

1 Transposable Elements activity and 2 role in *Meloidogyne incognita* genome 3 dynamics and adaptability

4 Djampa KL KOZLOWSKI, Rahim HASSANALY-GOULAMHOUSSEN, Martine DA-ROCHA,
5 Marc BAILLY-BECHET*, Etienne GJ DANCHIN*.

6

7 * co-last authors

8 Affiliation : Université Côte d'Azur, INRAE, CNRS, ISA, Sophia Antipolis, France

9 Abstract

10 Despite reproducing without sexual recombination, the root-knot nematode *Meloidogyne*
11 *incognita* is adaptive and versatile. Indeed, this species displays a global distribution, is able
12 to parasitize a large range of plants and can overcome plant resistance in a few generations.
13 The mechanisms underlying this adaptability without sex remain poorly known and only low
14 variation at the single nucleotide polymorphism level have been observed so far across
15 different geographical isolates with distinct ranges of compatible hosts. Hence, other
16 mechanisms than the accumulation of point mutations are probably involved in the genomic
17 dynamics and plasticity necessary for adaptability. Transposable elements (TEs), by their
18 repetitive nature and mobility, can passively and actively impact the genome dynamics. This
19 is particularly expected in polyploid hybrid genomes such as the one of *M. incognita*. Here,
20 we have annotated the TE content of *M. incognita*, analyzed the statistical properties of this
21 TE content, and used population genomics approach to estimate the mobility of these TEs
22 across 12 geographical isolates, presenting phenotypic variations. The TE content is more
23 abundant in DNA transposons and the distribution of sequence identity of TE occurrences to
24 their consensus suggest they have been at least recently active. We have identified loci in
25 the genome where the frequencies of presence of a TE showed variations across the
26 different isolates. Compared to the *M. incognita* reference genome, we detected the insertion
27 of some TEs either within coding regions or in the upstream regulatory regions. These
28 predicted TEs insertions might thus have a functional impact. We validated by PCR the
29 insertion of some of these TEs, confirming TE movements probably play a role in the
30 genome plasticity with possible functional impacts.

31 Introduction

32 Agricultural pests cause substantial yield loss to the worldwide life-sustaining production
33 (Savary et al. 2019) and threaten the survival of different communities in developing
34 countries. With a constantly growing human population, it becomes more and more crucial to
35 reduce the loss caused by these pests while limiting the impact on the environment. In this
36 context, understanding how pests evolve and adapt both to the control methods deployed
37 against them and to a changing environment is essential. In metazoa, nematodes and
38 insects are the most destructive agricultural pests. Nematodes alone are responsible for

39 crop yield losses of ca. 11% which represents up to 100 billion € economic loss annually
40 (Agrios 2005; McCarter 2009). The most problematic nematodes to worldwide agriculture
41 belong to the genus *Meloidogyne* (Jones et al. 2013) and are commonly named root-knot
42 nematodes (RKN) owing to the gall symptoms their infection leaves on the roots. Curiously,
43 the RKN species showing the wider geographical distribution and infecting the broadest
44 diversity of plants reproduce asexually via mitotic parthenogenesis (Trudgill and Blok 2001;
45 Castagnone-Sereno and Danchin 2014). In the absence of sexual recombination, the
46 genomes are supposed to irreversibly accumulate deleterious mutations, the efficiency of
47 selection is reduced due to linkage between conflicting alleles while the combination of
48 beneficial alleles from different individuals is impossible (Muller 1964; Hill and Robertson
49 1966; Kondrashov 1988; Glémin and Galtier 2012). For these reasons asexual reproduction
50 is considered an evolutionary dead end and is actually quite rare in animals (Rice 2002). In
51 this perspective, the parasitic success of the parthenogenetic RKN might represent an
52 evolutionary paradox.

53 Previous comparative genomics analyses have shown the genomes of the most devastating
54 RKN are polyploid as a result of hybridization events (Blanc-Mathieu et al. 2017; Szitenberg
55 et al. 2017). In the model RKN *M. incognita*, the resulting gene copies not only diverge at the
56 nucleotide level but also in their expression patterns, suggesting this peculiar genome
57 structure could support a diversity of functions and might be involved in the parasitic success
58 despite the absence of sexual reproduction (Blanc-Mathieu et al. 2017). This hypothesis
59 seems consistent with the ‘general purpose genotype’ concept, which proposes successful
60 parthenogens have a generalist genotype with good fitness in a variety of environments
61 (Vrijenhoek and Parker 2009). An alternative non mutually exclusive hypothesis is the
62 ‘frozen niche variation’ concept which proposes parthenogens are successful in stable
63 environments because they have a frozen genotype adapted to this specific environment
64 (Vrijenhoek and Parker 2009). Interestingly, the frequency of parthenogenetic invertebrates
65 is higher in agricultural pests, probably because the anthropized environments in which they
66 live are more stable and uniform (Hoffmann et al. 2008).

67 However, although a general purpose genotype brought by hybridization might partially
68 explain the wide host range and geographical distribution of these parthenogenetic RKNs,
69 this alone, cannot explain how these species evolve and adapt to new hosts or environments
70 without sex. For instance, initially, avirulent populations of some of these RKN, controlled by
71 a resistance gene in a tomato, are able to overcome the plant resistance in a few
72 generations, leading to virulent sub-populations, in controlled laboratory experiments
73 (Castagnone-Sereno et al. 1994; Castagnone-Sereno 2006). Emergence of virulent
74 populations, not controlled anymore by resistance genes have also been reported in the field
75 (Barbary et al. 2015).

76 The mechanisms underlying the adaptability of parthenogenetic RKN without sex remain
77 elusive. Recent population genomics analyses showed that only a few single nucleotide
78 variations (SNV) could be identified by comparing different Brazilian *M. incognita* isolates
79 showing distinct ranges of host compatibility (ie host races) (Koutsovoulos et al. 2020).
80 Addition of further isolates from different geographical locations across the world did not
81 substantially expand the number of variable positions in the genome. Furthermore, the few
82 identified SNV did not show correlation with either the geographical location, the host range
83 or the current crop species. However, these SNV could be used as markers to confirm the
84 absence of sexual meiotic recombination in *M. incognita*. Thus, the low nucleotide variability
85 that was observed between isolates is probably not the main player in the genomic plasticity
86 underlying the adaptability of *M. incognita*.

87 Consistent with these views, convergent gene copy number variations were observed
88 following resistance breaking down by two originally avirulent populations of *M. incognita*
89 from distinct geographic origins (Castagnone-Sereno et al. 2019). The mechanisms
90 supporting these gene copy numbers and other genomic variations possibly involved in the
91 adaptive evolution of *M. incognita* remain to be described.

92 Transposable elements (TEs), by their repetitive and mobile nature, can both passively and
93 actively impact genome plasticity and stability. Being repetitive, they can be involved in
94 illegitimate genomic rearrangements leading to loss of genomic portions or expansion of
95 gene copy numbers. Being mobile, they can insert in coding or regulatory regions and have
96 a functional impact on the gene expression or gene structure itself. In some fungal
97 phytopathogens, TEs are a major player of adaptive genome evolution by both passively and
98 actively impacting the genome structure and sequence (Faino et al. 2016). In parallel,
99 although TE movements can provide beneficial 'novelty' / plasticity, their uncontrolled activity
100 can also be highly detrimental and put the organism at risk. Whether TEs also play an
101 important role in the adaptive evolution of animal genomes and particularly in parasites,
102 engaged in a continuous 'arms race' with their hosts, remains poorly known. According to
103 the Red Queen hypothesis, host-parasites arms race is a major justification for the
104 prevalence of otherwise costly sexual reproduction (Lively 2010) and, in the absence of sex,
105 other mechanisms should provide the necessary plasticity to sustain this arms race.

106 From an evolutionary point of view, the parthenogenetic root-knot nematode *M. incognita*
107 represents an interesting model to study the activity of TEs and their impact. Indeed, being a
108 plant parasite, *M. incognita* is engaged in an arms race with the plant defence systems and
109 point mutations alone are not expected to be a major mechanism supporting adaptation in
110 this species (Koutsovoulos et al. 2020).

111 In this study, we have tested whether the TE activity could represent a mechanism
112 supporting genome plasticity and eventually adaptive evolution in *M. incognita*. We have re-
113 annotated the 185Mb triploid genome of *M. incognita* (Blanc-Mathieu et al. 2017) for TEs,
114 using stringent filters to identify canonical TEs, possibly active in the genome. We analyzed
115 the statistical properties of the TE content and the distribution of TE sequence identity levels
116 to their consensus was skewed towards high values, suggesting they might have
117 undergone recent multiplications in the genome. We have then tested whether the
118 frequencies of presence/absence of these TEs across the genome varied between different
119 isolates. To test for variations in frequencies, we have used population genomics data from
120 eleven *M. incognita* isolates collected on different crops and locations and differing in their
121 ranges of compatible hosts (Koutsovoulos et al. 2020). From the set of TE loci that
122 presented the most contrasted patterns of presence/absence across the isolates, we
123 investigated whether some could represent neo-insertions. To estimate the potential impact
124 of TE insertions, we checked whether some were inserted within coding or possible
125 regulatory regions. Finally, we validated some of the neo-insertions, predicted by population
126 genomics data by PCR assays. Overall, our study represents the first estimation of TE
127 activity as a mechanism possibly involved in the genome plasticity and the associated
128 functional impact in the most devastating nematode to worldwide agriculture. Because this
129 study focuses on an allopolyploid and parthenogenetic animal species, it also opens new
130 evolutionary perspectives on the fate and potential adaptive impact of TEs in these singular
131 organisms.

132

133 Results

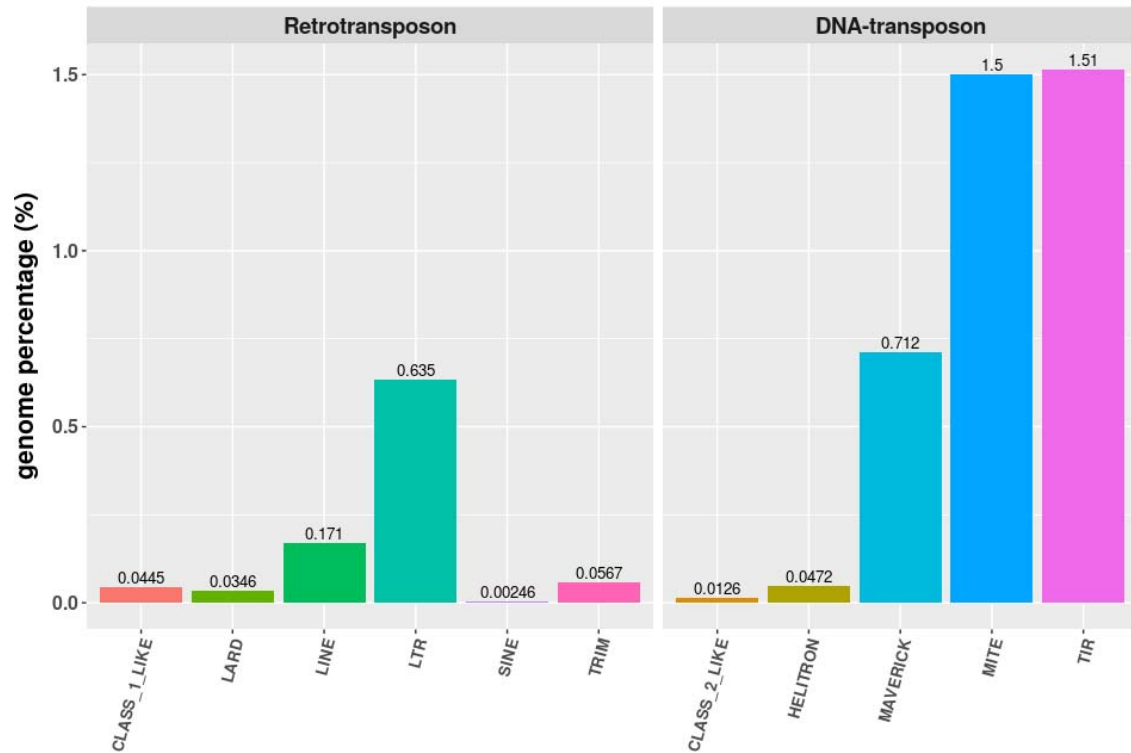
134 The *M. incognita* TE landscape is diversified but mostly
135 composed of DNA transposons.

