy* superfamily (9 families in total) and 623  
339 belonging to the *Copia* superfamily (3 families in total). We found a further 40 loci with an insertion  
340 of long interspersed repeat (LINE) elements. A total of 289 loci contain class II/DNA transposons  
341 with most belonging to Subclass 2 and order Helitron (249 loci), and to Subclass 1 (40 loci). TE  
342 families with a high total copy number across all isolates tend to also have a high copy number per  
343 genome (Figure 2E).



344

345 **Figure 2: Transposable element (TE) landscape across populations.** (A) Number of TE  
 346 insertions 1 kb up- and downstream of genes on core chromosomes including introns and exons (100  
 347 bp windows). (B) Allele frequencies of the TE insertions across all isolates. (C) TE insertions per  
 348 Mb on core chromosomes (dark) and accessory chromosomes (light). Dashed lines represent mean  
 349 values. Red: global mean of 75.65 insertions/Mb, dark: core chromosome mean of 58.00 TEs/Mb,  
 350 light: accessory chromosome mean of 102.24 insertions/Mb). (D) Number of TE insertions per  
 351 family. (E) TE frequencies among isolates and copy numbers across the genome. The red line  
 352 indicates the maximum number of isolates ( $n = 284$ ). (F) TE insertions into introns and exons that  
 353 are present in the reference genome and TEs absent from the reference genome but present in more  
 354 than two copies in the populations. A hexagon indicates that the insertion was found in only one  
 355 population, all other insertions were found in at least two populations. I = intron insertion, E = exon  
 356 insertion.

357

358 We found 153 loci where a TE inserted into a gene, with most of these insertions being singletons ( $n$   
359 = 68) or at very low frequency (Figure 2F). Overall, TE insertions into exonic sequences were less  
360 frequent than expected compared to insertions into up- and downstream regions, a pattern consistent  
361 with effective purifying selection (Figure 2A). Interestingly, insertions into introns were also  
362 strongly under-represented, likely due to the small size of most fungal introns (50-100 bp) and the  
363 high probability of disrupting splicing or adjacent coding sequences. We also found that insertions  
364 800-1000 bp away from coding sequences of a focal gene were under-represented. Given the high  
365 gene density, with an average spacing between genes of 1,744 kb, TE insertions within 800-1000 bp  
366 of a coding gene tend to be near adjacent genes already. Taken together, TEs in the species show a  
367 high degree of transposition activity and are subject to strong purifying selection.

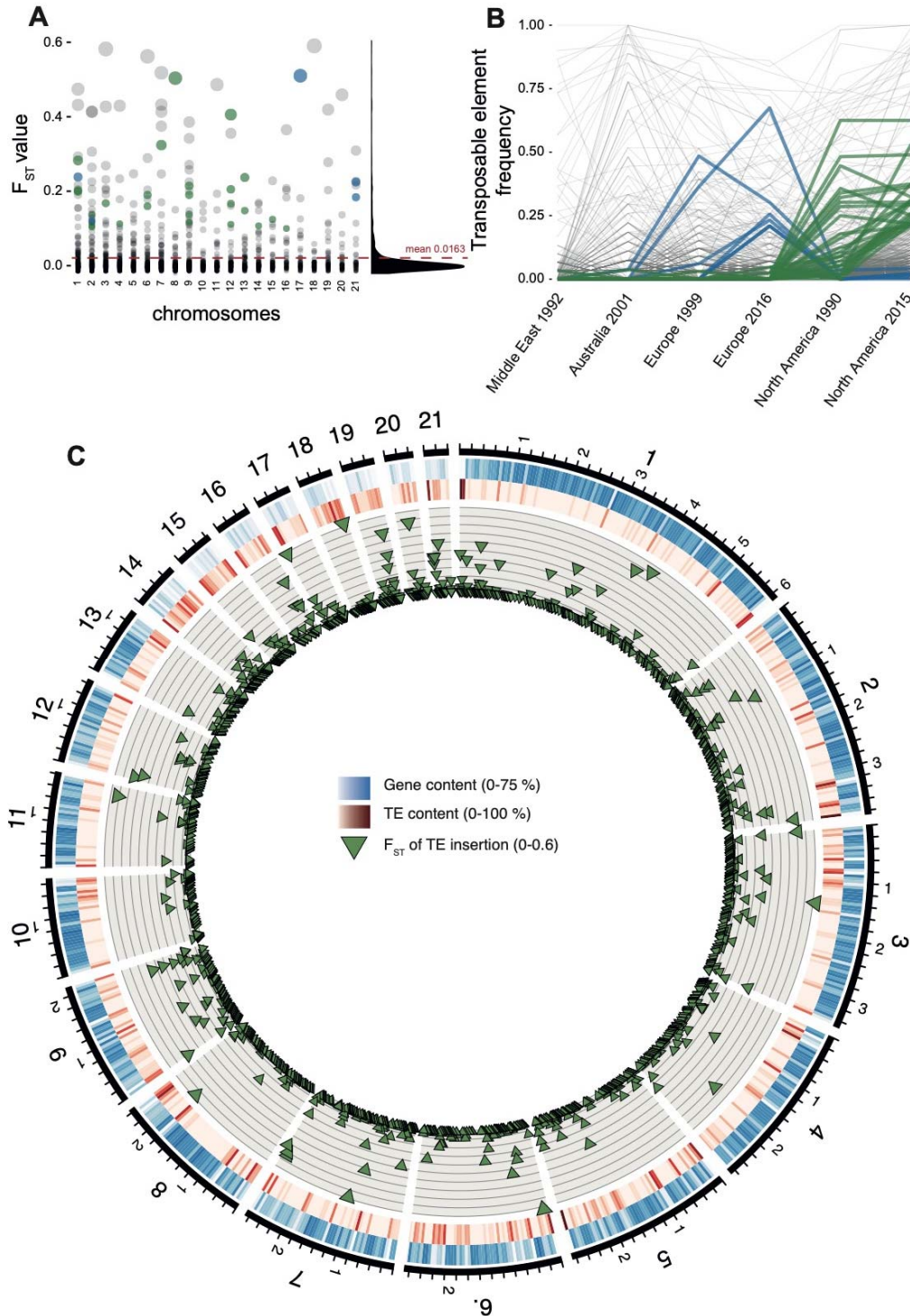
368

#### 369 DETECTION OF TE LOCI UNDER POSITIVE SELECTION

370 The dynamic transposition activity can potentially generate adaptive genetic variation. To identify  
371 potentially adaptive TE insertions, we calculated the fixation index ( $F_{ST}$ ) for each TE locus. Across  
372 all populations,  $F_{ST}$  was highly variable with a strong skew towards extremely low  $F_{ST}$  values (mean  
373  $F_{ST}$  = 0.0163; Figure 3A). High  $F_{ST}$  loci tend to have high TE frequencies in either the North  
374 American population from 2015 or the Australian population. Given our population sampling, we  
375 tested for the emergence of adaptive TE insertions either in the North American or European  
376 population pairs. We selected loci having low TE insertion frequencies (< 5%) in all populations  
377 except either the North American or European population pairs, respectively (Figure 3B). We  
378 required that the locus had a TE frequencies >20% in either the 2015 North American or 2016  
379 European population. Based on these criteria, we obtained 26 candidate loci possibly underlying  
380 local adaptation in the North American populations with 22 loci showing retrotransposon insertions,  
381 three *Helitron*, and one DNA TIR transposon (Figure 3C). In parallel, we found six loci of  
382 retrotransposons possibly underlying local adaptation in the European populations (Figure 4A and  
383 Supplementary Table S4). To further analyze evidence for TE-mediated adaptive evolution, we  
384 analyzed the whole-genome sequencing datasets for evidence of selective sweeps using selection



385 scans. Out of all 32 loci showing signatures of local adaptation in North American or European  
386 populations, we found five loci overlapping selective sweep regions.  
387



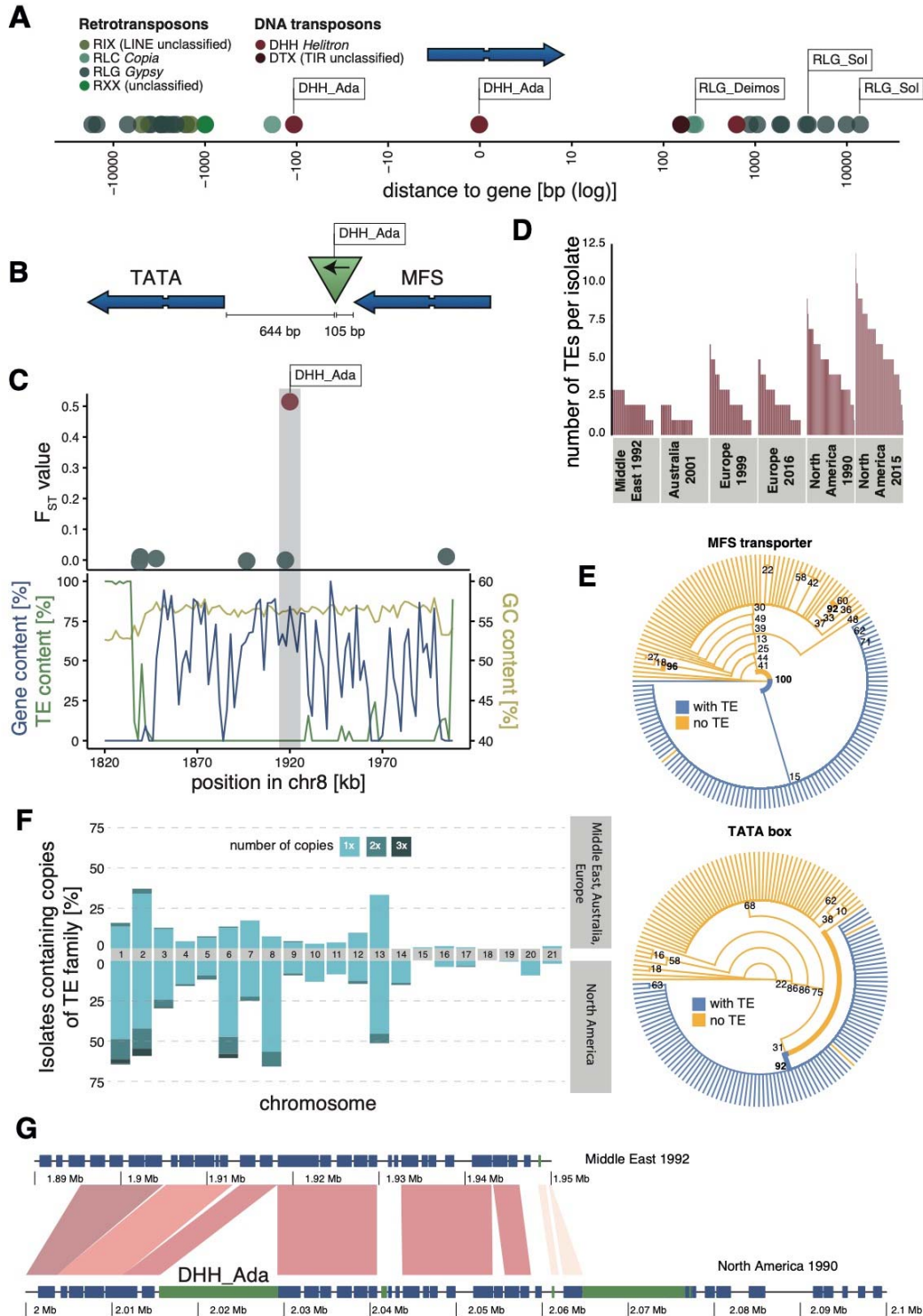
388

389 **Figure 3: Differentiation in transposable element insertions frequencies across the genome.** (A)  
390 Global pairwise  $F_{ST}$  distributions shown across the 21 chromosomes. The red horizontal line  
391 indicates the mean  $F_{ST}$  (= 0.0163). TEs with a strong frequency differences among populations are  
392 highlighted (blue: increase in Europe; green: increase in North America). (B) Allele frequency  
393 changes between the populations. Outlier TE loci are highlighted (colors as in panel A). (C) Circos  
394 plot describing from the outside to the inside: The black line indicates chromosomal position in Mb.  
395 Blue bars indicate the gene density in windows of 100 kb with darker blue representing higher gene  
396 density. Red bars indicate the TE density in windows of 100 kb with a darker red representing higher  
397 TE density. Green triangles indicate positions of TE insertions with among population  $F_{ST}$  value  
398 shown on the y-axis.  
399