136 We used the REPET pipeline (Quesneville et al. 2005; Flutre et al. 2011) to predict and
137 annotate the *M. incognita* repeatome (cf. methods). Here, we define the repeatome as all the
138 repeated sequences in the genome, excluding Simple Sequence Repeats (SSR or
139 microsatellites). The repeatome spans 26.38 % of the *M. incognita* genome length. As we
140 wanted to assess whether TEs actively contributed to genomic plasticity, we applied a series
141 of stringent filters on the whole repeatome to retain only repetitive elements presenting
142 canonical signatures of TEs (cf. methods). We identified 480 different TE-consensus
143 sequences that allowed annotation of 9,702 canonical TE, spanning 4.73% of the genome.
144 Both retro (Class I) and DNA (Class II) transposons (Wicker et al. 2007) compose the *M.*
145 *incognita* TE landscape with 5/7 and 4/5 of the known TE orders represented respectively,
146 showing a great diversity of elements (Fig 1). Retro-transposons and DNA-transposons
147 respectively cover 0.94 and 3.78 % of the genome. TIR (Terminal Inverted Repeats) and
148 MITEs (Miniature Inverted repeat Transposable Elements) DNA-transposons alone
149 represent almost two-thirds of the *M. incognita* TE content (63.64 %). Hence, the *M.*
150 *incognita* TE landscape is diversified but mostly composed of DNA-transposons.

151 As a technical validation of our annotation method, we used the same protocol to predict the
152 *C. elegans* TE genomic content (sup. Fig 1 & sup. Table 1), using the PRJNA13758
153 assembly (The *C. elegans* Genome Sequencing Consortium 1998), and compared our
154 results to the reference report of the TE landscape in this model nematode (Bessereau
155 2006). We estimated that the *C. elegans* repeatome spans 11.81% of its genome, which is
156 close to the 12 % described in (Bessereau 2006). The same resource also reported that
157 MITEs and LTR respectively compose ~2% and 0.4% of the *C. elegans* genomes while we
158 predicted 1.8% and 0.2%. Predictions obtained using our protocol are thus in the range of
159 previous predictions for *C. elegans*. In the same study, it was mentioned that most of *C.*
160 *elegans* TE sequences "are fossil remnants that are no longer mobile", and that active TEs
161 are DNA transposons. This suggests a stringent filtering process is necessary to isolate TEs
162 that are the most likely to be active (e.g. the 'canonical' ones). Using the same post-
163 processing protocol as for *M. incognita*, we estimated that canonical TEs span 3.99% of the
164 *C. elegans* genome and DNA-transposon alone span 3.13% of this genome.

165 Overall, the similarity of our results with the previous reports in *C. elegans* suggests our
166 filtered annotation of the TE content of *M. incognita* represents an accurate picture of the
167 potentially active TE landscape.

168



169
170
171
172
173
174
175
176
177

Fig 1: Canonical TE annotations distribution in *M. incognita* genome

CLASS_1_LIKE and CLASS_2_LIKE elements present sufficient evidence to class them as Retro or DNA transposons respectively, but not enough to further assign them to an order. Genome percentage is based on a *M. incognita* genome size of 183,531,997 bp

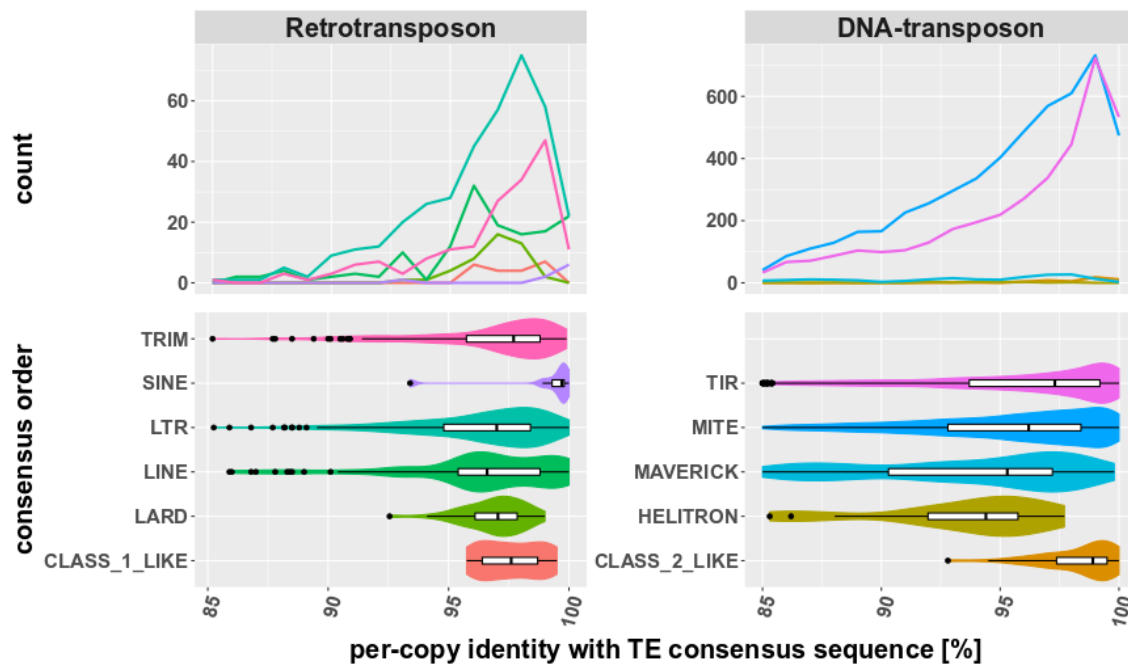
Table 1: Per-order summary of *M. incognita* TE annotations.

	order	nb. of features	total length (bp)	genome percentage (%)	median length (bp)	median of median identity with consensus (%)
Retro-transposon	CLASS_1_LIKE	21	81,609	0.044	3,871	97.6
	LARD	45	63,42	0.035	1,433	97.05
	LINE	145	313,224	0.171	1,971	96.6
	LTR	373	1,164,836	0.635	2,415	97
	SINE	9	4,522	0.002	528	99.7
	TRIM	174	104,018	0.057	525	97.7
Total		767	1,731,629	0.944		

DNA-transposon	CLASS_2_LIKE	48	23,132	0.013	508.5	98.9
	HELITRON	18	86,666	0.047	5,080	94.4
	MAVERICK	189	1,307,068	0.712	6,224	95.3
	MITE	5085	2,755,381	1.501	525	96.2
	TIR	3595	2,777,270	1.513	737	97.3
Total		8935	6,949,517	3.787		

178 Canonical TE annotations are highly identical to their
179 consensus sequences and some present evidence for
180 transposition machinery.

181 Canonical TE annotations have a median nucleotide identity of 97.12% with their respective
182 consensus sequences, but the distribution of identity values varies between TE orders (Fig
183 2, Table 1). Most of the TEs within an order share a high identity level with their
184 consensus, except HELITRON and MAVERICK. Even considering our inclusion threshold
185 at minimum 85% identity (cf methods), the overall distribution of average % identities peaks
186 at high values. DNA transposons show a wider dispersion in identity to their consensus
187 than retrotransposons, as indicated by the elongated boxes in Fig 2B and this distribution
188 peaks at higher values. As a consequence, it is among those elements that annotations
189 sharing a higher similarity to their consensus sequences are found. In particular among TIR,
190 one fourth (Fig 2; sup. Table 2) of the annotations share above 99% identity with their
191 consensus. SINE and CLASS_2_LIKE have similar profiles but are present in very low
192 numbers.
193



194 **Fig2: per-copy identity rate with consensus**
195

196 Top frequency plots show the distribution of TE copies count per order in function of the
197 identity % they share with their consensus sequence. To facilitate inter-orders comparison,
198 bottom violin plots display the same information as a density curve, but also encompass
199 boxplots. Each colour is specific to a TE order.

200
201

202 Higher identity of TE annotations to their consensus can be considered a proxy of their
203 recent activity (Bast et al. 2015; Lerat et al. 2019). To further investigate whether some TEs
204 might be (or have been recently) active, we searched for the presence of genes involved in
205 the transposition machinery within *M. incognita* canonical TEs (cf methods). Among the
206 canonical TE annotations, 6.26% (607/9,702) contain at least one predicted protein-coding
207 gene, with a total of 907 genes. Of these 907 genes, 328 code for proteins with at least one
208 conserved domain known to be related to transposition machinery. We found that 35.37%
209 (116/328) of the transposition machinery genes had substantial expression support from
210 RNA-seq data. In total, 111 canonical TE-annotations contain at least one substantially
211 expressed transposition machinery gene (see supplementary material 1). These 111 TE
212 annotations correspond to 41 different TE-consensuses, and as expected, only consensuses
213 from the autonomous TE orders, e.g. LTRs, LINEs, TIRs, HELITRON, and MAVERICKs
214 present TE-copies with substantially expressed genes coding for transposition machinery.
215 Also as expected, the non-autonomous TEs do not contain any transposition machinery
216 gene at all. This suggests that some of the detected TEs might have functional transposition
217 machinery, which in turn could be hijacked by the non-autonomous elements.

218 Overall, the presence of a substantial proportion of TE annotations highly similar to their
219 consensuses combined with the presence of genes coding for the transposition machinery
220 and supported by expression data suggest some TE might be active in the genome of *M.*
221 *incognita*.

222

223 Thousands of loci show variations in TE presence frequencies 224 across *M. incognita* isolates.

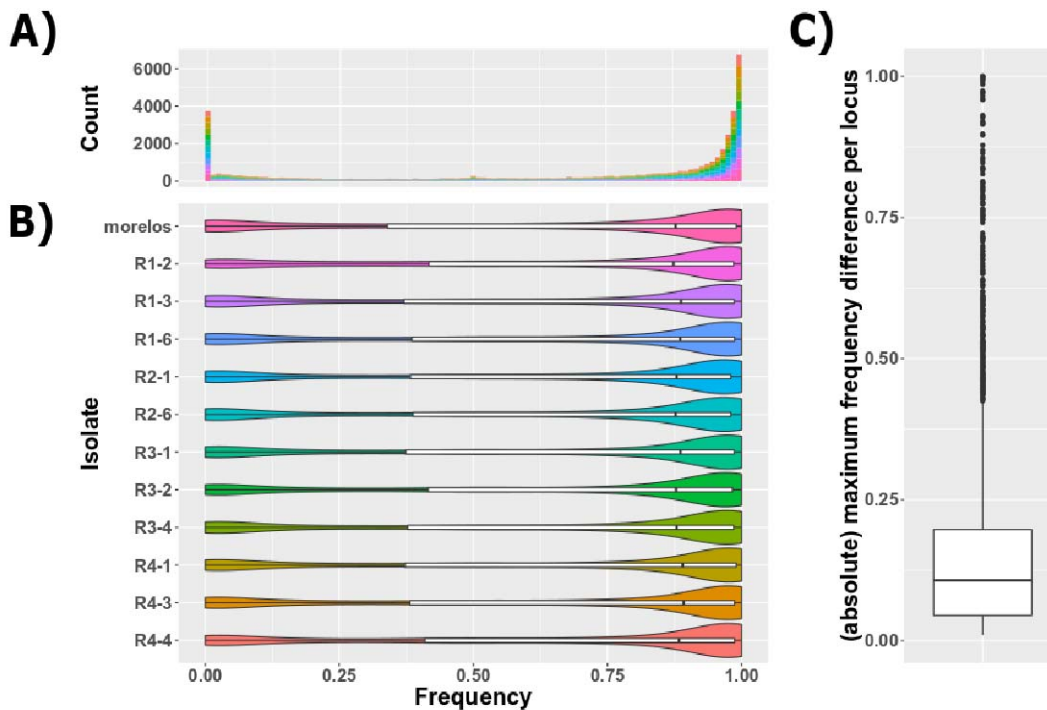
225 We used the PopoolationTE2 (Kofler et al. 2016) pipeline (joint algorithm) on the *M.*
226 *incognita* reference genome (Blanc-Mathieu et al. 2017) and the canonical TE annotation to
227 detect variations in TE frequencies across the genome between 12 geographical isolates (cf.
228 methods). One isolate comes from Morelos in Mexico, which is the isolate that was used to
229 produce the *M. incognita* reference genome. The 11 other isolates come from different
230 locations across Brazil, and present different ranges of compatible hosts (referred to as R1,
231 R2, R3, R4, see sup. Fig 2) currently infected crop species (Koutsovoulos et al. 2020). Pool-
232 seq paired-end Illumina data has been generated for all these isolates. For each locus, each
233 isolate has an associated frequency value representing the proportion of individuals in the
234 pool having the TE detected at this location.

235 We identified 3,524 loci where the frequency variation between at least two isolates was
236 higher than our estimated PopoolationTE2 error rate (0.00972 i.e less than 1%, see
237 methods).