400 We focused on five TE insertion loci in proximity to genes with a function that can be associated to  
401 fungicide resistance or adaptation to the host. One TE insertion on chromosome 8 is only 105 bp  
402 downstream of a major facilitator superfamily (MFS) transporter gene and 644 bp upstream of a  
403 TATA box gene (Figure 4B). MFS transporters can contribute to the detoxification of antifungal  
404 compounds in the species (Omrane *et al.*, 2017). The inserted *Helitron* TE was only found in North  
405 American populations (Figure 4G). The TE insertion occurred in a gene-rich, TE-poor region and the  
406  $F_{ST} = 0.51$  was one of the highest values of all TE loci (Figure 4C). Generally, the *Helitron* increased  
407 strongly in copy number from the Israel to the North American populations (Figure 4D, 4F). The  
408 phylogeny of the gene encoding the MFS showed a high degree of similarity for all isolates carrying  
409 the *Helitron* insertion compared to the isolates lacking the *Helitron* (Figure 4E). This is consistent  
410 with a rapid rise in frequency of the haplotype carrying the *Helitron* driven by positive selection.  
411 Another TE locus also contains a *Helitron* of the family Ada, which was found only in the two North  
412 American populations. The TE was inserted into an intron of a Phox domain-encoding gene  
413 (Supplementary Figure S7). Phox homologous domain proteins contribute to sorting membrane  
414 trafficking (Odorizzi *et al.*, 2000). A further North American possible adaptive insertion of a *Copia*  
415 Deimos TE was 229 bp upstream of a gene encoding a SNARE domain protein and 286 bp upstream  
416 of a gene encoding a flavin amine oxidoreductase and located in a region of selective sweep  
417 (Supplementary Figure S8). SNARE domains play a role in vesicular transport and membrane fusion  
418 (Bonifacino & Glick, 2004). A TE insertion locus on chromosome 12 was both upstream 1'977 bp  
419 of a gene encoding another MFS transporter and 2'672 bp of a gene encoding an alpha/beta  
420 hydrolase fold. The inserted TE belongs to the *Gypsy* family Sol (Supplementary Figure S9). A

421 *Gypsy* Sol on chromosome 2 increased in frequency in both the North American (1.8 to 75%) as well  
422 as in the European site (0 to 66.7%). This TE was inserted between a gene encoding a potential  
423 virulence factor (*i.e.* an effector) and a second gene of unknown function. Interestingly, the same  
424 locus contains a multitude of additional inserted TEs across populations (Supplementary Figure  
425 S10). We experimentally tested whether the TE insertion loci in proximity to genes could contribute  
426 to higher levels of fungicide resistance. For this, we measure growth rates of the fungal isolates in  
427 the presence or absence of an azole fungicide widely deployed against the pathogen. We found that  
428 the insertion of TEs at three loci was positively associated with higher levels of fungicide resistance  
429 suggesting TE-mediated adaptations.

430



431

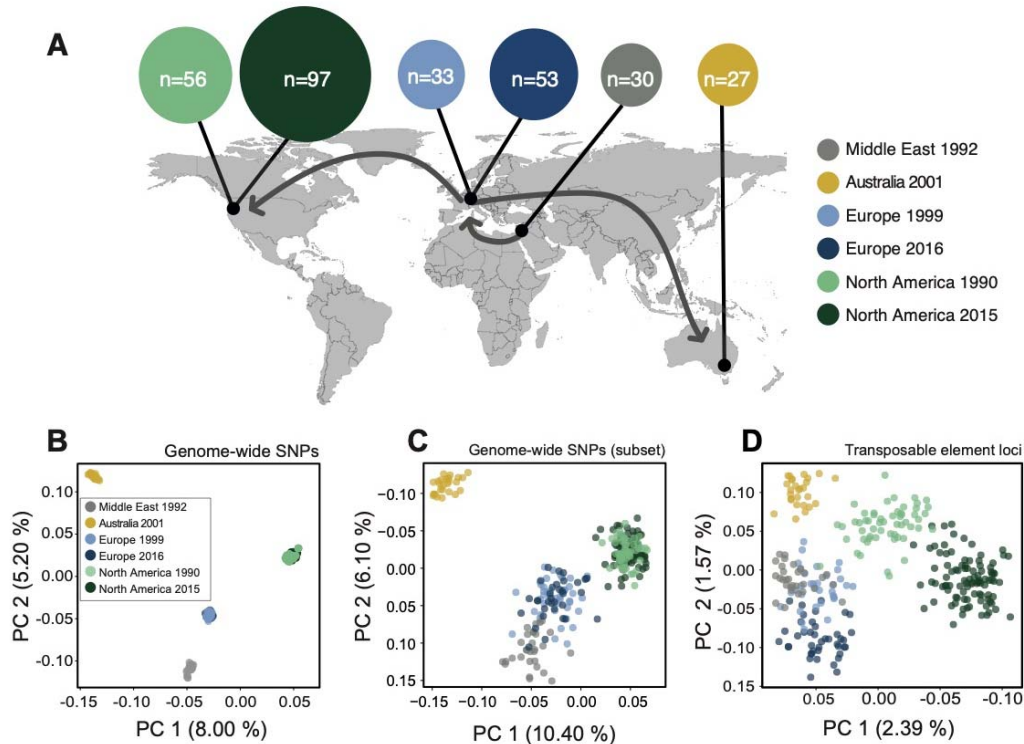
432 **Figure 4: Candidate adaptive transposable element (TE) insertions.** (A) Distribution of all  
 433 extremely differentiated TEs and their distance to the closest gene. Color indicates the superfamily.  
 434 TE sites potentially under selection according to  $F_{ST}$  are flagged. (B) Location of the *Helitron* Ada  
 435 TE insertion on chromosome 8 corresponding to its two closest genes. (C) Genomic niche of the  
 436 *Helitron* Ada TE insertion on chromosome 8:  $F_{ST}$  values for each TE insertion, gene content (blue),

437 TE content (green) and GC-content (yellow). The grey section highlights TE loci with extremely  
438 differentiated population frequencies. (D) Number of Ada copies per isolate and population. (E)  
439 Phylogenetic trees of the coding sequences of each the MFS transporter upstream and the TATA box  
440 downstream of the TE insertion. Isolates of the two North American populations and an additional  
441 11 isolates from other populations not carrying the insertion are shown. Blue color indicates TE  
442 presence, yellow indicates TE absence. (F) Frequency changes of the TE family Ada between the  
443 two North American populations compared to the other populations. Colors indicate the number of  
444 copies per chromosome. (G) Synteny plot of the Ada insertion locus on chromosome 8 between two  
445 complete genomes from the Middle East (TE missing) and North America (TE present). Figures S7-  
446 S11 show additional candidate regions.

447

#### 448 POPULATION-LEVEL EXPANSIONS IN TE CONTENT

449 If TE insertion dynamics are largely neutral across populations, TE frequencies across loci should  
450 reflect neutral population structure. To test this, we performed first a principal component analysis  
451 based on a set of six populations on four continents that represent the global genetic diversity of the  
452 pathogen (Figure 5A). The SNP set contained 900'193 genome-wide SNPs (Figure 5B). The  
453 population structure reflected the demographic history of the pathogen with clear continental  
454 differentiation and only minor within-site differentiation. In stark contrast, TE frequencies across  
455 loci showed only weak clustering by geographic origin with the Australian population being the  
456 most distinct (Figure 5D). We found a surprisingly strong differentiation of the two North American  
457 populations sampled at a 25-year interval in the same field in Oregon. To account for the lower  
458 number of TE loci, we performed an additional principal component analysis using a comparably  
459 sized SNP set to number of TE loci. Genome-wide SNPs retained the geographic signal of the  
460 broader set of SNPs (Figure 5C).



461

462 **Figure 5: Population differentiation at transposable element (TE) and genome-wide SNP loci.**  
463 (A) Sampling locations of the six populations. Middle East represents the region of origin of the  
464 pathogen. In North America, the two populations were collected at an interval of 25 years in the  
465 same field in Oregon. In Europe, two populations were collected at an interval of 17 years from two  
466 fields in Switzerland <20 km apart. Dark arrows indicate the historic colonization routes of the  
467 pathogen. (B) Principal component analysis (PCA) of 284 *Zymoseptoria tritici* isolates, based on  
468 900,193 genome-wide SNPs. (C) PCA of a reduced SNP data set with randomly selected 203 SNPs  
469 matching approximately the number of analyzed TE loci. (D) PCA based on 193 TE insertion loci.

470

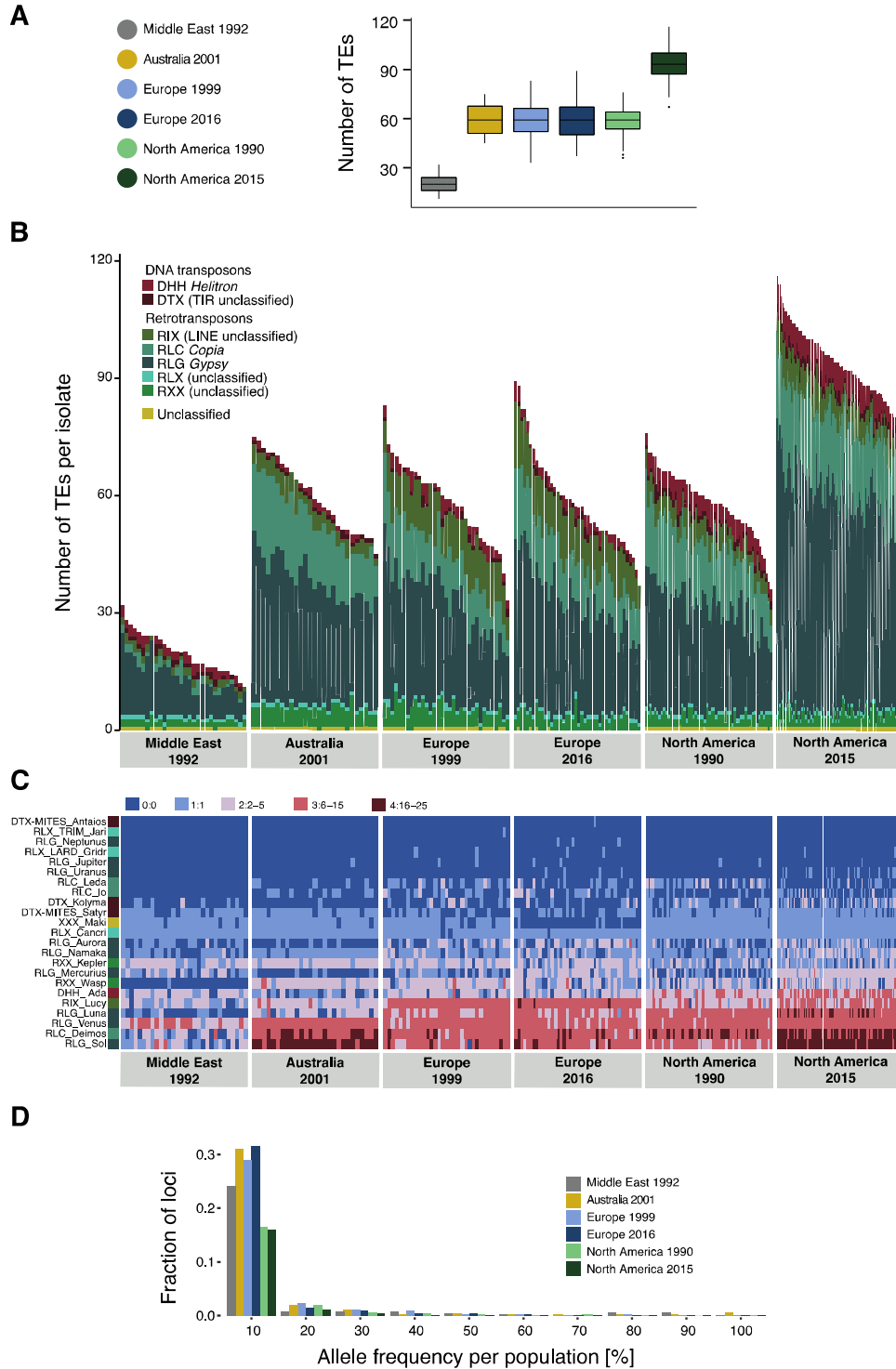
471 Unusual patterns in population differentiation at TE loci suggests that TE activity may drastically  
472 vary across populations (Figure 6A). To analyze this, we first identified the total TE content across  
473 all loci per isolate. We found generally low amounts of TEs in the Middle Eastern population from  
474 Israel (Figure 6B), which is close to the pathogen's center of origin (Stukenbrock *et al.*, 2007).  
475 Populations that underwent at least one migration bottleneck showed a substantial burst of TEs  
476 across all major superfamilies. These populations included the two populations from Europe  
477 (Switzerland) collected in 1999 and 2016 and the North American population from 1990, as well as  
478 the Australian population. We found a second stark increase in TE content in the North American  
479 population sampled in 2015 at the same site as the population from 1990. Strikingly, the isolate with

480 the lowest number of analyzed TEs collected in 2015 was comparable to the isolate with the highest  
481 number of TEs at the same site in 1990. We tested whether sequencing coverage may explain  
482 variation in the detected TEs across isolates, but we found no meaningful association  
483 (Supplementary Figure S3B). We analyzed variation in TE copy numbers across families and found  
484 that the expansions were mostly driven by *Gypsy* elements including the families Luna, Sol and  
485 Venus, the *Copia* family Deimos and the LINE family Lucy (Figure 6C; Supplementary Figures S4-  
486 6). We also found a burst specific to the two North American populations in *Helitron* elements  
487 (Ada), an increase specific to Swiss populations in LINE elements, and an increase in *Copia*  
488 elements in the Australian and the two North American populations. Analyses of complete *Z. tritici*  
489 genomes from the same populations revealed high TE contents in Australia and North America  
490 (Oregon 1990) (Badet et al. 2019). The complete genomes confirmed also that the increase in TEs  
491 was driven by LINE, *Gypsy* and *Copia* families in Australia and *Helitron*, *Gypsy* and *Copia* families  
492 in North America (Badet *et al.*, 2019).

493

494 Finally, we analyzed whether the population-specific expansions were accompanied by shifts in the  
495 allele frequency spectra (Figure 6D). We found that the first step of expansions observed in  
496 Australia and Europe were associated with a downwards shift in allele frequencies. This is consistent  
497 with transposition activity creating new copies in the genomes and stronger purifying selection. In  
498 contrast, the North American populations showed an upwards shift in allele frequencies indicating  
499 relaxation of selection against TEs.

500



501

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

503 (A) The number of transposable elements (TEs) per population. (C) Total TE copies per isolate.

504 Colors stand identify TE superfamilies. (D) TE family copy numbers per isolate. (E) TE insertion

505 frequency spectrum per population.

506



507

## 508 **DISCUSSION**

509 TEs play a crucial role in generating adaptive genetic variation within species but are also drivers of  
510 deleterious genome expansions. We analyzed the interplay of positive selection and incipient  
511 genome expansions using a large-scale population genomics dataset. TEs have substantial  
512 transposition activity in the genome but are strongly counter-selected and maintained at low  
513 frequency. TE dynamics showed distinct trajectories across populations with more recently  
514 established populations having higher TE content in the genome. This strongly suggests that  
515 population specific TE expansions are leading to changes in genome size. In parallel, individual TE  
516 loci possibly contributed to recent local adaptation related to fungicide resistance and host  
517 adaptation.

518

### 519 TRANSPOSITION ACTIVITY IS COUNTER-ACTED BY STRONG PURIFYING SELECTION

520 TE frequencies show a strong skew towards singleton TE insertions across the genome, consistent  
521 with transposition activity creating new insertions and purifying selection maintaining frequencies at  
522 a low level. This is a broadly-known pattern across plants and animals, including *Drosophila*  
523 *melanogaster*, *Zea mays*, *Brachypodium distachyon*, and *Arabidopsis thaliana* (Cridland *et al.*, 2013;  
524 Stuart *et al.*, 2016; Lai *et al.*, 2017; Stritt *et al.*, 2017). TE insertions were under-represented in or  
525 near coding regions, showing that purifying selection acts against TE insertions that disrupt genes.  
526 Coding sequences in the *Z. tritici* genome are densely spaced with an average distance of only ~ 1 kb  
527 (Goodwin *et al.*, 2011). Consistent with this high gene density, TE insertions close to genes peaked  
528 at a distance of 200-400 bp away from coding sequences. A rapid decay in linkage disequilibrium in  
529 the *Z. tritici* populations (Croll *et al.*, 2015; Hartmann *et al.*, 2018) likely contributed to the  
530 efficiency of removing deleterious insertions. The large number of low frequency insertion loci  
531 suggests that at least some TE families are transpositionally active and can create new copies.  
532 Analyses of sequence similarities and transcriptional activity suggest that several TE families are  
533 actively creating new copies in the *Z. tritici* genome (Dhillon *et al.*, 2014; Fouché *et al.*, 2019). The

534 transposition activity in a genome and counter-acting purifying selection is expected to establish an  
535 equilibrium over time (Charlesworth & Charlesworth, 1983). However, changes in population size  
536 due to bottlenecks or founder events are likely to shift the equilibrium. Furthermore, genetic drift  
537 may also impact the prevalence of active TEs or the fixation of mutations contributing to TE control.

538

#### 539 TE INSERTIONS POTENTIALLY UNDERPINNING ADAPTIVE EVOLUTION IN POPULATIONS

540 An emerging theme across kingdoms is that a substantial fraction of adaptive genetic variation in  
541 populations is generated by the insertion of TEs (Chuong *et al.*, 2017). Population genomic datasets  
542 can be used to identify the most likely candidate loci underlying recent adaptation. The shallow  
543 genome-wide differentiation of *Z. tritici* populations provides a powerful background to test for  
544 outlier loci (Hartmann *et al.*, 2018). We focused on two scenarios for an adaptive TE insertion to  
545 arise across populations. We analyzed TE insertions to arise from a globally low frequency to a high  
546 frequency either in the most recent North American population or the most recent European  
547 population. The strongest candidate loci for TE-mediated adaptation were two TE insertions in close  
548 proximity to genes encoding major facilitator superfamily (MFS) transporters. For both loci, the  
549 frequency increase occurred in the North American populations which experienced the first  
550 systematic fungicide applications in the decade prior to the last sampling (Estep *et al.*, 2015). TE-  
551 mediated overexpression of the MFS1 transporter is a known resistance mechanism of *Z. tritici* and  
552 acts by increasing efflux of fungicides out of the cell (Omrane *et al.*, 2017). TE-mediated fungicide  
553 resistance adaptation in the North American population is further supported by a significant  
554 association of levels of fungicide resistance in the population and the presence of the *Gypsy* insertion  
555 near the MFS gene. Furthermore, the locus experienced a selective sweep following the insertion of  
556 the TE. We found that the same TEs experienced genome-wide copy number expansions, suggesting  
557 that the availability of adaptive TE insertions may be a by-product of a TE burst in individual  
558 populations.

559

560 POPULATION-LEVEL TE INVASIONS AND RELAXED SELECTION

561 Across the surveyed populations from four continents, we identified substantial variation in TE  
562 counts per individual. The increase in TEs matches the global colonization history established for *Z.*  
563 *tritici* (Zhan *et al.*, 2003; Stukenbrock *et al.*, 2007). Compared to the Israeli population located  
564 nearest the center of origin in the Middle East, the European populations showed a three-fold  
565 increase in TE counts. The Australian and North American populations established from European  
566 descendants retained high TE counts. We identified a second increase at the North American site  
567 where TE counts nearly doubled again over a 25-year period. Interestingly, the first TE expansion  
568 was caused by a broad increase in copy numbers across the spectrum of TE families. The second  
569 expansion at the North American site was driven by a small number of TE families. Analyses of  
570 completely assembled genomes from the same populations confirmed that genome expansions were  
571 primarily driven by *Gypsy*, *Copia* and *Helitron* superfamilies (Badet *et al.*, 2019). Consistent with  
572 the contributions from individual TEs, we found that the first expansion in Europe led to an increase  
573 in low-frequency variants, suggesting higher transposition activity of many TEs in conjunction with  
574 strong purifying selection. The second expansion at the North American site shifted TE frequencies  
575 upwards, suggesting relaxed selection against TEs. The population-level context of TEs in *Z. tritici*  
576 shows how heterogeneity in TE control interacts with demography to determine extant levels of TE  
577 content and, ultimately, genome size.

578

579 The activity of TEs is controlled by complex selection regimes within species. Actively transposing  
580 elements may accelerate genome evolution and underpin expansions. Hence, genomic defenses  
581 should evolve to efficiently target recently active TEs. Yet, TE-mediated adaptation involves  
582 selection favoring genotypes carrying active TEs, hence counteracting overall negative selection. In  
583 the case of *Z. tritici*, TE expansion activity and counteracting genomic defenses established a  
584 unstable equilibrium across the species range. Furthermore, we show that population subdivisions  
585 are at the origin of highly differentiated TE content within a species. The variability in TE content  
586 emerging over the span of only a few decades or centuries is largely sufficient to explain large

587 genome size expansion across deeper time scales and species. In conclusion, population-level  
588 analyses can recapitulate incipient genome expansions driven by TEs.

589

590

591 **Acknowledgments**

592 We thank Andrea Sánchez Vallet, Anne Roulin and Luzia Stadler for helpful discussions and  
593 comments on previous versions of the manuscript. DC is supported by the Swiss National Science  
594 (grants 31003A\_173265 and IZCOZO\_177052) and the Fondation Pierre Mercier pour la Science.