238 Overall, the distribution of frequencies is bimodal (Fig 3-A), and this pattern is common to all
239 the isolates, including the reference Morelos isolate (Fig 3-B). On average per isolate, 21.1%
240 of the loci have frequencies < 25%, 60.7% have frequencies > 75%, and only 18.2% show
241 intermediate frequencies Hence, in every isolate, most of the TE frequency values pack
242 around extreme values e.g. <25% or >75%.

243 Nevertheless, these statistics provide no information about the frequency variability between
244 isolates for a given locus. To address this question, for each locus, we computed the
245 absolute maximum frequency difference between isolates (Fig 3-C). We found that the
246 maximum frequency variation across the isolates is smaller than 20% in 75% of the loci
247 (2,643/3,524). Hence, most of the loci show little to moderate variations in frequencies
248 between isolates. Combined to the previous result, this implies that for most loci, the TEs are
249 present either at a high or a low frequency among all isolates. However, some TE loci show
250 more contrasted variations and will be the focus of further studies in our pipeline.

251



252

253 **Fig. 3: TE frequency distribution.**

254 The histogram (A) and violin plot (B) represent the TE frequency distribution per isolate. The
255 colour chart is identical between the two figures. Both representations reveal that in all the
256 isolates, only a few TE are found with intermediate frequencies. Right boxplot (C) represents
257 the frequency absolute maximum difference per locus. For a given locus, it illustrates the
258 frequency variability between isolates. The higher is the value; the more important is the
259 frequency difference between at least two isolates. A value of 1 implies that the TE is absent
260 in at least one isolate while it is present in 100% of the individuals of at least another isolate.

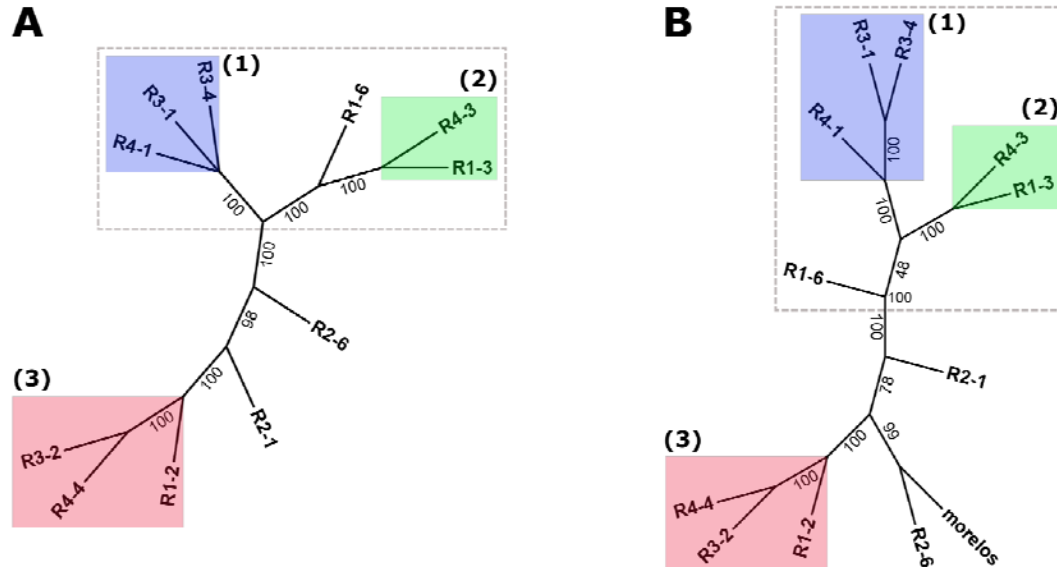
261

262 Variations of TE frequencies across isolates recapitulate their
263 divergence at the sequence level

264 We performed a Neighbour-joining phylogenetic analysis of *M. incognita* isolates based on a
265 distance matrix constructed from TE frequencies (cf methods). We then compared the

266 obtained phylogenetic tree to the phylogenetic tree based on SNPs in coding regions
267 performed as part of previous analysis (Koutsovoulos et al. 2020).
268 The TE-based and SNP-based tree topologies are almost identical. In particular, the two
269 trees allowed defining three highly supported clades, with maximum support value (Fig.4).
270 Clades (2) and (3) were identical, including branching orders. Because clade (1) is a
271 polytomy in the SNP tree (A), relative branching of the three isolates cannot be compared
272 although clustering remains unchanged. In both trees, isolates R2-1 and R2-6 are outside of
273 the clusters but their relative positions differ. Similarly, although R1-6 is more closely related
274 to clusters (1) and (2) than the rest of the isolates in the two trees, its position also slightly
275 differs between the SNP-based (A) and TE-based (B) trees. Differences in relative positions
276 of R2-1 and R2-6 can be explained by the low distance observed between these two isolates
277 in the SNP-based analysis (version with branch length in sup. Fig 3).

278 Overall, the similarity between the SNP-based and TE frequency-based trees indicates that
279 most of the phylogenetic signal coming from variations in TE-frequencies between isolates
280 recapitulates the SNP-based genomic divergence between isolates.
281



282
283

284 **Fig 4: Phylogenetic tree for *M. incognita* isolates.**

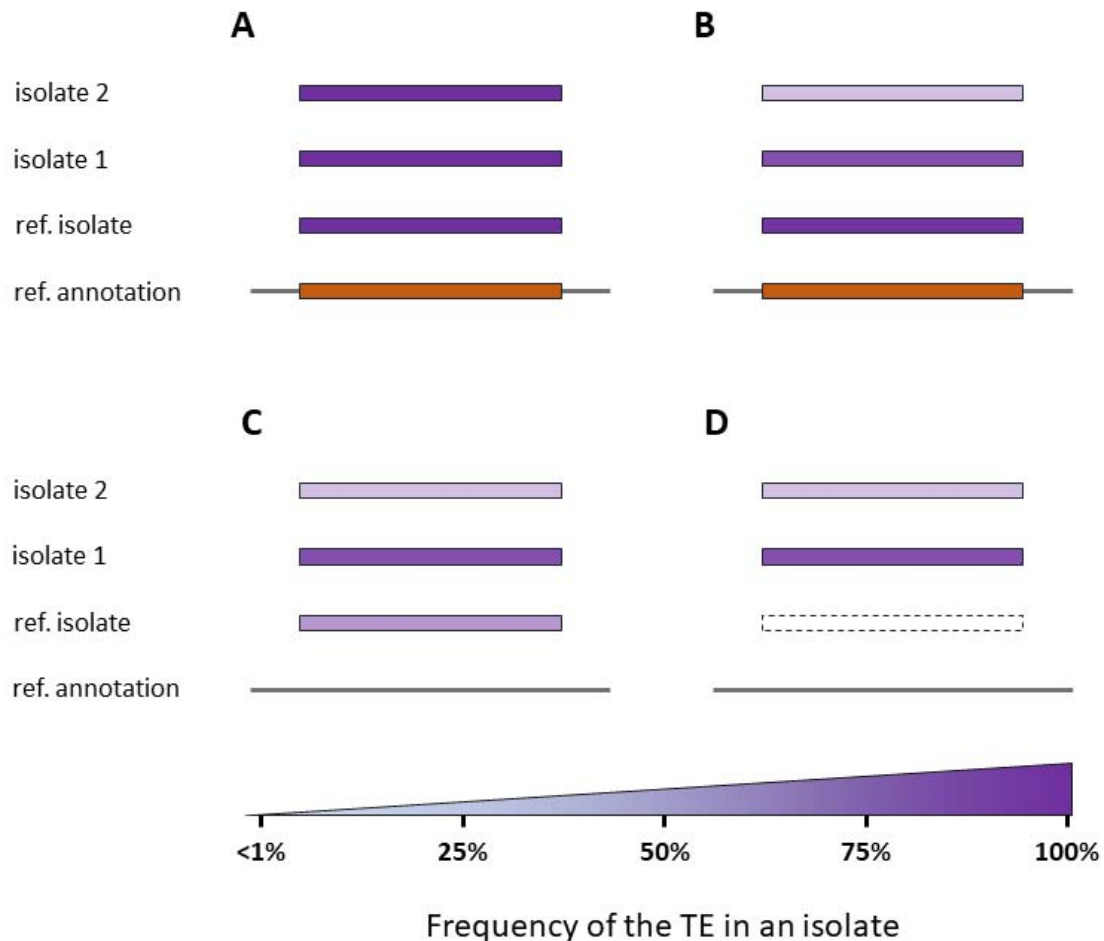
285 A- Phylogenetic tree based on SNP present in coding sequences. Maximum Likelihood (ML) tree
286 reconstruction. Branch length not displayed.

287 B- Phylogenetic tree based on TE-frequencies euclidean distances between isolates. Neighbor-
288 Joining (NJ) tree reconstruction. Branch length not displayed. In both trees, bootstrap support
289 values are indicated on the branches. Isolates enclosed in the dashed area form a super-
290 cluster composed of the clusters (1) and (2), and the isolate R1-6. Clades with bootstrap
291 support values ≤ 20 were collapsed and represented as a polytomy.

292

293 Some loci with TE frequency variations across isolates
294 correspond to neo-insertions.

295 As explained below (and cf methods), we categorized all the loci with TE frequency
296 variations across the isolates by (i) comparing their position to the TE annotation in the
297 reference genome, (ii) analysing TE frequency in the reference isolate Morelos, (iii)
298 comparing TE-frequencies detected for each isolate to the reference isolate Morelos. This
299 allowed defining, on the one hand, non-polymorphic and hence stable reference annotation,
300 and on the other hand, 3 categories of polymorphic (variable) loci (see Fig 5).



301

302 **Fig 5: Categories of polymorphic TE loci**

303 Orange boxes illustrate the presence of a TE at this position in the reference genome
304 annotation. Purple boxes illustrate the percentage of individuals in the isolates for which the
305 TE is present at this position (i.e. frequency). Frequencies values are reported as colour
306 gradients. A - non-polymorphic TE locus: a TE is predicted in the reference annotation
307 (orange box), is detected in all the isolates (purple box), but the presence frequency does

308 not vary substantially between the isolates (frequency variation < 1% among all the isolates).
309 B - polymorphic reference locus: a TE is predicted in the reference annotation, is detected in
310 the reference isolate Morelos with a frequency > 75%, and the presence frequency varies
311 between the isolates (frequency variation \geq 1% among all the isolates). This category
312 encompasses 2,096 loci. C - extra-detection: no TE is predicted at this locus in the reference
313 annotation but one is detected at a frequency >25% in the reference isolate Morelos, and
314 optionally in other isolates. This category counts 210 loci. D - neo-insertion: no TE is
315 predicted at this locus in the reference genome annotation and none is detected in the
316 reference isolate (dashed box, frequency < 1%), but a TE is detected in at least one other
317 isolate with a frequency \geq 25%. This category counts 287 loci.

318 Only 931 loci could not be assigned to the categories in Fig 5 and were discarded.
319 Uncategorised loci are cases where all the isolates (Morelos included) show frequencies <
320 25% or ambiguous cases with a low Morelos frequency (between 1% and 25%) and at least
321 one isolate showing a frequency > 25%.

322 Overall, 73.6% (2,593/3,524) of the loci with TE frequency variations could be assigned to
323 one of the 3 categories of TE-polymorphisms (B, C, D in Fig 5) and the decomposition per
324 TE order is given in Fig 6 and sup. Table 3.

325 The vast majority of the polymorphic loci (80.83 %; 2,096/2,593) corresponds to an already
326 existing TE-annotation in the reference genome and the corresponding TE is fixed
327 (frequency > 75 %) at least in the reference isolate Morelos but varies in at least another
328 isolate. These polymorphic loci cover ~21.6% (2,096/9,702) of the canonical TE annotations,
329 in total. These loci will be referred to as 'polymorphic reference loci' from now on (Fig 5B)
330 and they encompass both DNA- and Retro-transposons.

331 Then, we considered as 'neo-insertion' TEs present at a frequency >25% in at least one
332 isolate at a locus where no TE was annotated in the reference genome and the frequency of
333 TE presence was < error rate (~1%) in the reference Morelos isolate (Fig 5D). In total, 11.07
334 % (287/2,593) of the detected TE polymorphisms correspond to such neo-insertions. It
335 should be noted here that we consider neo-insertions as regard to the reference Morelos
336 isolate only and some of these so-called neo-insertions might represent TE loss in Morelos.
337 Comparison with the phylogenetic pattern of presence / absence will allow distinguishing
338 further the most parsimonious of these two possibilities.