595

596

597 **REFERENCES**

- 598 **Andrews S, Lindenbaum P, Howard B, Ewels H. 2013.** FastQC: a quality control tool for high  
599 throughput sequence data.
- 600 **Badet T, Oggenfuss U, Abraham L, McDonald BA, Croll D. 2019.** A 19-isolate reference-quality  
601 global pangenome for the fungal wheat pathogen *Zymoseptoria tritici*.
- 602 **Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM,**  
603 **Nikolenko SI, Pham S, Prjibelski AD, et al. 2012.** SPAdes: a new genome assembly algorithm and  
604 its applications to single-cell sequencing. *Journal of computational biology*: a journal of  
605 *computational molecular cell biology* **19**: 455–77.
- 606 **Bao W, Kojima KK, Kohany O. 2015.** Repbase Update, a database of repetitive elements in  
607 eukaryotic genomes. *Mobile DNA* **6**: 4–9.
- 608 **Barnett DW, Garrison EK, Quinlan AR, Strömberg MP, Marth GT. 2011.** Bamtools: A C++  
609 API and toolkit for analyzing and managing BAM files. *Bioinformatics* **27**: 1691–1692.
- 610 **Barron MG, Fiston-Lavier AS, Petrov DA, González J. 2014.** Population Genomics of  
611 Transposable Elements in *Drosophila*. In: Bassler BL, ed. Annual Review of Genetics, Vol 48. 561–  
612 581.
- 613 **Blumenstiel JP, Chen X, He M, Bergman CM. 2014.** An age-of-allele test of neutrality for  
614 transposable element insertions. *Genetics* **196**: 523–538.
- 615 **Bolger AM, Lohse M, Usadel B. 2014.** Trimmomatic: a flexible trimmer for Illumina sequence  
616 data. *Bioinformatics* **30**: 2114–2120.
- 617 **Bonifacino JS, Glick BS. 2004.** The Mechanisms of Vesicle Budding and Fusion. *Cell* **116**: 153–  
618 166.
- 619 **Butelli E, Licciardello C, Zhang Y, Liu J, Mackay S, Bailey P, Reforgiato-Recupero G, Martin**  
620 **C. 2012.** Retrotransposons control fruit-specific, cold-dependent accumulation of anthocyanins in  
621 blood oranges. *Plant Cell* **24**: 1242–1255.
- 622 **Charlesworth B, Charlesworth D. 1983.** THE POPULATION-DYNAMICS OF  
623 TRANSPOSABLE ELEMENTS. *Genetical Research* **42**: 1–27.
- 624 **Chuong EB, Elde NC, Feschotte C. 2017.** Regulatory activities of transposable elements: from  
625 conflicts to benefits. *Nature Reviews Genetics* **18**: 71–86.
- 626 **Cridland JM, Macdonald SJ, Long AD, Thornton KR. 2013.** Abundance and distribution of  
627 transposable elements in two *Drosophila* QTL mapping resources. *Molecular Biology and Evolution*  
628 **30**: 2311–2327.
- 629 **Croll D, Lendenmann MH, Stewart E, McDonald BA. 2015.** The Impact of Recombination  
630 Hotspots on Genome Evolution of a Fungal Plant Pathogen. *Genetics* **201**: 1213-U787.
- 631 **Croll D, Zala M, McDonald BA. 2013.** Breakage-fusion-bridge Cycles and Large Insertions  
632 Contribute to the Rapid Evolution of Accessory Chromosomes in a Fungal Pathogen. *Plos Genetics*  
633 **9**.
- 634 **Darriba D, Taboada GL, Doallo R, Posada D. 2012.** jModelTest 2: more models, new heuristics  
635 and parallel computing. *Nature Methods* **9**: 772.
- 636 **Dhillon B, Gill N, Hamelin RC, Goodwin SB. 2014.** The landscape of transposable elements in the  
637 finished genome of the fungal wheat pathogen *Mycosphaerella graminicola*. *Bmc Genomics* **15**.
- 638 **Dobin A, Davis CA, Schlesinger F, Drenkow J, Zaleski C, Jha S, Gingeras TR, Batut P,**  
639 **Chaisson M. 2012.** STAR: ultrafast universal RNA-seq aligner. *Bioinformatics* **29**: 15–21.
- 640 **Eichler EE, Sankoff D. 2003.** Structural dynamics of eukaryotic chromosome evolution. *Science*  
641 **301**: 793–797.
- 642 **Estep LK, Torriani SFF, Zala M, Anderson NP, Flowers MD, McDonald BA, Mundt CC,**  
643 **Brunner PC. 2015.** Emergence and early evolution of fungicide resistance in North American  
644 populations of *Zymoseptoria tritici*. *Plant Pathology* **64**: 961–971.

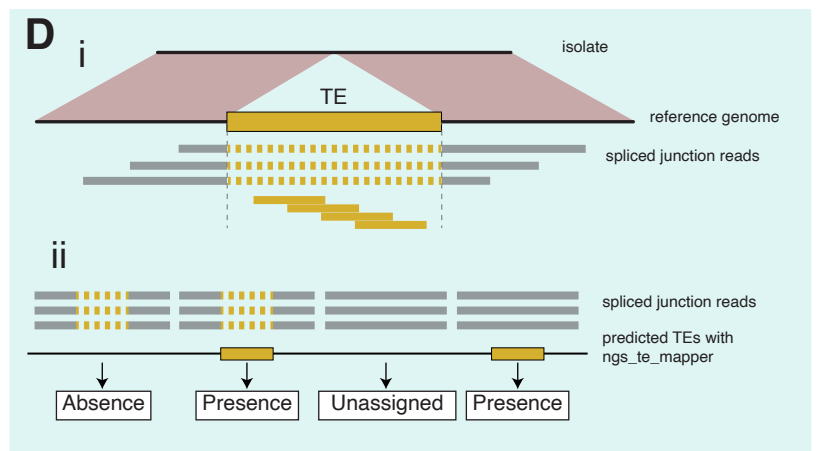
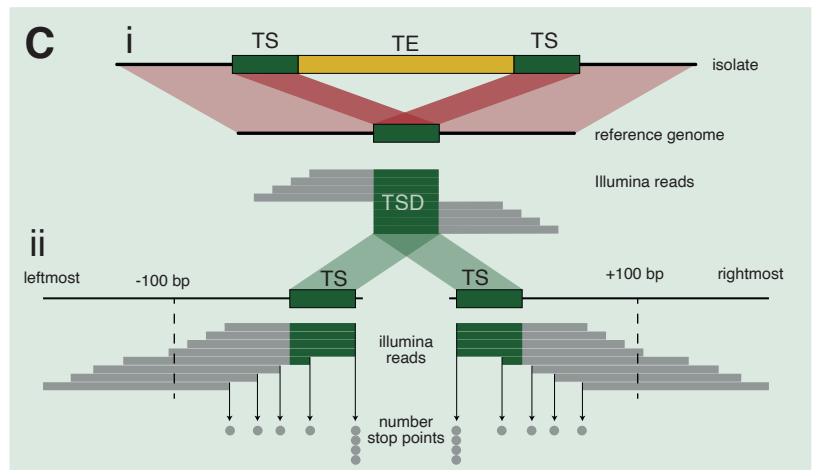
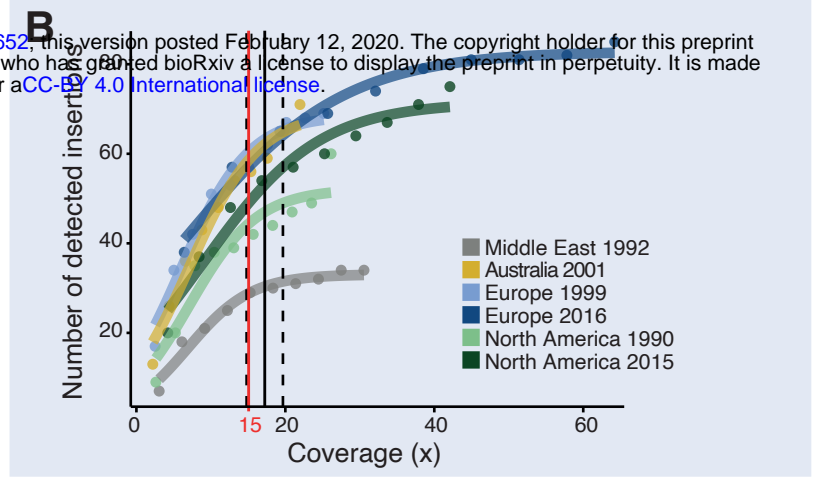
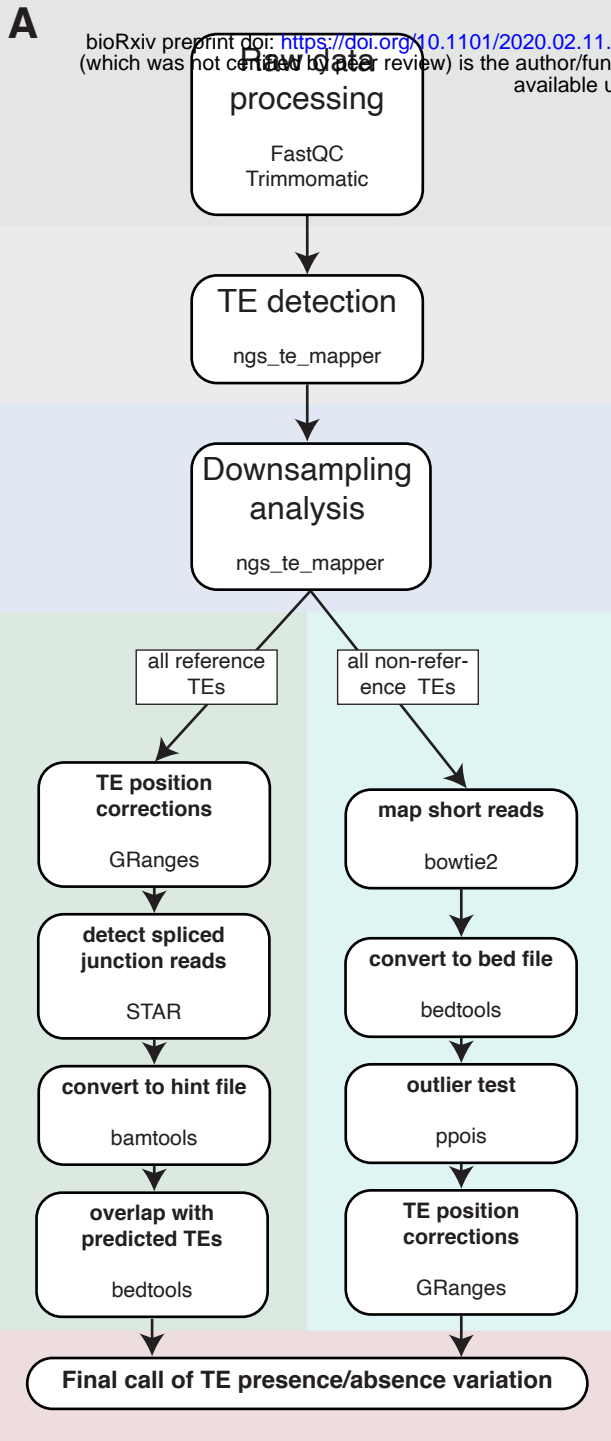
- 645 **Feschotte C. 2008.** Transposable elements and the evolution of regulatory networks. *Nature*  
646 *Reviews Genetics* **9**: 397–405.
- 647 **Fouché S, Badet T, Oggenfuss U, Plissonneau C, Francisco CS, Croll D. 2019.** Stress-driven  
648 transposable element de-repression dynamics in a fungal pathogen. *Molecular Biology and*  
649 *Evolution*.
- 650 **Frantzeskakis L, Kracher B, Kusch S, Yoshikawa-Maekawa M, Bauer S, Pedersen C, Spanu**  
651 **PD, Maekawa T, Schulze-Lefert P, Panstruga R. 2018.** Signatures of host specialization and a  
652 recent transposable element burst in the dynamic one-speed genome of the fungal barley powdery  
653 mildew pathogen. *BMC Genomics* **19**: 1–23.
- 654 **Galagan JE, Selker EU. 2004.** RIP: the evolutionary cost of genome defense. *Trends in Genetics*  
655 **20**: 417–423.
- 656 **Gautier M, Vitalis R. 2012.** Rehh An R package to detect footprints of selection in genome-wide  
657 SNP data from haplotype structure. *Bioinformatics* **28**: 1176–1177.
- 658 **González J, Macpherson JM, Petrov DA. 2009.** A recent adaptive transposable element insertion  
659 near highly conserved developmental loci in *Drosophila melanogaster*. *Molecular Biology and*  
660 *Evolution* **26**: 1949–1961.
- 661 **Goodwin SB, Ben M'Barek S, Dhillon B, Wittenberg AHJ, Crane CF, Hane JK, Foster AJ,**  
662 **Van der Lee TAJ, Grimwood J, Aerts A, et al. 2011.** Finished Genome of the Fungal Wheat  
663 Pathogen *Mycosphaerella graminicola* Reveals Dispensome Structure, Chromosome Plasticity, and  
664 Stealth Pathogenesis. *Plos Genetics* **7**.
- 665 **Goudet J, Jombart T. 2015.** Estimation and Tests of Hierarchical F-Statistics. R package version  
666 0.04-22.
- 667 **Grandaubert J, Bhattacharyya A, Stukenbrock EH. 2015.** RNA-seq-Based Gene Annotation and  
668 Comparative Genomics of Four Fungal Grass Pathogens in the Genus *Zymoseptoria* Identify Novel  
669 Orphan Genes and Species-Specific Invasions of Transposable Elements. *G3-Genes Genomes*  
670 *Genetics* **5**: 1323–1333.
- 671 **Guindon S, Dufayard J-F, Lefort V, Anisimova M, Hordijk W, Gascuel O. 2010.** New  
672 Algorithms and Methods to Estimate Maximum-Likelihood Phylogenies: Assessing the Performance  
673 of PhyML 3.0. *Systematic Biology* **59**: 307–321.
- 674 **Guindon S, Gascuel O. 2003.** A simple, fast, and accurate algorithm to estimate large phylogenies  
675 by maximum likelihood. *Systematic Biology* **52**: 696–704.
- 676 **Guy L, Kultima JR, Andersson SGE. 2010.** GenoPlotR: comparative gene and genome  
677 visualization in R. *Bioinformatics* **26**: 2334–2335.
- 678 **Hartmann F, Croll D. 2017.** Distinct Trajectories of Massive Recent Gene Gains and Losses in  
679 Populations of a Microbial Eukaryotic Pathogen. *Molecular Biology and Evolution*.
- 680 **Hartmann F, McDonald M, Croll D. 2018.** Genome-wide evidence for divergent selection  
681 between populations of a major agricultural pathogen. *Molecular Ecology* **27**: 2725–2741.
- 682 **Hartmann F, Sanchez-Vallet A, McDonald B, Croll D. 2017.** A fungal wheat pathogen evolved  
683 host specialization by extensive chromosomal rearrangements. *ISME J*.
- 684 **Hollister JD, Gaut BS. 2009.** Epigenetic silencing of transposable elements: A trade-off between  
685 reduced transposition and deleterious effects on neighboring gene expression. *Genome Research* **19**:  
686 1419–1428.
- 687 **Huson DH, Richter DC, Rausch C, DeZulian T, Franz M, Rupp R. 2007.** Dendroscope: An  
688 interactive viewer for large phylogenetic trees. *BMC Bioinformatics* **8**: 1–6.
- 689 **Jiao W-B, Schneeberger K. 2019.** Chromosome-level assemblies of multiple *Arabidopsis thaliana*  
690 accessions reveal hotspots of genomic rearrangements. *bioRxiv*: 738880.
- 691 **Jombart T. 2008.** Adegnet: A R package for the multivariate analysis of genetic markers.  
692 *Bioinformatics* **24**: 1403–1405.
- 693 **Jombart T, Ahmed I. 2011.** adegenet 1.3-1: New tools for the analysis of genome-wide SNP data.  
694 *Bioinformatics* **27**: 3070–3071.

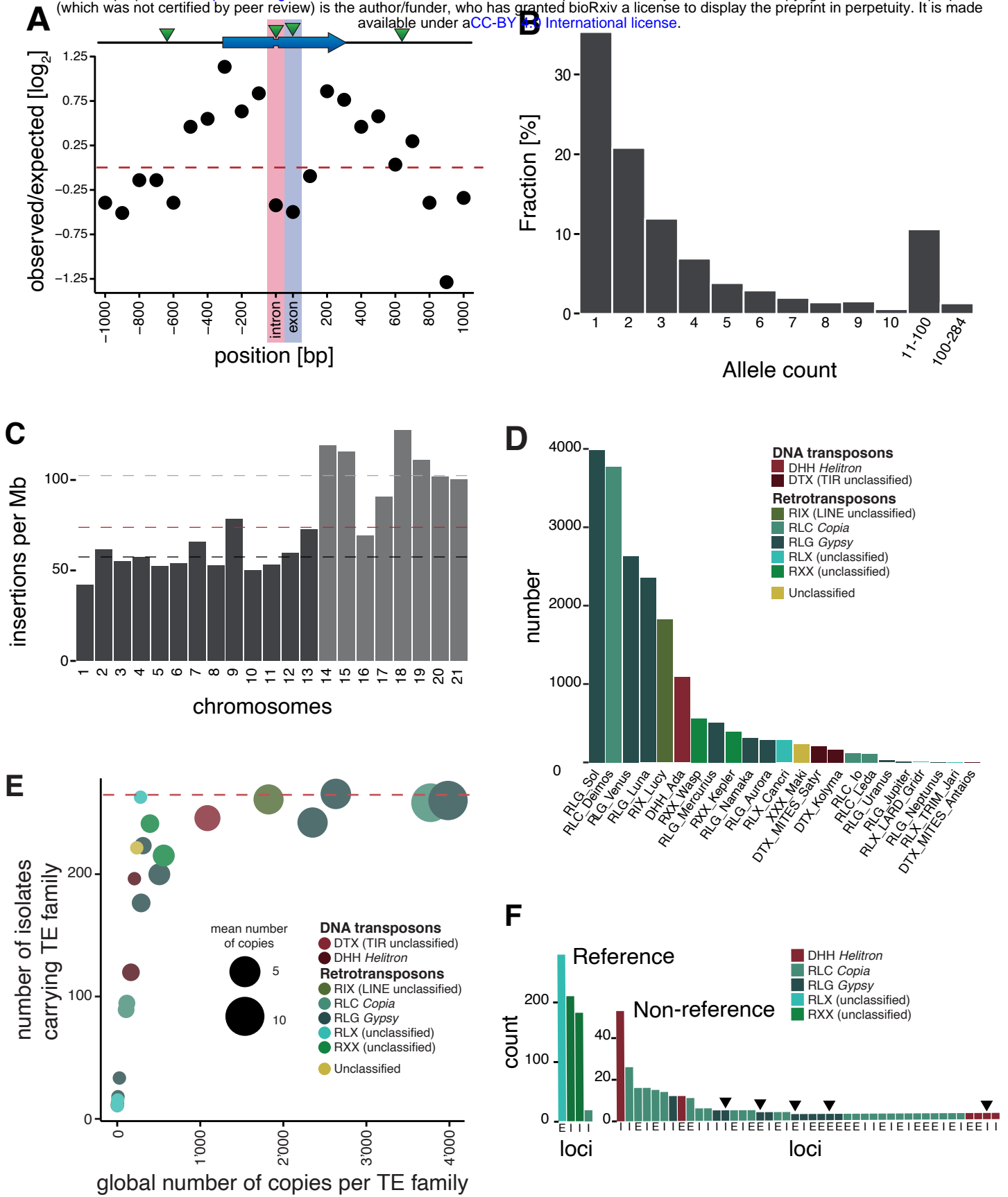
- 695 **Jurka J, Bao W, Kojima KK. 2011.** Families of transposable elements, population structure and  
696 the origin of species. *Biology Direct* **6**: 44.
- 697 **Katoh K, Standley DM. 2013.** MAFFT multiple sequence alignment software version 7:  
698 Improvements in performance and usability. *Molecular Biology and Evolution* **30**: 772–780.
- 699 **Kidwell MG. 2002.** Transposable elements and the evolution of genome size in eukaryotes.  
700 *Genetica* **115**: 49–63.
- 701 **Krishnan P, Meile L, Plissonneau C, Ma X, Hartmann FE, Croll D, McDonald BA, Sánchez-  
702 Vallet A. 2018.** Transposable element insertions shape gene regulation and melanin production in a  
703 fungal pathogen of wheat. *BMC Biology* **16**: 1–18.
- 704 **Lai X, Schnable JC, Liao Z, Xu J, Zhang G, Li C, Hu E, Rong T, Xu Y, Lu Y. 2017.** Genome-  
705 wide characterization of non-reference transposable element insertion polymorphisms reveals  
706 genetic diversity in tropical and temperate maize. *BMC Genomics* **18**: 1–13.
- 707 **Langmead B, Salzberg SL. 2012.** Fast gapped-read alignment with Bowtie 2. *Nature Methods* **9**:  
708 357–359.
- 709 **Lawrence M, Huber W, Pagès H, Aboyoun P, Carlson M, Gentleman R, Morgan MT, Carey  
710 VJ. 2013.** Software for Computing and Annotating Genomic Ranges. *PLoS Computational Biology*  
711 **9**: 1–10.
- 712 **Li H, Durbin R. 2009.** Fast and accurate short read alignment with Burrows-Wheeler transform.  
713 *Bioinformatics* **25**: 1754–1760.
- 714 **Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R.  
715 2009.** The Sequence Alignment/Map format and SAMtools. *Bioinformatics* **25**: 2078–2079.
- 716 **Lim JK. 1988.** Intrachromosomal rearrangements mediated by hobo transposons in *Drosophila*  
717 *melanogaster*. *PNAS* **85**: 9153–9157.
- 718 **Linde CC, Zhan J, McDonald BA. 2002.** Population Structure of *Mycosphaerella graminicola*:  
719 From Lesions to Continents. *Phytopathology* **92**: 946–955.
- 720 **Linheiro RS, Bergman CM. 2012.** Whole Genome Resequencing Reveals Natural Target Site  
721 Preferences of Transposable Elements in *Drosophila melanogaster*. *PLOS ONE* **7**.
- 722 **Lu L, Chen J, Robb SMC, Okumoto Y, Stajich JE, Wessler SR. 2017.** Tracking the genome-  
723 wide outcomes of a transposable element burst over decades of amplification. *Proceedings of the*  
724 *National Academy of Sciences*: 201716459.
- 725 **Lynch M. 2007.** *The Origins of genome architecture*. Sunderland MA: Sinauer Associates.
- 726 **McDonald BA, Mundt CC, Chen R. 1996.** The role of selection on the genetic structure of  
727 pathogen populations: Evidence from field experiments with *Mycosphaerella graminicola* on  
728 wheat. : 73–80.
- 729 **McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernytsky A, Garimella K,  
730 Altshuler D, Gabriel S, Daly M, et al. 2010.** The Genome Analysis Toolkit: A MapReduce  
731 framework for analyzing next-generation DNA sequencing data. *Genome Research* **20**: 1297–1303.
- 732 **Meile L, Croll D, Brunner PC, Plissonneau C, Hartmann FE, McDonald BA, Sánchez-Vallet A.  
733 2018.** A fungal avirulence factor encoded in a highly plastic genomic region triggers partial  
734 resistance to septoria tritici blotch. *New Phytologist*.
- 735 **Nelson MG, Linheiro RS, Bergman CM. 2017.** McClintock: An Integrated Pipeline for Detecting  
736 Transposable Element Insertions in Whole-Genome Shotgun Sequencing Data. *G3 & #58;  
737 Genes/Genomes/Genetics* **7**: 2763–2778.
- 738 **Odorizzi G, Babst M, Emr SD. 2000.** Phosphoinositide signaling and the regulation of membrane  
739 trafficking in yeast. *Trends in Biochemical Sciences* **25**: 229–235.
- 740 **Oliver KR, McComb JA, Greene WK. 2013.** Transposable elements: Powerful contributors to  
741 angiosperm evolution and diversity. *Genome Biology and Evolution* **5**: 1886–1901.
- 742 **Omrane S, Audéon C, Ignace A, Duplaix C, Aouini L, Kema G, Walker A-S, Fillingier S. 2017.**  
743 Plasticity of the MFS1 promoter leads to multi drug resistance in the wheat pathogen *Zymoseptoria*  
744 *tritici*. *mSphere*: 1–42.