339 Finally, we classified as 'extra-detection' (Fig 5C) (8.10%; 210/2,593) the loci where no TE
340 was initially annotated by REPET in the reference genome, but a TE was detected at a
341 frequency >25% at least in the ref isolate Morelos by PopoolationTE2. It should be noted
342 that 57.62% (121/210) of these loci correspond to draft annotations that have been
343 discarded during the filtering process to only select the canonical annotations. These draft
344 annotations might represent truncated or diverged versions of TE that exist in a more
345 canonical version in another locus in the genome. Half of the remaining 'extra-detections'
346 (45/89) are detected with low to moderate frequency (<43.5%) in the reference isolate
347 Morelos. We hypothesise that because they represent the minority form, these regions were
348 not taken into account during the assembly of the genome. This would explain why these
349 TEs could not be detected in the genome assembly by REPET (assembly-based approach)
350 but were identified with a read mapping approach on the genome plus reappearome by

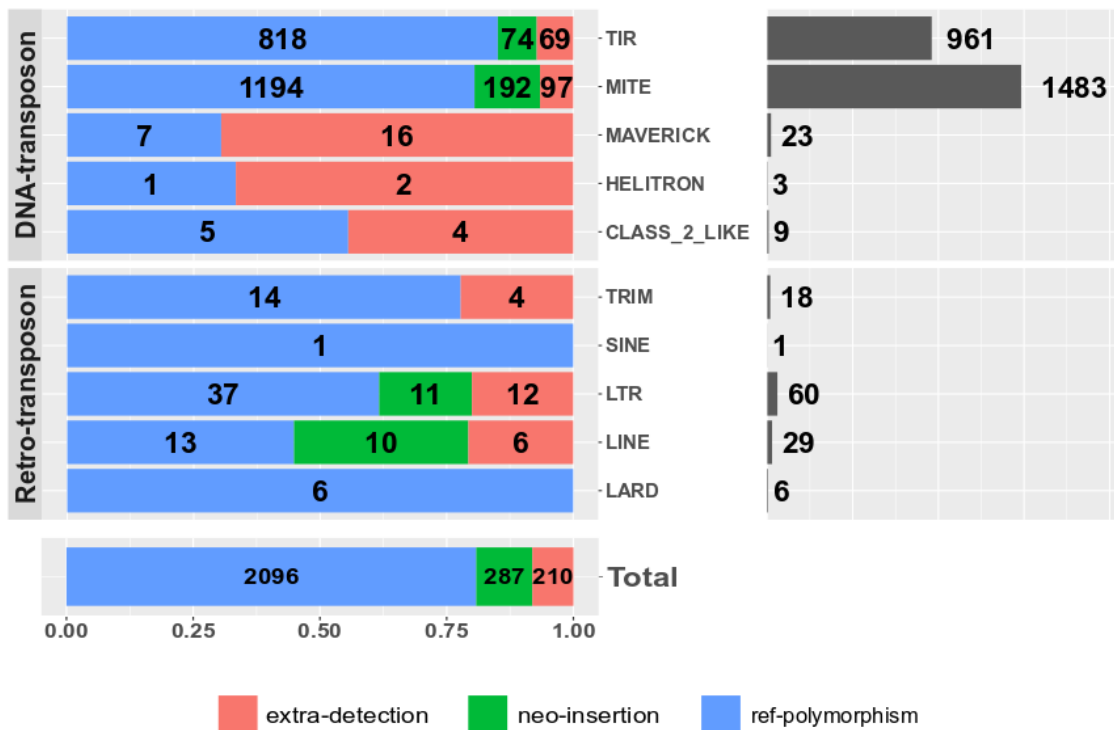
351 PopoolationTE2. The remaining 'extra-detections' might correspond to REPET false
 352 negatives, PopoolationTE false positives, or a combination of the two. Nonetheless, we can
 353 notice these cases only represent 1.69% (44/2,593) of the detected polymorphic TEs.
 354 Loci of variable frequency that could not be assigned to the 3 above-mentioned B, C, D categories
 355 (26.4% of them) were discarded from the rest of the analysis.

356

357 TIR and MITEs elements are overrepresented among TE- 358 polymorphisms.

359 By themselves, MITEs and TIRs elements encompass 94.25% (2,444/2,593) of the
 360 categorized TE-polymorphisms (Fig 6).

361 We showed that the polymorphism distribution varies significantly between the four
 362 categories presented in Fig. 5 (Chi-square test, p-value < 2.2e-16), indicating that some TE
 363 orders are characterised by specific polymorphisms types.
 364



365

366

367 **Fig 6: TE polymorphisms count per orders and types.**

368 Top left barplot shows TE polymorphisms distribution per type and per order. Bottom-left
 369 barplot summarizes TE polymorphisms distribution per type. In both barplots, the values in
 370 black represent the count per polymorphism type. Top-right barplot illustrates the total
 371 number of polymorphisms per order. Orders are sorted identically in both plots.

372

373

374 The analysis of the chi-square residuals (sup. Fig 4) shows MITEs and TIRs are the only two
375 orders presenting a relative lack of non-polymorphic TEs. Hence, in addition to being the
376 most abundant in the genome, these two TE orders are significantly enriched among
377 polymorphic loci. MITEs are over-represented in both TE polymorphisms types (polymorphic
378 ref. loci and neo-insertions, Fig5 B and D), suggesting a variety of activities within this order.
379 On the other hand, TIRs are found in excess in ref-polymorphisms but lack in neo-insertions.
380 This lack of neo-insertions in TIRs may indicate a lower activity in this order, or a more
381 efficient negative selection.

382 Finally, we observed a strong excess of Maverick among the extra-detection as almost 70%
383 of Maverick polymorphisms (16/23) (Fig 6) fell into this category. Consistent with the
384 observation that, globally, >50% of the extra detections were actually draft annotations
385 eliminated afterwards during filtering steps, $\frac{3}{4}$ (12/16) of these Maverick elements were also
386 actually present in the draft annotations but eliminated during filtering and thus only appear
387 here due to the stringency of our filtering.

388 Overall, in proportion, MITEs and TIRs elements are significantly over-represented in TE-
389 polymorphisms. This observation suggests TEs from MITEs and TIRs orders, in addition to
390 being the most numerous canonical TEs, might have been more active in the genome of *M.*
391 *incognita* than elements from other TE-orders.

392

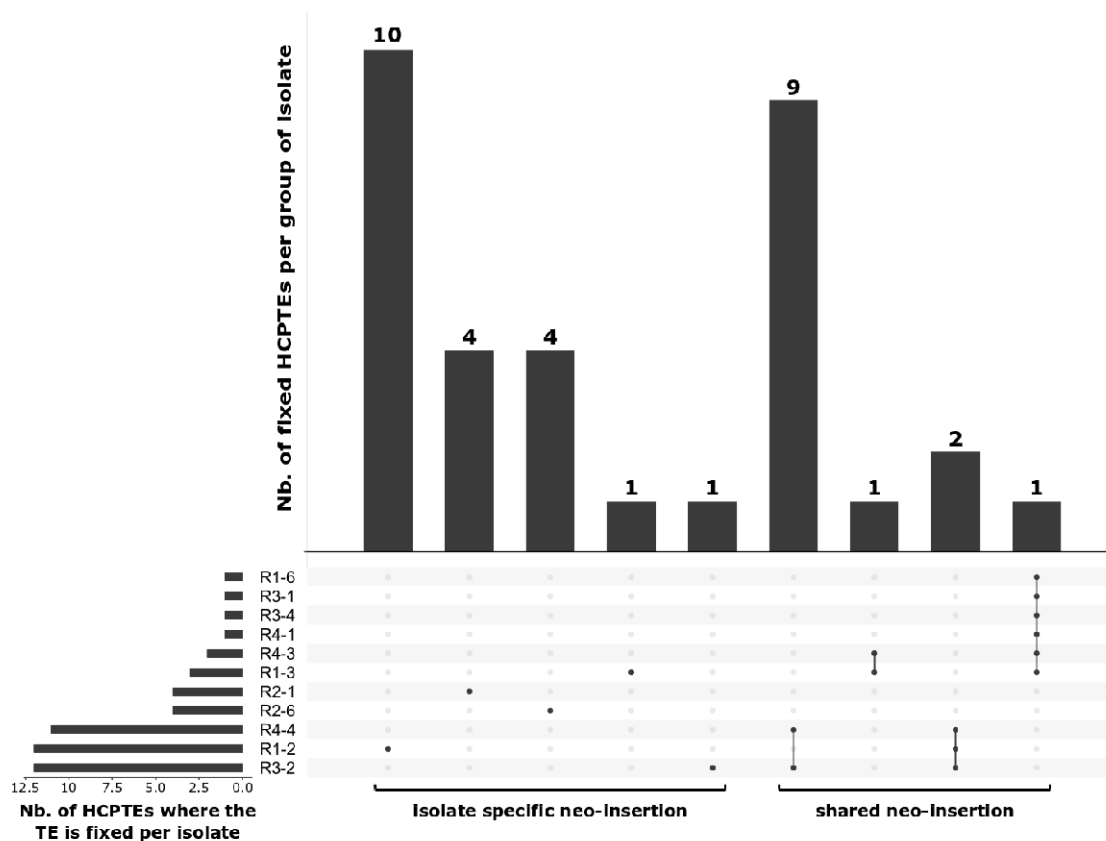
393 Some polymorphic loci showing contrasted frequency variations 394 between isolates represent neo-insertions.

395 We investigated the variability of TE presence frequency per locus over the 12 isolates for all
396 the categorized polymorphic loci in the genome.

397 In $\sim \frac{3}{4}$ (1,919/2,593) of the categorized polymorphic TE loci, the TE presence frequency is
398 homogeneous between isolates (cf methods). Said differently, it means that although we
399 observe variations in frequencies above the estimated error rate (<1%) between isolates,
400 these variations remain at low amplitude (maximum frequency variation between isolates
401 $\leq 25\%$ for a given locus). The vast majority (97.97%; 1,880/1,919) concerns loci where the
402 TE is present at a high frequency in all isolates (> 75%). These loci might be considered as
403 fixed in all the isolates. In the remaining 2.03% (39/1,919), the TE frequency is either
404 between 25 and 50% or between 50 and 75% in all isolates. As expected given our
405 methodology, all the high-frequency loci correspond to ref-polymorphisms while all the
406 intermediate frequency loci belong to extra-detections. Also, we did not detect loci where the
407 TE was present with low frequency (<25%) in all isolates as they did not meet the
408 categorisation criteria ($1\% < \text{freq} \leq 25$ in Morelos isolate).

409 In the 674 remaining polymorphic TE loci, TE frequency is heterogeneous, meaning the
410 frequency difference between at least two isolates is > 25%. Among the most extreme cases
411 of frequency variation per locus, we identified 33 loci for which the presence of the TE is
412 found with high frequencies (> 75%) for some isolate(s) while it is absent or rare (frequency
413 <25 %) in the other(s). These loci will be from now on referred to as HCPTEs standing for
414 "Highly Contrasted Polymorphic TE" loci. Because they are highly contrasted, these loci
415 might represent differential fixation/loss across isolates and will be the focus of the following
416 analyses.

417 HCPTes encompass 19 MITEs elements, 12 TIRs and 2 LINEs. We can also notice that
 418 some consensuses are more involved in HCPTes as 4 TE consensuses are responsible for
 419 72.72% (24/33) of these polymorphisms.
 420 Interestingly, all the HCPTes loci correspond to neo-insertions regarding the reference
 421 genome, meaning that no TE was annotated in the reference genome at this location and
 422 the TE presence frequency is < 1% in the Morelos reference isolate. As described in Fig. 7,
 423 most of these fixed neo-insertions (20/33) are specific to an isolate. However, we also found
 424 neo-insertions shared by 2 (10/33), 3 (2/33) or even 6 isolates (1/33).
 425 Interestingly, all the shared neo-insertions were between isolates present in a same cluster
 426 in the phylogenetic trees (TE-based and SNP-based in Fig. 4), suggesting they might have
 427 been fixed in a common ancestor and then inherited. For example, two neo-insertions are
 428 shared by isolates R4-4, R1-2 and R3-2 which belong to the same cluster 1 and one neo-
 429 insertion is shared by isolates R4-3 and R1-3 which belong to the same cluster 2. Even the
 430 neo-insertion shared by 6 isolates follows this pattern as all the concerned isolates belong to
 431 the same super-cluster composed of the cluster 2 and 3 (dashed line in Fig 4).
 432 Hence, the phylogenetic distribution reinforces the idea that these cases are more likely to
 433 represent branch-specific neo-insertions than multiple independent losses, including in the
 434 reference isolate Morelos.
 435 R1-2, R3-2, and R4-4 show the highest number of neo-insertions among isolates. However,
 436 their profiles are quite different. In 10/12 HCPTes involving R1-2, the TE is present only in
 437 this isolate while most of the HCPTes involving R3-2 and R4-4 are neo-insertions shared
 438 with close isolates. This may be related to the relative phylogenetic divergence of those
 439 isolates (sup Fig 3), which shows that R1-2 is the most divergent isolate, while R3-2 is quite
 440 close to its neighbour, and especially to R4-4.