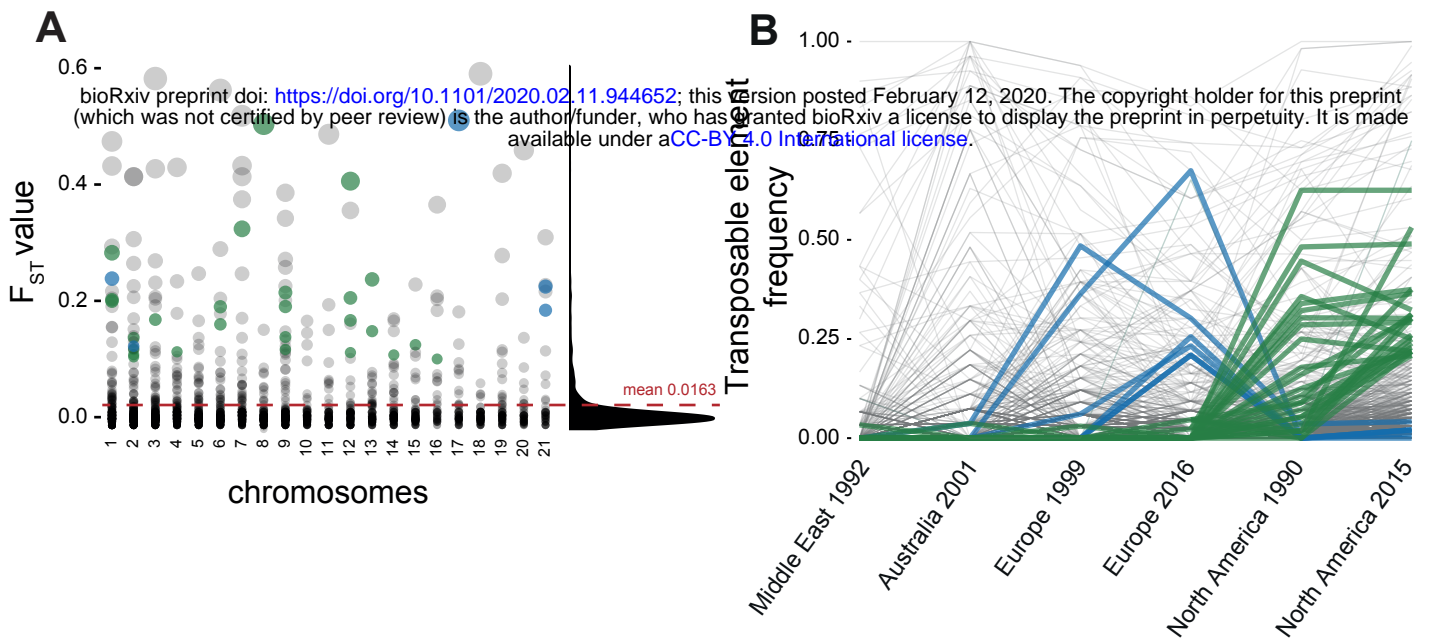
- 745 **Omrane S, Sghyer H, Audeon C, Lanen C, Duplaix C, Walker AS, Fillinger S. 2015.** Fungicide  
746 efflux and the MgMFS1 transporter contribute to the multidrug resistance phenotype in  
747 *Zyoseptoria tritici* field isolates. *Environmental Microbiology* **17**: 2805–2823.
- 748 **Peter M, Kohler A, Ohm RA, Kuo A, Krutzmann J, Morin E, Arend M, Barry KW, Binder M,**  
749 **Choi C, et al. 2016.** Ectomycorrhizal ecology is imprinted in the genome of the dominant symbiotic  
750 fungus *Cenococcum geophilum*. *Nature Communications* **7**.
- 751 **Petrov DA, Aminetzach YT, Davis JC, Bensasson D, Hirsh AE. 2003.** Size matters: Non-LTR  
752 retrotransposable elements and ectopic recombination in *Drosophila*. *Molecular Biology and*  
753 *Evolution* **20**: 880–892.
- 754 **Plissonneau C, Sturchler A, Croll D. 2016.** The Evolution of Orphan Regions in Genomes of a  
755 Fungal Pathogen of Wheat. *Mbio* **7**.
- 756 **Quinlan AR, Hall IM. 2010.** BEDTools: A flexible suite of utilities for comparing genomic  
757 features. *Bioinformatics* **26**: 841–842.
- 758 **R Core Team. 2017.** R: A language and environment for statistical computing. R Foundation for  
759 Statistical Computing, Vienna, Austria.
- 760 **Raffaele S, Kamoun S. 2012.** Genome evolution in filamentous plant pathogens: why bigger can be  
761 better. *Nature Reviews Microbiology* **10**: 417–430.
- 762 **Rice P, Longden L, Bleasby A. 2000.** EMBOSS: The European Molecular Biology Open Software  
763 Suite. *Trends in Genetics* **16**: 276–277.
- 764 **Ritz C, Streibig JC. 2005.** Bioassay analysis using R. *Journal of Statistical Software* **12**: 1–22.
- 765 **Rizzon C, Martin E, Marais G, Duret L, Ségalat L, Biéumont C. 2003.** Patterns of Selection  
766 Against Transposons Inferred from the Distribution of Tc1, Tc3 and Tc5 Insertions in the mut-7  
767 Line of the Nematode *Caenorhabditis elegans*. *Genetics* **165**: 1127–1135.
- 768 **Rouxel T, Grandaubert J, Hane JK, Hoede C, van de Wouw AP, Couloux A, Dominguez V,**  
769 **Anthouard V, Bally P, Bourras S, et al. 2011.** Effector diversification within compartments of the  
770 *Leptosphaeria maculans* genome affected by Repeat-Induced Point mutations. *Nature*  
771 *communications* **2**: 202.
- 772 **Sabeti PC, Varilly P, Fry B, Lohmueller J, Hostetter E, Cotsapas C, Xie X, Byrne EH,**  
773 **McCarroll SA, Gaudet R, et al. 2007.** Genome-wide detection and characterization of positive  
774 selection in human populations. *Nature* **449**: 913–918.
- 775 **SanMiguel P, Gaut BS, Tikhonov A, Nakajima Y, Bennetzen JL. 1998.** The paleontology of  
776 intergene retrotransposons of maize. *Nature Genetics* **20**: 43–45.
- 777 **Shen RM, Batzer MA, Deininger PL. 1991.** Evolution of the master Alu gene(s). *Journal of*  
778 *Molecular Evolution* **33**: 311–320.
- 779 **Slotkin RK, Martienssen R. 2007.** Transposable elements and the epigenetic regulation of the  
780 genome. *Nature Reviews Genetics* **8**: 272–285.
- 781 **Smit A, Hubley R.** RepeatModeler Open-1.0.
- 782 **Stöver BC, Müller KF. 2010.** TreeGraph 2: Combining and visualizing evidence from different  
783 phylogenetic analyses. *BMC Bioinformatics* **11**: 1–9.
- 784 **Stritt C, Gordon SP, Wicker T, Vogel JP, Roulin AC. 2017.** Recent activity in expanding  
785 populations and purifying selection have shaped transposable element landscapes across natural  
786 accessions of the Mediterranean grass *Brachypodium distachyon*. *Genome Biology and Evolution*  
787 **10**: 1–38.
- 788 **Stuart T, Eichten SR, Cahn J, Karpievitch Y V, Borevitz JO, Lister R. 2016.** Population scale  
789 mapping of transposable element diversity reveals links to gene regulation and epigenomic variation.  
790 *Elife* **5**.
- 791 **Stukenbrock EH, Banke S, Javan-Nikkhah M, McDonald BA. 2007.** Origin and domestication of  
792 the fungal wheat pathogen *Mycosphaerella graminicola* via sympatric speciation. *Molecular Biology*  
793 *and Evolution* **24**: 398–411.
- 794 **van't Hof AE, Campagne P, Rigden DJ, Yung CJ, Lingley J, Quail MA, Hall N, Darby AC,**