442 **Fig 7: HCPTes Neo-insertions specificity among the isolates.**

443 The central plot shows how many and which isolate(s) share common HCPTes neo-
444 insertion(s), every line representing an isolate. Columns with several dots linked by a line
445 indicate shared HCPTes neo-insertion(s) between isolates. Each dot represents which
446 isolate is involved. Columns with a single dot design isolate-specific HCPTes neo-
447 insertion(s). The top bar plot indicates how many HCPTes neo-insertions the corresponding
448 group of isolate shares. The left side barplot specifies how many HCPTes neo-insertion(s)
449 occurred in a given isolate.

450 **Functional impact of TE neo-insertion and validation of in silico**
451 **predictions**

452 Interestingly, the vast majority (22/33) of the fixed HCPTes are inserted inside a gene or in a
453 possible regulatory region (i.e. 1 kb region upstream of a gene). These fixed neo-insertions
454 might have a functional impact in *M. incognita*. Overall, 27 different protein-coding genes are
455 possibly impacted by the 22 neo-insertions, some genes being in the opposite direction at a
456 neo-insertion point (overlapping this insertion point or being at max 1kb downstream). More
457 than 80% of these genes (22/27) show a substantial expression level during at least one life
458 stage of the nematode life cycle, suggesting the impacted genes are functional in the *M.*
459 *incognita* genome (cf. methods). Some of the impacted genes (37.04%, 10/27) are specific
460 to the *Meloidogyne* genus. Despite being all conserved in multiple *Meloidogyne* species,
461 reinforcing their importance in the genus, they have no predicted orthologs in other
462 nematodes. Among the remaining genes, one is specific to *M. incognita* (even other
463 *Meloidogyne* species have no ortholog). Another one is present in multiple *Meloidogyne*
464 species and otherwise only found in other Plant Parasitic Nematodes species (PPN)
465 (*Ditylenchus destructor*, *Globodera rostochiensis*) (see sup. Table 4). Conservation of these
466 genes across multiple PPN but exclusion from the rest of the nematodes suggest these
467 genes might be involved in important functions relative to these organisms' lifestyle,
468 including plant parasitism itself.

469 To experimentally validate in-silico predictions of TE neo-insertions with potential functional
470 impact, we performed PCR experiments on 5 of the 24 HCPTes loci falling in coding or
471 possible regulatory regions. To perform these PCR validations, we used the DNA remaining
472 from previous extractions performed on the *M. incognita* isolates for population genomics
473 analysis (Koutsovoulos et al. 2020). Basically, the principle was to validate whether the
474 highly contrasted frequencies (>75% / <25%) obtained by PopoolationTE2 actually
475 corresponded to absence/presence of a TE at the locus under consideration (cf methods).
476 One isolate (R3-1) presented no amplification in any of the tested loci nor in the positive
477 control. After testing the DNA concentration in the sample, we concluded that the DNA
478 quantity was limiting in this isolate and decided to discard it from the analysis.

479 For four of the five tested HCPTes loci, we could validate by PCR the differential
480 presence/absence of a sequence at this position, predicted by PopoolationTE2 across the
481 different isolates (Fig 8; supplementary material 4).

482 In one of the five tested loci, named locus 1, we could i) validate by PCR the presence of a
483 sequence at this position for the isolates presenting a PopoolationTE2 frequency >75% and
484 absence for those having a frequency <25%; ii) also validate by sequencing that the
485 sequence itself corresponded to the TE (a MITE) under consideration. This case is further
486 explained in detail below and in Fig. 8.

487 According to PopoolationTE2 frequencies, the locus 1 MITE is inserted and fixed in 3
488 isolates (R1-2, R3-2, R4-4) as the estimated frequencies are higher than 75% in these
489 isolates. We assumed the TE is absent from the rest of the isolates as all of them display
490 frequencies <5%. To validate this differential presence across the isolates, we designed
491 specific primers from each side of the estimated insertion point so that the amplicon should
492 measure 973 bp with the TE insertion and 180 bp without.

493 The PCR results are coherent with the frequency predictions as only R1-2, R3-2, and R4-4
494 display a ~1 kb amplicon while all the other isolates show a ~0.2 kb amplicon (Fig 8). Hence,
495 as expected, only the 3 isolates with a predicted TE frequency >75% at this locus exhibit a
496 longer region, compatible with the MITE insertion.

497 To make sure the amplified regions corresponded to the expected MITE, we sequenced the
498 amplicons for the 3 predicted insertions and aligned the sequences to the TE consensus and
499 the genomic region surrounding the estimated insertion point (see supplementary material
500 4). R-1_2, R-3_2, and R-4_4 amplicon sequences all covered a significant part of the TE
501 consensus sequence length (> 78%) with high % identity (> 87%) and only a few gaps
502 (<5%). These results confirm that the inserted sequence corresponds to the predicted TE
503 consensus. Moreover, all the 3 amplicons aligned on the genomic region downstream of the
504 insertion point with high % identity ($\geq 99\%$), which helped us to further determine the real
505 position of the insertion point. The real insertion point is 26 bp upstream of the one predicted
506 by PopoolationTE2 and falls in the forward primer sequence. This explains why the amplicon
507 sequences do not align on the region upstream the insertion point.

508 We also noticed that the inserted TE sequences slightly diverged between the isolates while
509 the genomic region surrounding the insertion point remains identical. Interestingly, the level
510 of divergence in the TE sequence does not follow the phylogeny as R-4_4 is closer to R-1_2
511 than to R-3_2 (see sup. Table 5).

512 Finally, in the Morelos isolate, as well as R-2_1, and R-2_6 isolates, the sequencing of the
513 amplicon validated the absence of insertions as the sequences aligned on the genomic
514 region surrounding the insertion point with high % identity (99, 97, 87 % respectively) but not
515 with the MITE consensus.

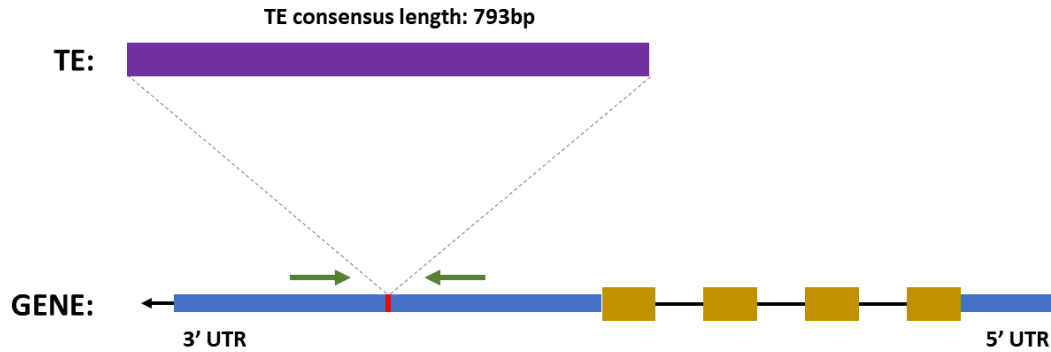
516 Hence, we fully validated experimentally the presence/absence profile across isolates
517 predicted in silico at this locus.

518 In the *M. incognita* genome, this neo-insertion is predicted to occur in the 3' UTR region of a
519 gene (Minc3s00026g01668). This gene has no obvious predicted function, as no conserved
520 protein domain is detected and no homology to another protein with an annotated function
521 could be found. However, orthologs were found in the genomes of several other
522 *Meloidogyne* species (*M. arenaria*, *M. javanica*, *M. floridensis*, *M. enterolobii*, and *M.*
523 *graminicola*), ruling out the possibility that this gene results from a prediction error from gene
524 calling software. The broad conservation of this gene in the *Meloidogyne* genus suggests
525 this gene might be important for *Meloidogyne* biology and survival.

526 In the Morelos isolate, for which no TE was inserted at this position, this gene is substantially
527 supported by transcriptomic RNA-seq data during the whole life cycle of the nematode (see
528 supplementary material 1), suggesting this gene is probably functionally important in *M.*
529 *incognita* and other root-knot nematodes. Consequently, the insertion of the TE in R-1_2, R-
530 3_2, and R-4_4 genome at this locus could have functional impacts.

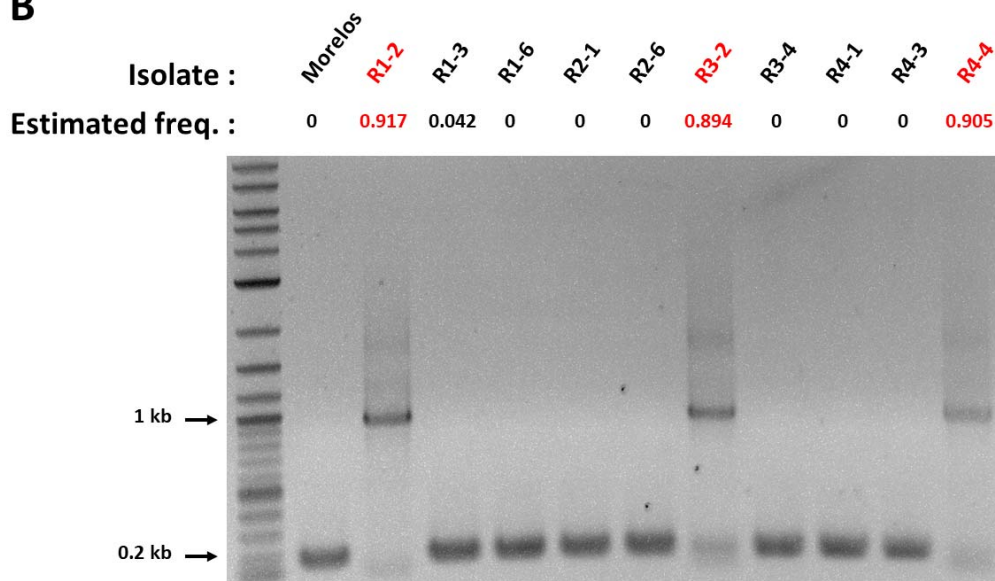
531
532
533

A



534

B



535

536

537 **Fig 8: Experimental validation of a predicted neo-insertion.**

538 A- Diagram of the TE neo-insertion. The neo-insertion of the MITE element occurs in the
539 3'UTR region of the gene (Minc3s00026g01668). Blue boxes illustrate the 3' and 5' UTR
540 regions of the gene while the yellow boxes picture the exons. Green arrows represent the
541 primers used to amplify the region. Gene subparts and TE representations are not at scale.
542 Predicted size of the amplicon: 973 bp with the TE insertion, 180 bp without.

543

544 B- PCR validation of the TE neo-insertion. Estimated freq. values correspond to the
545 proportion of individuals per isolate predicted to have the TE at this position
546 (PopulationTE2). Isolates in red were predicted to have the TE inserted at this locus. Only
547 these isolates show an amplicon with a size suggesting an insertion. See supplementary
548 material 4 for sequences.

549 Discussion

550 The role of TE in genome plasticity and adaptive evolution of 551 root-knot nematodes

552 *M. incognita* is a parthenogenetic mitotic nematode with a major agronomic impact. How this
553 pest adapts to its environment in the absence of sexual recombination remains unresolved.
554 In this study, we investigated whether TE movements could constitute a mechanism of
555 genome plasticity compatible with adaptive evolution.

556 In *M. javanica*, a closely related root-knot nematode, comparison between an avirulent line
557 unable to infest tomato plants carrying a nematode resistance gene and another virulent line
558 that overcame this resistance, led to the identification of a gene present in the avirulent
559 nematodes but absent from the virulent ones. Interestingly, the gene under consideration is
560 present in a TIR-like DNA transposon and its absence in the virulent line suggests this is due
561 to excision of the transposon and thus that TE activity plays a role in *M. javanica* adaptive
562 evolution (Gross and Williamson 2011). However, this is so far the sole report and no large-
563 scale analysis of the global role of TE in the adaptive evolution of root-knot nematodes has
564 been performed.

565 In *M. incognita*, convergent gene losses at the whole genome level between two virulent
566 populations compared to their avirulent populations of origin were recently reported
567 (Castagnone-Sereno et al. 2019). However, although TE activity might be involved in these
568 convergent gene losses, this has never been shown.

569 Hence, the importance of TE activity both in the genome plasticity and adaptive evolution of
570 the root-knot nematodes has never been assessed so far.

571 More broadly, except for *C. elegans* (Bessereau 2006; Laricchia et al. 2017), little to nothing
572 is known yet about the TE dynamics and its impacts on nematode genomes and possibly in
573 their adaptive evolution.

574 Here, we analysed the *M. incognita* TE dynamics using TE presence frequencies variations
575 between isolates as a reporter of TE activity. We also assessed whether TE activity could
576 have a functional impact, a necessary prerequisite for natural selection and adaptive
577 evolution.

578 Our results established that TEs were likely recently active in *M. incognita* and that
579 thousands of loci in the genome show substantial variations in the frequency of reads
580 supporting a TE at this position, across different geographical isolates. It should be noted
581 here that the total impact of TE on the genome dynamics is probably underestimated, in part
582 because of our strategy to eliminate false positives as much as possible by applying a series
583 of stringent filters, in another part because of the intrinsic limitations of the tools, such as the
584 incapacity of PopoolationTE2 to detect nested TEs (Kofler et al. 2016). Despite this, we
585 identified highly contrasted TE loci across the populations and some of these represented
586 neo-insertions of TE, specifically in some isolates or some branches ancestral to several
587 isolates. Because for certain TE, the insertion took place within a coding sequence or in a
588 possible regulatory region, these insertions have a potential impact either by disrupting the
589 encoded protein or modifying the gene expression pattern.

590 We confirmed insertions of the TEs for some of these functionally-important regions by PCR
591 and thus a possible functional impact. However, this functional impact itself would need to be
592 evaluated in the future either by generating transcriptomics data for the different isolates

593 enabling study of difference in gene expression patterns or transcript length or by proteomic
594 studies to directly search for differences at the encoded protein level.

595 Regardless of the future experimental validation of the functional impact, one important
596 question concerns the current preliminary evidence for a possible role in the nematode
597 adaptive evolution. Because some of the impacted genes are specific to plant-parasitic
598 species and yet conserved in several of these phytoparasites, a role in plant parasitism is
599 possible. However, in the absence of known protein domains or functional characterization
600 of these genes, the exact biochemical activity or biological processes in which they might be
601 involved is unknown. Furthermore, for the functional impact, we have exclusively considered
602 neo-insertions as the most evident cases. Some TE showing less contrasted changes in
603 frequencies in coding or possible regulatory regions might be associated to genes with
604 characterized functions but as these cases are less clear cut, they were not considered.

605 It should be noted that *M. incognita* is a relatively recent model with a first version of the
606 genome available only since 2008 (Abad et al. 2008) and only up to a dozen of genes
607 functionally characterized. In comparison, the model nematode *C. elegans* was the first
608 animal genome to be sequenced in 1998 (The *C. elegans* Genome Sequencing Consortium
609 1998) and as early as 2003, more than 85% of the genes had already been inactivated one
610 by one via RNAi to monitor the effect on the worm phenotype (Kamath et al. 2003).
611 Therefore, although no evident role in adaptive evolution for the *M. incognita* genes
612 impacted by TE insertions could be reported so far, future functional characterization might
613 bring more evidence.

614 Taking into account all the TE loci and regardless of the potential functional impact on
615 genes, we found that the pattern of variations of TE frequencies across the loci between the
616 different populations recapitulated almost exactly the phylogeny of the isolates built on SNP
617 in coding regions. Thus, it seems that most of the divergence in terms of TE pattern follows
618 the divergence at the nucleotide level and thus the phylogeny of the isolates. Almost the
619 same conclusion was drawn by comparing SNV data to TE variation data across different *C.*
620 *elegans* populations (Laricchia et al. 2017). In *M. incognita*, the phylogeny of isolates does
621 not significantly correlate with the biological traits that had been surveyed, namely
622 geographical distribution, range of compatible host plants and nature of the crop currently
623 infected (Koutsovoulos et al. 2020). Interestingly, no correlation was also observed between
624 variations in TE frequencies and geographical distribution for European *Drosophila*
625 populations (Lerat et al. 2019). As for the other traits considered here (range of compatible
626 host plants and nature of the plant infected), the lack of evident correlation between those
627 traits and phylogenetic signals regardless whether it is TE-based or SNV-based suggest that
628 most of the variations follow the drift between isolates and are not necessarily adaptive,
629 which is not surprising.

630 It should be noted that, in terms of functional impact, we have so far only considered the
631 active role of TE and not analyzed yet their possible passive roles. Indeed, being repetitive,
632 TE can be involved in illegitimate recombination events, and also generate loops in the DNA
633 which can eventually be excised and lead to gene loss. In parallel, genes can be hitchhiked
634 by TE and multiplied in the genome. These different properties can clearly impact gene copy
635 number variations (CNV). CNV are known to be involved in genomic plasticity and in
636 adaptive evolution (Katju and Bergthorsson 2013). As stated above, convergent gene CNV
637 have been associated with plant resistance breaking down by *M. incognita*, although neither
638 direct functional link between these CNV and mechanisms of resistance breaking down nor
639 evidence for involvement of TEs have been shown so far (Castagnone-Sereno et al. 2019).
640 Actually, these analyses have been done in a previous version of the *M. incognita* genome

641 assembly which was partially incomplete (Abad et al. 2008) in comparison to the most recent
642 available genome assembly (Blanc-Mathieu et al. 2017). Reinvestigating CNV and the
643 possible involvement of TEs in association to an adaptive process such as resistance
644 breaking down on this more complete and more recent assembly would be interesting.
645 However, the current genome is still fragmentary and limits structural studies such as the
646 identification of TE-rich vs. TE-poor genome regions in possible association with CNV loci.
647 Future more contiguous versions of the genomes of root-knot nematodes will undoubtedly
648 enable such interesting perspectives to be undertaken.
649

650 TE-load and composition in a clonal allopolyploid species

651 *M. incognita* is an asexual (mitotic parthenogenetic), polyploid, and hybrid species. These
652 three features are expected to impact TE load in the genome with various intensities and
653 possibly conflicting effects.

654 Contradictory theories exist concerning the activity/proliferation of TEs as a function of the
655 reproductive mode. The higher efficacy of selection under sexual reproduction can be
656 viewed as an efficient system to purge TEs and control their proliferation. Supporting these
657 views, in parasitoid wasps, TE load was shown to be higher in asexual lineages induced by
658 the endosymbiotic *Wolbachia* bacteria than in sexual lineages (Kraaijeveld et al. 2012).
659 However, whether this higher load is a consequence of the shift in reproductive mode or of
660 *Wolbachia* infection remains to be clarified.

661 In an opposite theory, sexual reproduction can also be considered as a way for TEs to
662 spread across individuals within the population whereas in clonal reproduction the
663 transposons are trapped exclusively in the offspring of the holding individual. Under this
664 view, asexual reproduction is predicted to reduce TE load as TE are unable to spread in
665 other individuals, and are thus removed by genetic drift and/or purifying selection in the
666 long term (Wright and Finnegan 2001). Consistent with this theory, comparison of sexual
667 and asexual *Saccharomyces cerevisiae* populations showed that the TE loads decrease
668 rapidly under asexual reproduction (Bast et al. 2019).

669 Hence, whether the TE-load is expected to be higher or lower in species with clonal vs.
670 sexual reproduction remains unclear and other conflicting factors such as TE excision rate
671 and the effective size of the population probably blur the signal (Glémin et al. 2019).
672 Interestingly, at a broader scale, a comparative analysis of different lineages of sexual and
673 asexual arthropods revealed no evidence for differences in TE load according to the
674 reproductive modes (Bast et al. 2015). Similar conclusions were drawn for nematodes
675 (Szitenberg et al. 2016), although only one asexually-reproducing species was present in the
676 comparative analysis.

677 Polyploidy, in contrast, is commonly accepted as a major event initially favouring the
678 multiplication and activity of TEs. This is clearly described with numerous examples in plants
679 (Vicent and Casacuberta 2017) and some examples are also emerging in animals
680 (Rodriguez and Arkhipova 2018). When hybridization and polyploidy are combined, this can
681 lead to TE bursts in the genome. As proposed by Barbara McClintock, allopolyploidization
682 produces a "genomic shock", a genome instability associated with the relaxation of the TE
683 silencing mechanisms and the reactivation of ancient TEs (McClintock 1984; Mhiri et al.
684 2019).

685 Hybridization, polyploidy and asexual reproduction are combined in *M. incognita* with relative
686 effects on the TE load extremely challenging, if not impossible, to disentangle. Initial
687 comparisons of the TE loads in three allopolyploid clonal *Meloidogyne* against a diploid
688 facultative sexual relative suggested a higher TE load in the clonal species (Blanc-Mathieu
689 et al. 2017). However, to differentiate the relative contribution of each of these three features
690 to the *M. incognita* TE load, it would be necessary to conduct comparative analysis with a
691 same method on diploid asexuals, on polyploid sexuals as well as on diploid asexuals in the
692 genus *Meloidogyne*, and ideally with and without hybrid origin. So far, genomic sequences
693 are only available for other polyploid clonal species, which are all suspected to have a hybrid
694 origin (Blanc-Mathieu et al. 2017; Szitenberg et al. 2017; Koutsovoulos et al. 2019; Susič
695 et al. 2020), and, apart from that, only two diploid facultative sexual species (Opperman et al.
696 2008; Somvanshi et al. 2018). Hence, further sampling of *Meloidogyne* species with diverse
697 ploidy levels and reproductive modes will be necessary to answer these questions.

698

699 Regardless of the TE load, we found that DNA transposons were majoritary in *M. incognita*.
700 Using the same annotation method, we found a similar result in *C. elegans*. Interestingly,
701 even if the methodology used was different, a similar observation was made at the whole
702 nematoda level (Szitenberg et al. 2016), suggesting a higher abundance of DNA
703 transposons might be a general feature of nematode genomes.

704

705 TE show signs of recent activity in *M. incognita* and they might 706 still be active

707 In the current analysis, we used variations in TE frequencies between geographical isolates
708 across loci in the *M. incognita* genome as a reporter of their activity. We have shown 75% of
709 the polymorphic TE loci display moderate frequency variations between isolates (<25%), a
710 majority being found with high frequencies (> 75%) in all the isolates simultaneously. Hence,
711 a substantial part of the TE can be considered as stable and fixed among the isolates.

712 Nevertheless, the remaining quarter of polymorphic TE loci present frequency variations
713 across the isolates higher than 25%. This observation concerns both the TE already present
714 in the reference genome, but also the neo-insertions. We even detected loci where the TE
715 frequencies were so contrasted between the isolates (HCPTes) that we could predict the TE
716 presence/absence pattern among the isolates. Such frequency variations between isolates,
717 and the fact that part of the HCPTes are isolate-specific neo-insertions, constitute strong
718 evidence for TE activity in the *M. incognita* genome.

719 We then evaluated how recent this activity could be, using % identity of the TE copies with
720 their respective consensus as a proxy for their age. We showed that a substantial
721 proportion of the canonical TE annotations were highly similar to their consensus, indicating
722 most of these TE copies were recent in the genome. This result suggests a TE burst in the
723 *M. incognita* genome, which would be consistent with its likely recent hybrid origin (Blanc-
724 Mathieu et al. 2017). Indeed, as evoked previously, it is well established that hybridization
725 events can lead to a relaxation of the TE silencing mechanisms and consequently to a TE
726 expansion (Belyayev 2014; Guerreiro 2014; Rodriguez and Arkhipova 2018).

727 However, as suggested in (Bourgeois and Boissinot 2019), the extent of this phenomenon
728 might differ depending on the TE order. In *M. incognita*, MITEs and TIRs alone account for
729 ~2/3 of the canonical TE annotations, but their fate in the genome seems to have followed

730 different paths. Indeed, MITE copy numbers almost linearly increase as a function of the
731 identity rate with their consensus, which suggests they might have progressively invaded the
732 genome being uncontrolled or poorly controlled as suggested for the rice genome (Lu et al.
733 2017). On the opposite, almost all the TIR copies share high percentage identity with their
734 consensus which could be reminiscent of a rapid and recent burst. Nevertheless, this
735 burst could have quickly been under control as we observed that the TIR neo-insertions are
736 less numerous than expected owing to their abundance in the genome.
737 Because no molecular clock is available for *M. incognita*, it is impossible to evaluate more
738 precisely when this burst would have happened and how fast each TE from each order
739 would have spread in the genome. However, while an absolute dating of TE activities is
740 currently not possible, a relative timing of the events regarding speciation and diversification
741 can still be deduced from distribution of TE loci frequencies across the isolates. Indeed, we
742 have shown that some neo-insertion were shared between isolates and that in each case,
743 the concerned isolates were related according to the phylogeny. This observation indicates
744 that these neo-insertions occurred after *M. incognita* speciation, but before the diversification
745 of the phylogenetically-related isolates, in a common ancestor. Other TE neo-insertions, in
746 contrast, were so far isolate-specific, suggesting some TE movements were more recent and
747 that TE mobility might be a continuous phenomenon.
748 As there is no correlation between life history traits and geography among those isolates
749 (Koutsovoulos 2019), we can make the hypothesis that the Brazilian geographical isolates
750 we compared were recently spread by human intervention across different cultivated fields
751 during the modern era of extensive agriculture. Hence we could conclude that these neo-
752 insertions happened in the last centuries.
753 Overall, these observations, the distribution of percent identities of some TE copies to their
754 consensus shifted towards high value, as well as support for transcriptional activity of
755 some of the genes involved in the transposition machinery, suggest TE have recently been
756 active in *M. incognita* and are possibly still active.
757

758 Concluding remarks

759 In this study we used population genomics technique and statistical analyses of the results
760 to assess whether TE might contribute to the genome dynamics of *M. incognita* and possibly
761 to its adaptive evolution. Overall, we provided a body of evidence suggesting TE have been
762 at least recently active and might still be active. With thousands of loci showing variations in
763 TE presence frequencies across geographical isolates, there is a clear impact on the *M.*
764 *incognita* genome plasticity. Some TE being neo-inserted in coding or regulatory regions
765 might have a functional impact. Although no clear connection with a role in adaptive
766 evolution could be made so far, based on the few impacted coding loci we experimentally
767 checked in this study, this is not to be excluded given the current lack of large-scale
768 functional information for this species. This pioneering study constitutes a valuable resource
769 and opens new perspectives for future targeted investigation of the potential effect of TE
770 dynamics on the evolution, fitness and adaptability of *M. incognita*.