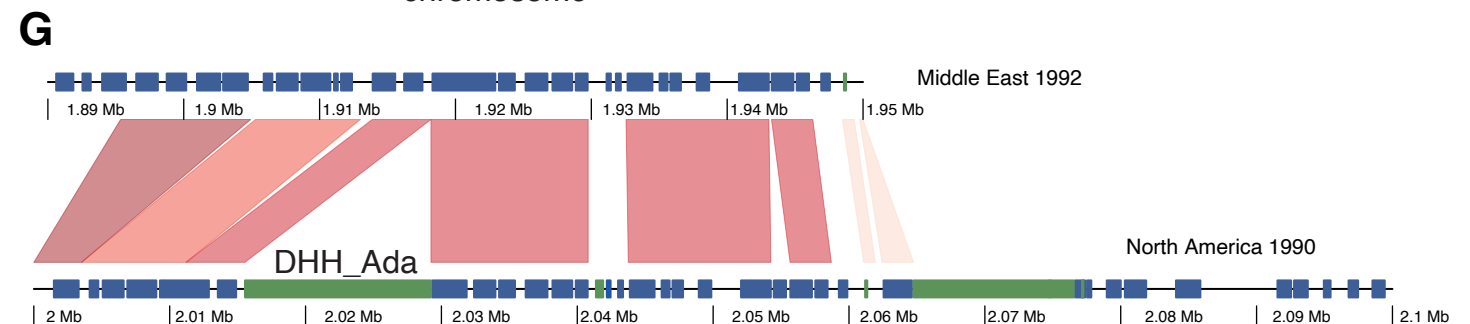
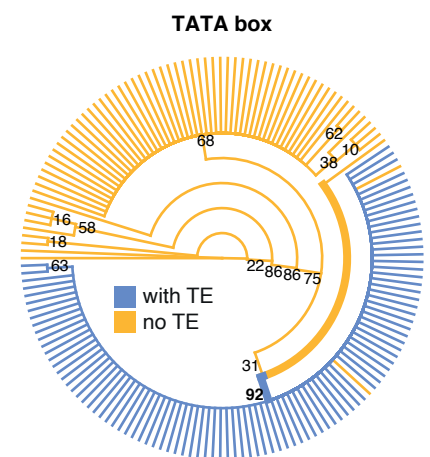
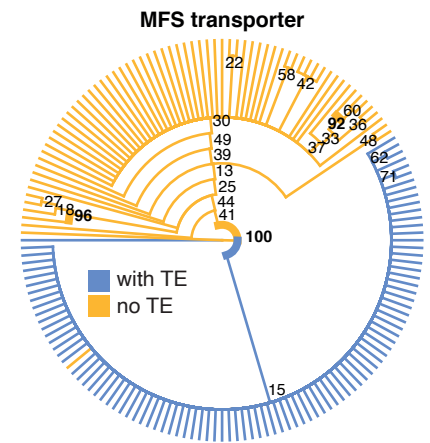
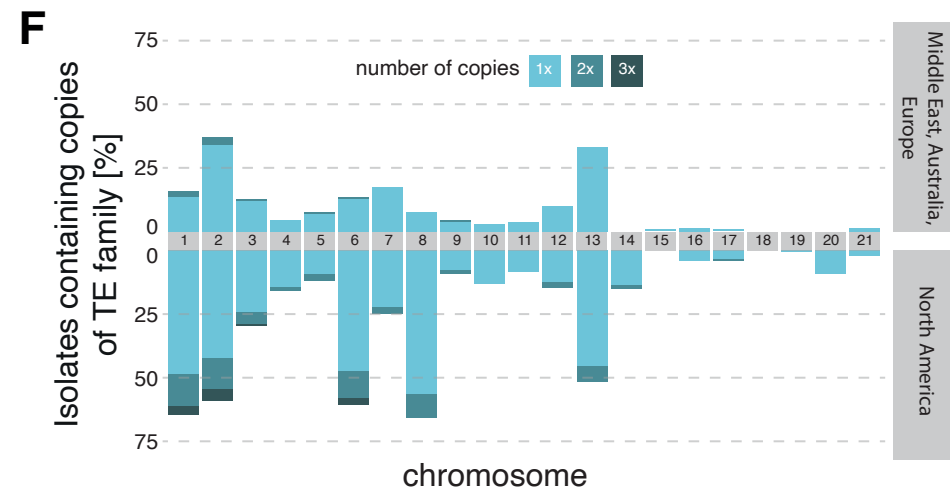
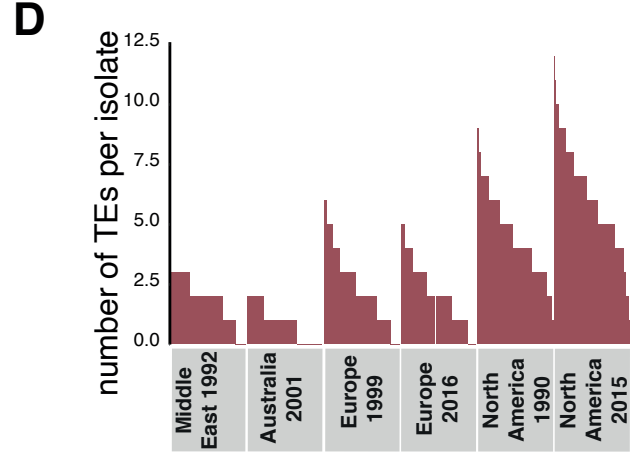
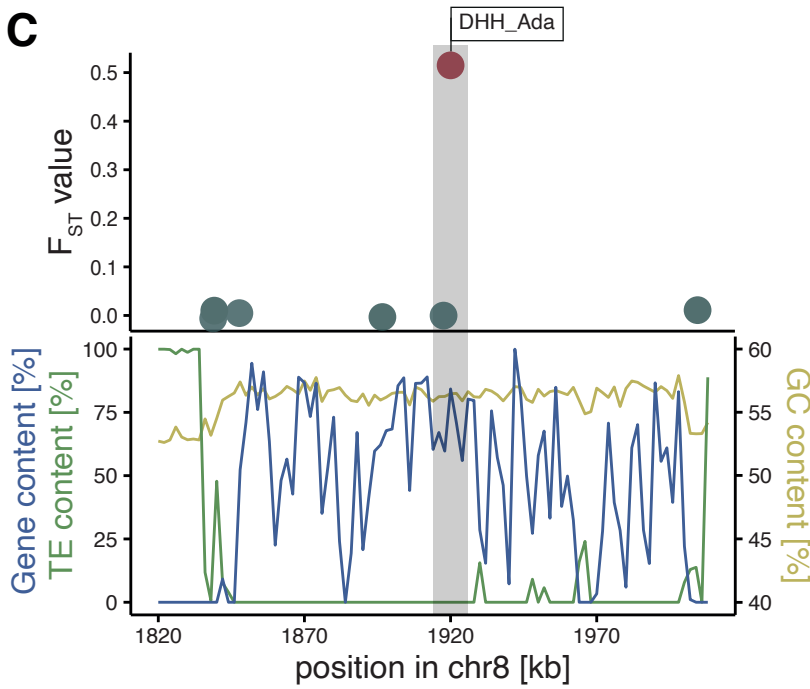
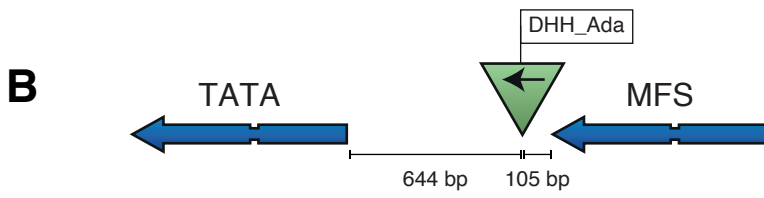
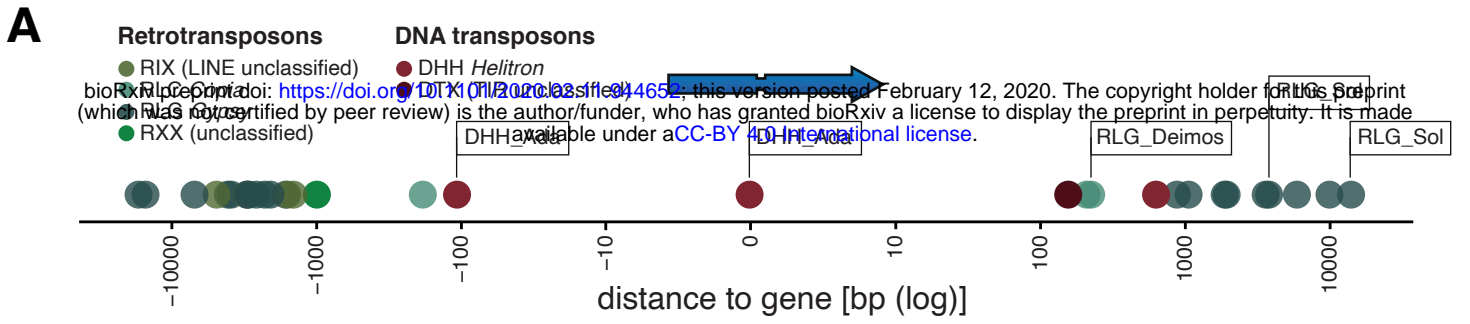


- 795 **Saccheri IJ. 2016.** The industrial melanism mutation in British peppered moths is a transposable  
796 element. *Nature* **534**: 102–105.
- 797 **Walser J-C, Chen B, Feder ME. 2006.** Heat-Shock Promoters: Targets for Evolution by P  
798 Transposable Elements in *Drosophila*. *PLoS Genetics* **2**: e165.
- 799 **Waterhouse AM, Procter JB, Martin DMA, Clamp M, Barton GJ. 2009.** Jalview Version 2-A  
800 multiple sequence alignment editor and analysis workbench. *Bioinformatics* **25**: 1189–1191.
- 801 **Wickham H. 2016.** *ggplot2: Elegant Graphics for Data Analysis*. New York: Springer-Verlag.
- 802 **Wickham H, Chang W. 2016.** devtools: Tools to Make Developing R Packages Easier.
- 803 **Wickham H, Francois R, Henry L, Müller K. 2017.** dplyr: A Grammar of Data Manipulation.
- 804 **Wong WY, Simakov O, Bridge DM, Cartwright P, Bellantuono AJ, Kuhn A, Holstein TW,  
805 David CN, Steele RE, Martínez DE. 2019.** Expansion of a single transposable element family is  
806 associated with genome-size increase and radiation in the genus *Hydra*. *PNAS*: 1–3.
- 807 **Zhan J, Kema GHJ, Waalwijk C, McDonald BA. 2002.** Distribution of mating type alleles in the  
808 wheat pathogen *Mycosphaerella graminicola* over spatial scales from lesions to continents. *Fungal  
809 Genetics and Biology* **36**: 128–136.
- 810 **Zhan J, Linde CC, Jurgens T, Merz U, Steinebrunner F, McDonald BA. 2005.** Variation for  
811 neutral markers is correlated with variation for quantitative traits in the plant pathogenic fungus  
812 *Mycosphaerella graminicola*. *Mol Ecol* **14**: 2683–2693.
- 813 **Zhan J, Pettway RE, McDonald BA. 2003.** The global genetic structure of the wheat pathogen  
814 *Mycosphaerella graminicola* is characterized by high nuclear diversity, low mitochondrial diversity,  
815 regular recombination, and gene flow. *Fungal Genetics and Biology* **38**: 286–297.
- 816 **Zhang L, Hu J, Han X, Li J, Gao Y, Richards CM, Zhang C, Tian Y, Liu G, Gul H, et al. 2019.**  
817 A high-quality apple genome assembly reveals the association of a retrotransposon and red fruit  
818 colour. *Nature Communications* **10**: 1494.
- 819 **Zheng X, Gogarten SM, Lawrence M, Stilp A, Conomos MP, Weir BS, Laurie C, Levine D.  
820 2017.** SeqArray-a storage-efficient high-performance data format for WGS variant calls.  
821 *Bioinformatics* **33**: 2251–2257.
- 822 **Zheng X, Levine D, Shen J, Gogarten SM, Laurie C, Weir BS. 2012.** A high-performance  
823 computing toolset for relatedness and principal component analysis of SNP data. *Bioinformatics* **28**:  
824 3326–3328.
- 825
- 826
- 827
- 828