771 Materials and Methods

772 Material

773 The genome of *M. incognita*

774 We used the genome assembly published in (Blanc-Mathieu et al. 2017) as a reference for
775 TE prediction and annotation as well as for read-mapping of the different geographical
776 isolates (Koutsovoulos et al. 2020), used for prediction of TE presence frequencies.

777 Briefly, the triploid *M. incognita* genome is 185Mb long with ~12,000 scaffolds and a N50
778 length of ~38 kb. Although the genome is triploid, because of the high nucleotide divergence
779 between the genome copies (8% on average), most of these genome copies have been
780 correctly separated during genome assembly, which can be considered effectively haploid
781 (Blanc-Mathieu et al. 2017; Koutsovoulos et al. 2020). This reference genome originally
782 came from a *M. incognita* population from the Morelos region of Mexico and was reared on
783 tomato plants from the offspring of one single female in our laboratory.

784 The genome of *C. elegans*

785 We used the *C. elegans* genome (The *C. elegans* Genome Sequencing Consortium 1998)
786 assembly (PRJNA13758) to perform its repeatome prediction and annotation and compare
787 our results to the literature as a methodological validation.

788 Genome reads for 12 *M. incognita* geographical isolates

789 To predict the presence frequencies at TE loci across different *M. incognita* isolates, we
790 used whole-genome sequencing data from pools of individuals from 12 different
791 geographical regions (sup. Fig 2 & sup. Table 6). One pool corresponds to the Morelos
792 isolates used to produce the *M. incognita* reference genome itself, as described above. The
793 11 other pools correspond to different geographical isolates across Brazil as described in
794 (Koutsovoulos et al. 2020).

795 All the samples were reared from the offspring of one single female and multiplied on tomato
796 plants. Then, approximately 1 million individuals were pooled and sequenced by Illumina
797 paired-end reads (2*150bp). Libraries sizes vary between 74 and 76 million reads
798 (Koutsovoulos et al. 2020).

799 We used cutadapt-1.15 ([Martin 2011](#)) to trim adapters, discard small reads, and trim low-
800 quality bases in reads boundaries (`-max-n=5 -q 20,20 -m 51 -j 32 -a`
801 `AGATCGGAAGAGCACACGTCTGAACTCCAGTCA` `-A`
802 `AGATCGGAAGAGCGTCGTGTAGGGAAAGAGTGT`). Then, for each library, we performed
803 a `fastqc` v-0.11.8 (Andrew S., 2010:
804 <http://www.bioinformatics.babraham.ac.uk/projects/fastqc>) analysis to evaluate the quality of
805 the reads. FastQC results analyses showed that no additional filtering or cleaning step was
806 needed and no further read was discarded.

807

808 Methods

809 We performed the statistical analysis and the graphical representation using R' v-3.6.3 and
810 the following libraries: ggplot2, reshape2, dplyr, ggpubr, phangorn, plyr, and UpSetr. All
811 codes and analysis workflows are available in the supplementary materials 1 to 3. For
812 experimental validations, see supplementary material 4.

813

814

815 *M. incognita* and *C. elegans* repeatome predictions and annotations.

816 We predicted and annotated the *M. incognita* and *C. elegans* repeatomes following the
817 same protocol as thoroughly explained in (Koutsovoulos et al. 2019). We define the
818 repeatome as all the repeated sequences in the genome, excluding Simple Sequence
819 Repeats (SSR) and microsatellites. Then, following the above-mentioned protocol, we
820 further analysed each repeatome to isolate annotations with canonical signatures of
821 Transposable Elements (TEs).

822 Below, we briefly explain each step and describe protocol adjustments.

823

824 Genome pre-processing.

825 Unknown nucleotides 'Ns' encompass 1.81% of the *M. incognita* reference genome and
826 need to be trimmed before repeatome predictions. We created a modified version of the
827 genome by splitting it at N stretches of length 11 or more and then trimming all N, using
828 dbchunk.py from the REPET package (Quesneville et al. 2005; Flutre et al. 2011). As this
829 increases genome fragmentation and may, in turn, lead to false positives in TE detection, we
830 only kept chunks of length above the L90 chunk length threshold, which is 4,891 bp. This
831 modified version of the genome was only used to perform the *de novo* prediction of the TE
832 consensus library (below). The TE annotation was performed on the whole reference
833 genome.

834 The *C. elegans* reference genome was entirely resolved (no N), at the chromosome-scale.
835 Hence, we used the whole assembly as is to perform the *de novo* prediction analysis.

836

837 *De novo* prediction: constituting draft TE-consensus libraries.

838 For each species, we used the TEdenovo pipeline from the REPET package to generate a
839 draft TE-consensus library..

840 Briefly, TEdenovo pipeline i) realises a self-alignment of the input genome to detect
841 repetitions, ii) clusters the repetitions, iii) performs multiple alignments from the clustered
842 repetitions to create consensus sequences, and eventually, iv) classify the consensus
843 sequence following the Wicker's classification (Wicker et al. 2007) using structural and
844 homology based information. One of the most critical steps of this process concerns the
845 clustering of the repetitions as it requires prior knowledge about assembly ploidy and
846 phasing quality.

847 We ran the analysis considering the modified *M. incognita* reference assembly previously
848 described as triploid and set the 'minNbSeqPerGroup' parameter to 7 (*i.e* 2n+1). As the *C.*
849 *elegans* assembly was haploid, we set the same parameter to 3.

850 All the remaining parameters values set in these analyses can be found in the TEdenovo
851 configuration files (see supplementary material 2).

852

853 Automated curation of the TE-consensus libraries.

854 To limit the redundancy in the previously created TE consensus libraries and the false
855 positives, we performed an automated curation step. Briefly, for each species, i) we
856 performed a minimal annotation (steps 1, 2, 3, 7 of TEannot) of their genome with their
857 respective draft TE-consensus libraries, and ii) only retained consensus sequences with at
858 least one Full-Length Copy (FLC) annotated in the genome. All parameters values are
859 described in the configuration files produced in the supplementary material 2.

860

861 Repeatome annotation

862 For each species, we performed a full annotation (steps 1, 2, 3, 4, 5, 7, and 8) of their
863 genome with their respective cleaned TE-consensus libraries using TEannot from the
864 REPET package. The obtained repeatome annotations (excluding SSR and microsatellites)
865 were exported for further analyses. All parameters values are described in the configuration
866 files produced in the supplementary material 2.

867

868 Repeatome post-processing: identifying annotations with canonical signatures of
869 TEs.

870 Using in house scripts (see supplementary material 2), we analysed REPET outputs to retain
871 annotations with canonical signatures of Transposable Elements (TEs) from the rest of the
872 repeatomes. The same parameters were set for *M. incognita* and *C. elegans*. Briefly, for
873 each species, we only conserved TE annotations i) classified as retro-transposons or DNA-
874 transposons, ii) longer than 250 bp, iii) sharing more than 85% identity with their consensus
875 sequence, iv) covering more than 33% of their consensus sequence length, v) first aligning
876 with their consensus sequence in a BLAST analysis against the TE-consensus library, and
877 vi) not overlapping with other annotations. TE annotations respecting all the described
878 criterion were referred to as canonical TE annotations.

879

880 Putative transposition machinery identification (*M. incognita* only)

881 We analysed the *M. incognita* predicted proteome and transcriptome (Blanc-Mathieu et al.
882 2017) and crossed the obtained information with the canonical TE-annotation to identify TE
883 containing genes putatively involved in the transposition machinery and evaluate TE-related
884 gene expression levels in comparison to the rest of the genes in the genome.

885

886 Finding genes coding for proteins with TE-related HMM profiles

887 We performed an exhaustive HMMprofile search analysis on the whole *M. incognita*
888 predicted proteome and then looked for proteins with TE-related domains. First, we
889 concatenated two HMMprofile libraries into one: Pfram32 (Finn et al. 2016) library and
890 Gypsy DB 2.0 (Llorens et al. 2011), a curated library of HMMprofiles linked to viruses, mobile
891 genetic elements, and genomic repeats. Then, using this concatenated HMM profile library,
892 we performed an exhaustive but stringent HMM profile search on the *M. incognita* proteome
893 using hmmscan (-E 0.00001 --domE 0.001 --noali).

894 Eventually, using in house script (see supplementary material 1), we selected the best non-
895 overlapping HMM profiles for each protein and then tagged corresponding genes with TE-
896 related HMM profiles thanks to a knowledge-based function from the REPET tool

897 'profileDB4Repet.py'. We kept as genes with TE-related profiles all the genes with at least
898 one TE-related HMM-profile identified.

899

900

901

902 Genes expression level

903 To determine the *M. incognita* protein-coding genes expression patterns, we used data from
904 a previously published life-stage specific RNA-seq analysis of *M. incognita* transcriptome
905 during tomato plant infection (Blanc-Mathieu et al. 2017). This analysis encompassed four
906 different life stages: (i) eggs, (ii) pre-parasitic second stage juveniles (J2), (iii) a mix of late
907 parasitic J2, third stage (J3) and fourth stage (J4) juveniles and (iv) adult females, all
908 sequenced in triplicates.

909 The cleaned RNA-seq reads were retrieved from the previous analysis and re-mapped to the
910 *M. incognita* annotated genome assembly (Blanc-Mathieu et al. 2017) using a more recent
911 version of STAR (2.6.1) (Dobin et al. 2013) and the more stringent end-to-end option (i.e. no
912 soft clipping) in 2-passes. Expected read counts were calculated on the predicted genes
913 from the *M. incognita* GFF annotation as FPKM values using RSEM (Li and Dewey 2011) to
914 take into account the multi-mapped reads via expectation maximization. To reduce
915 amplitude of variations, raw FPKM values were transformed to Log₁₀(FPKM+1) and the
916 median value over the 3 replicates was kept as a representative value in each life stage. The
917 expression data are available at
918 <https://data.inra.fr/dataset.xhtml?persistentId=doi:10.15454/YM2DHE> .

919 Then, for each life stage independently, i) we ranked the gene expression values, and ii)
920 defined gene expression level corresponding to the gene position in the ranking. We
921 considered as substantially expressed all the genes that presented an expression level \geq
922 1st quartile in at least one life stage.

923

924 TE annotations with potential transposition machinery

925 To identify TE-annotations including predicted genes involved in transposition machinery
926 (inclusion \geq 95% of the gene length), we performed the intersection of the canonical TE
927 annotation and the genes annotation BED files (see supplementary material 1) using the
928 intersect tool (-wo -s -F 0.95) from the bedtools v-2.27.1 suite (Quinlan and Hall 2010).

929 We then cross-referenced the obtained file with the list of the substantially expressed genes
930 and the list of the TE-related genes previously elaborated to identify the TEs containing
931 potential transposition machinery genes and their expression levels.

932

933 Evaluation of TE presence frequencies across the different *M. incognita*
934 isolates

935 We used the popoolationTE2 v-1.10.04 pipeline (Kofler et al. 2016) to compute isolate-
936 related support frequencies of both annotated, and *de novo* TE-loci across the 12 *M.*
937 *incognita* geographical isolates previously described. To that end, we performed a 'joint'
938 analysis as recommended by the popoolationTE2 manual. Briefly, popoolationTE2 uses both
939 quantitative and qualitative information extracted from paired-end (PE) reads mapping on the
940 TE-annotated reference genome and a set of reference TE sequences to detect signatures
941 of TE polymorphisms and estimate their frequencies in every analysed isolate. Frequency

942 values correspond to the proportion of individuals in an isolate for which a copy of the TE is
943 present at a given locus.