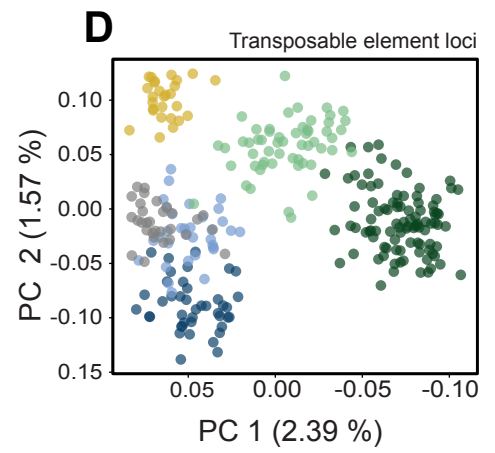
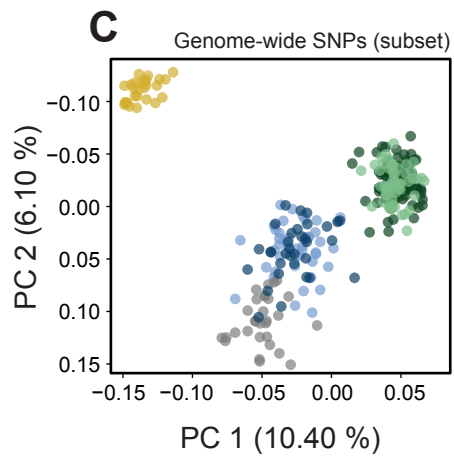
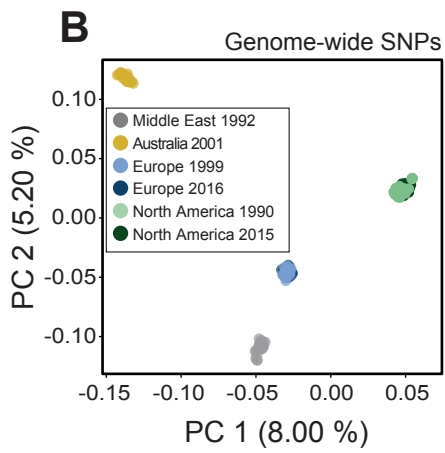
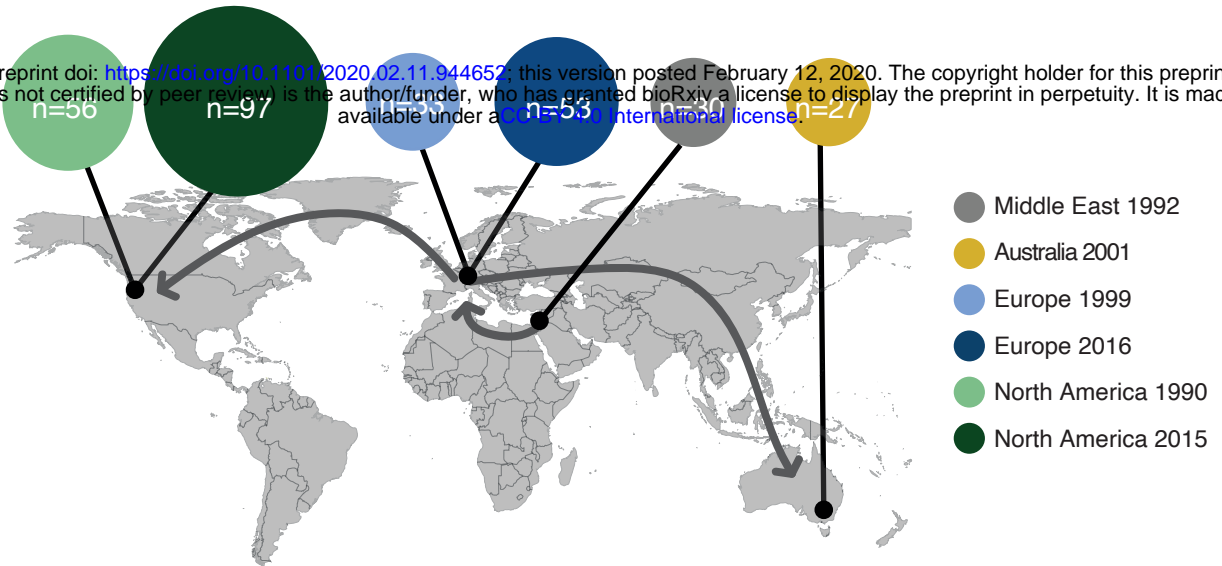


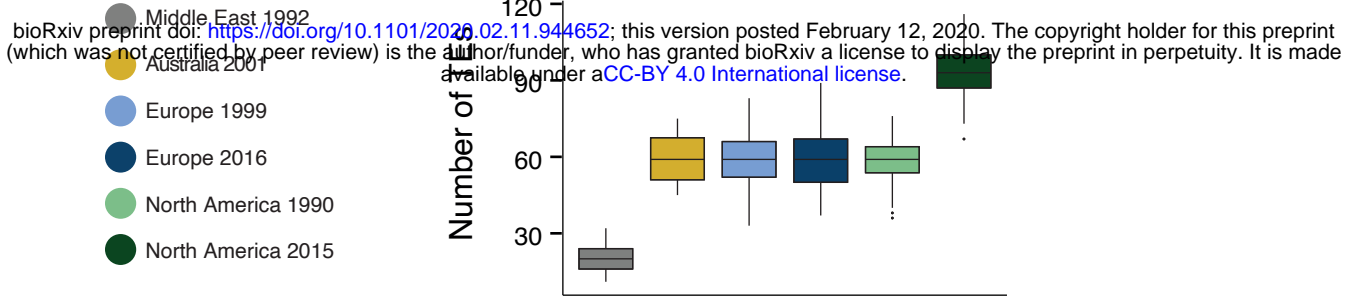
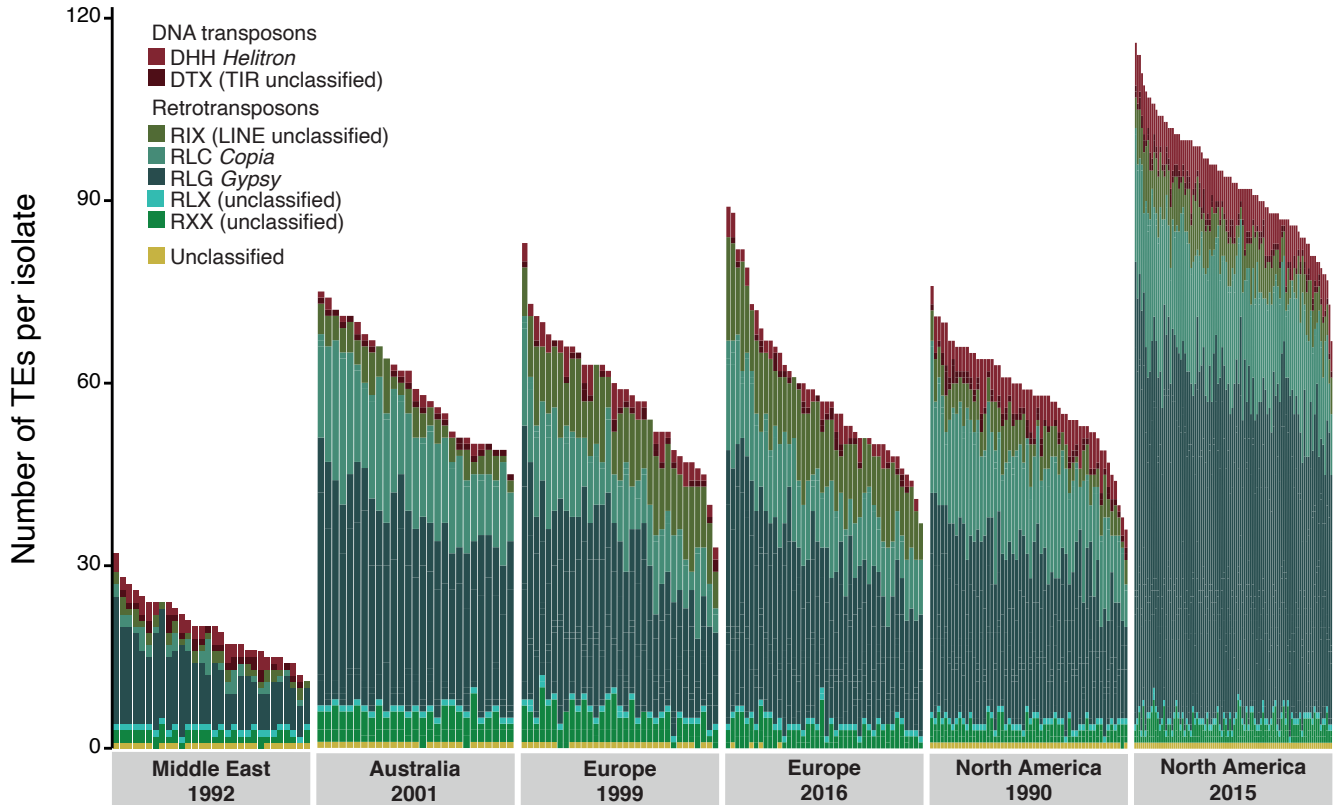
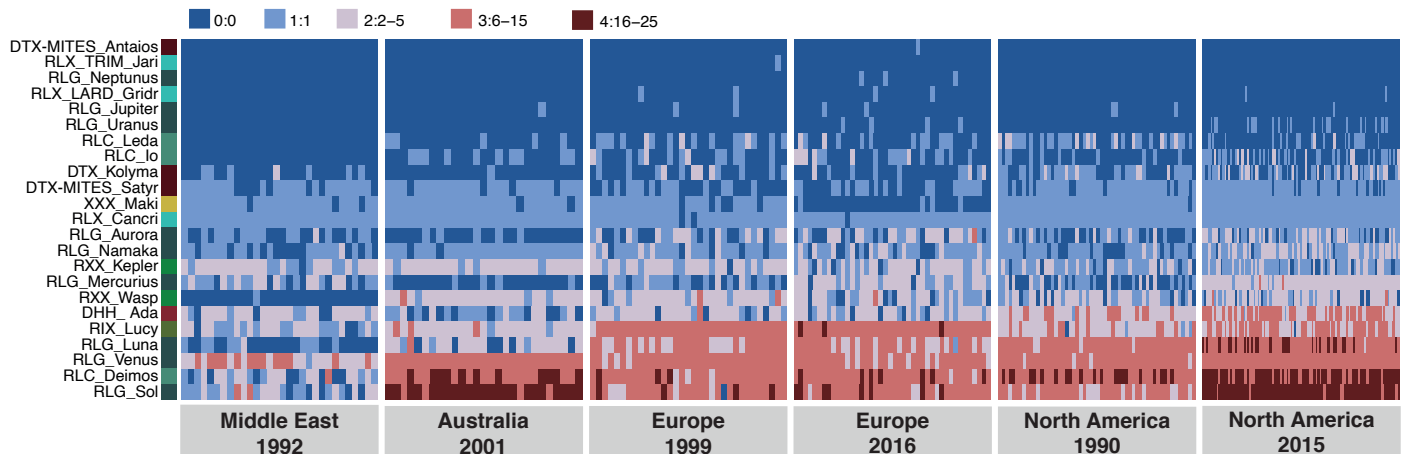






**A** bioRxiv preprint doi: <https://doi.org/10.1101/2020.02.11.944652>; this version posted February 12, 2020. The copyright holder for this preprint (which was not certified by peer review) is the author/funder, who has granted bioRxiv a license to display the preprint in perpetuity. It is made available under aCC-BY-NC-ND 4.0 International license.



**A****B****C****D**