944

945

946

947

948 Preparatory work: creating the TE-hierarchy and the TE-merged-reference files.

949 We used the canonical TE-annotation set created above (see supplementary material 3) and
950 the *M. incognita* reference genome to produce the TE-merged reference file and the TE-
951 hierarchy file necessary to perform the popoolationTE analysis.

952 We used getfasta and maskfasta commands (default parameters) from the bedtools suite to
953 respectively extract and mask the sequences corresponding to canonical TE-annotations in
954 the reference genome. Then we concatenated both resulting sequences in a 'TE-merged
955 reference' multi fasta file. The 'TE-hierarchy' file was created from the TE-annotation file from
956 which it retrieves and stores the TE sequence name, the family, and the TE-order for every
957 entry.

958

959 Reads mapping

960 For each *M. incognita* isolate library, we mapped forward and reverse reads separately on
961 the "TE-merged-references" genome-TE file using the local alignment algorithm bwa bwasmw
962 v-0.7.17-r1188 (Li and Durbin 2009) with the default parameters. The obtained sam
963 alignment files were then converted to bam files using samtools view v-1.2 (Li et al. 2009).

964

965 Restoring paired-end information and generating the ppileup file.

966 We restored paired-end information from the previous separate mapping using the sep2pe (-
967 -sort) tool from popoolationTE2. Then, we created the ppileup file using the 'ppileup' tool
968 from popoolationTE2 with a map quality threshold of 15 (--map-qual 15).

969 For every base of the genome, this file summarises the number of PE reads inserts
970 spanning the position (physical coverage) but also the structural status inferred from paired-
971 end read covering this site.

972

973 Estimating target coverage and subsampling the ppileup to a uniform coverage

974 As noticed by R. Kofler, heterogeneity in physical coverage between populations may lead to
975 discrepancies in TE frequency estimation. Hence, we flattened the physical coverage across
976 the *M. incognita* isolates by a subsampling and a rescaling approach.

977 We first estimated the optimal target coverage to balance information loss and homogeneity
978 using the 'stats-coverage' tool from PopoolationTE2 (default parameter) and set this value to
979 15X. We then used the 'subsamplePpileup' tool (--target-coverage 15) to discard positions
980 with a physical coverage below 15X and rescale the coverage of the remaining position to
981 that value.

982

983 Identify signatures of TE polymorphisms

984 We identified signatures of TE polymorphisms from the previously subsampled file using the
985 'identifySignature' tool following the joint algorithm (--mode joint; --min-count 2; --signature-
986 window minimumSampleMedian; --min-valley minimumSampleMedian).

987 Then, for each identified site, we estimated TE frequencies in each isolate using the
988 'frequency' tool (default parameters). Eventually, we paired up the signatures of TE

989 polymorphisms using 'pairupSignatures' tool (--min-distance -200; --max-distance -- 300 as
990 recommended by R. Kofler), yielding a final list of potential TE-polymorphisms positions in
991 the reference genome with their associated frequencies for each one of the isolates.

992 Evaluation of PopoolationTE2 systematic error rate in the TE-frequency estimation.

993 To estimate PopoolationTE2 systematic error rate in the TE-frequency estimation, we ran
994 the same analysis (from the PE information restoration step) but comparing each isolate
995 against itself (12 distinct analyses).

996 We then analysed each output individually, measuring the frequency difference between the
997 two 'replicates' in all the detected loci with FR signatures (see below for more explanations).

998 We tested the homogeneity of the frequency-difference across the 12 analyses with an
999 ANOVA and concluded that the mean values of the frequencies differences between the
1000 analysis were not significantly heterogeneous (p. value = 0.102 > 0.05). Hence, we
1001 concatenated the 12 analysis frequency-difference and set the systematic error rate in the
1002 TE-frequency estimation to 2 times the standard deviation of the frequency differences, a
1003 value of 0.97 % .
1004

1005 TE polymorphism analysis

1006 Isolating TE loci with frequency variation across *M. incognita* isolates.

1007 We parsed PopoolationTE2 analysis output to identify TE loci with enough evidence to
1008 characterise them as polymorphic in frequency across the isolates.

1009 PopoolationTE2 output informs for each detected locus i) its position on the reference
1010 genome, ii) its frequency value for every sample of the analysis (e.g each isolate), and iii)
1011 qualitative information about the reads mapping signatures supporting a TE insertion.

1012 In opposition to separate Forward ('F') or Reverse ('R') signatures, 'FR' signatures mean the
1013 locus both boundaries are supported by significant physical coverage. Entries with such type
1014 of signature are more accurate in terms of frequency and position estimation. Hence, we
1015 only retained candidate loci with 'FR' signatures. Then, for each locus, we computed the
1016 maximal frequency variation between all the isolates and discarded the loci with a frequency
1017 difference smaller than the PopoolationTE2 systematic error rate in the TE-frequency
1018 estimation we computed (0.97 %; see above). We also discarded loci where different TEs
1019 were predicted to be inserted. We considered the remaining loci as polymorphic in frequency
1020 across the isolates.

1021

1022 Isolates phylogeny

1023 We reconstructed *M. incognita* isolates phylogeny according to their patterns of
1024 polymorphism in TE frequencies.

1025 We first computed a euclidean distance matrix from the isolates TE frequencies of all the
1026 detected polymorphic loci. We then used the distance matrix to construct the phylogenetic
1027 tree using the Neighbor Joining (NJ) method (R' phangorn package v-2.5.5). We computed
1028 nodes support values with a bootstrap approach (n=500 replicates). We compared the
1029 resulting tree with the topology described in (Koutsovoulos et al. 2020) using ItoI v-4.0
1030 viewer (Letunic and Bork 2019).

1031 Polymorphisms characterisation.

1032 We exported the polymorphic TE positions as an annotation file, and we used bedtools
1033 intersect (-wao) to perform their intersection with the reference canonical TE annotation. We
1034 then cross-referenced the results with the filtered popoolationTE2 output and defined a
1035 decision tree to characterise the TE-polymorphism detected by popoolationTE2 as
1036 'reference-TE polymorphism' (ref-polymorphism), 'extra-detection', or 'neo-insertion'.

1037 We considered a reference TE-annotation as polymorphic (e.g. ref-polymorphism locus) if:

1038 i) The position of the polymorphism predicted by PoPoolationTE2 falls between the
1039 boundaries of the reference TE-annotation

1040 ii) Both the reference TE-annotation and the predicted polymorphism belong to the same TE-
1041 consensus sequence.

1042 iii) The TE has a predicted frequency > 75% in the reference isolate Morelos.

1043 Canonical TE-annotations that did not intersect with polymorphic loci predicted by
1044 PopoolationTE2 were considered as non-polymorphic.

1045 We classified as 'neo-insertions' all the polymorphic loci for which no canonical TE was
1046 predicted in the reference annotation (polymorphism position is not included in a reference
1047 TE-annotation), but which were detected with a frequency > 25% in at least one isolate
1048 different from the reference isolate Morelos, in which the TE frequency should be inferior to
1049 1% and thus considered truly absent in the reference genome.

1050 Finally, we classified as 'extra-detection' all the polymorphic loci which did not correspond to
1051 a reference annotation but which were detected with a frequency > 25% in the reference
1052 isolate Morelos (at least). Polymorphic loci having a frequency between 1% and 25% in
1053 Morelos isolate were considered ambiguous and were discarded.

1054 Then, for each TE polymorphism, we investigated the homogeneity of the TE frequency
1055 between the isolates. We considered TE frequency was homogeneous between isolates
1056 when the maximum frequency variation between isolate was <= to 25%. Above this value,
1057 we considered the TE presence frequency was heterogeneous between isolates.

1058

1059

1060

1061 Highly Contrasted Polymorphic TE loci (HCPTEs): isolation, 1062 characterisation and experimental validation.

1063 HCPTEs isolation

1064 We considered as highly contrasted all the polymorphic loci for which i) all the isolates had
1065 frequency values either < 25% or > 75%, ii) at least one isolate showed a frequency < 25 %
1066 while another presented a frequency > 75%. Polymorphic loci fitting with these requirements
1067 were exported as an annotation file in the bed format.

1068

1069 HCPTEs possible functional impact

1070 We first identified the genes potentially impacted by the HCPTEs by cross-referencing the
1071 HCPTEs annotation file with the gene annotation file, using the bedtools suite. We used the
1072 'closest' program (-D b -fu -io; b being the gene annotation file) to identify the closest (but not
1073 intersecting) gene downstream each HCPTe. We only retained the entries with a maximum
1074 distance of 1 kb between the HCPTe and gene boundaries. We identified the insertions in
1075 the gene using the 'intersect' tool (-wo).

1076

1077 Then, we performed a manual bioinformatic functional analysis for each gene potentially
1078 impacted by HCPTEs. Protein sequences were extracted from the *M. incognita* predicted
1079 proteome (Blanc-Mathieu et al. 2017) and blasted (blastp; default parameters) against the
1080 Non-Redundant protein sequences database (NR) from the NCBI
1081 (<https://blast.ncbi.nlm.nih.gov/>). The same sequences were also used on the InterProScan
1082 website (<https://www.ebi.ac.uk/interpro/>) to perform an extensive search on all the available
1083 libraries of conserved protein domains and motifs.

1084 Then, for each gene potentially impacted by HCPTEs, we performed an orthology search on
1085 the Wormbase Parasite website (<https://parasite.wormbase.org/>) using genes accession
1086 numbers and the pre-computed ENSEMBL Compara orthology prediction (Herrero et al.
1087 2016).

1088 Finally, we analysed the expression levels of the genes potentially impacted by HCPTEs
1089 extracting the information from the RNA-seq analysis of four *M. incognita* life-stages
1090 performed previously (see Putative transposition machinery identification section).
1091

1092 Experimental validation of Highly Contrasted Polymorphic TE loci

1093 To experimentally validate in-silico predictions of TE neo-insertions with potential functional
1094 impact, we selected 5 candidates among the HCPTEs loci and performed a PCR
1095 experiment. To run this experiment, we used DNA remaining from extractions performed on
1096 the *M. incognita* isolates for a previous population genomics analysis (Koutsovoulos et al.
1097 2020).

1098 **Primer design and PCR amplification.**

1099 We designed primers for the PCR analysis using the Primer3Plus web interface
1100 (Untergasser et al. 2007). The set of 10 primers with the corresponding sequence and
1101 expected amplicon sizes with, or without TE insertion, is shown in (sup. Table 7 &
1102 supplementary material 4). We used primers amplifying the whole actin-encoding gene
1103 (Minc3s00960g19311) as positive control.

1104 PCR experiments were performed on *M. incognita* Morelos isolate and 11 Brazilian isolates:
1105 R1-2, R1-3, R1-6, R2-1, R2-6, R3-1, R3-2, R3-4, R4-1, R4-3 and R4-4.

1106 R3-1 presented no amplification in any of the tested loci nor the positive control (actin) and
1107 was thus discarded from this analysis.

1108 PCR mixture contained 0.5µmol of each primer, 1x MyTaq™ reaction buffer and 1.0 U of
1109 MyTaq™ DNA polymerase (Bioline Meridian Bioscience) adjusted to a total volume of 20µL.
1110 PCR amplification was performed with a TurboCycler2 (Blue-Ray Biotech Corp.). PCR
1111 conditions were as follows: initial denaturation at 95°C for 5 min, followed by 35 cycles of
1112 95°C for 30 s, 56°C for 30 s of annealing, and 72°C for 3 min of extension, the program
1113 ending with a final extension at 72°C for 10 min. Aliquots of 5µL were migrated by
1114 electrophoresis on a 1% agarose gel (Sigma Chemical Co.) for 70 min at 100 V. The size
1115 marker used is 1kb Plus DNA Ladder (New England Biolabs Inc.), containing the following
1116 size fragments in bp: 100, 200, 300, 400, 500, 600, 700, 900, 1000, 1200, 1500, 2000, 3000,
1117 4000, 5000, 6000, 8000 and 10000.

1118 **Purification and sequencing of PCR amplicons.**

1119 Amplicon bands were revealed using ethidium bromide and exposure to ultraviolet radiation.
1120 PCR products bands were excised from the agarose gel with a scalpel and purified using
1121 MinElute Gel Extraction Kit (Qiagen) before sequencing, following the manufacturer's
1122 protocol. PCR products were sequenced by Sanger Sequencing (Eurofins Genomics).

1123 Forward (F) and Reverse (R) sequences were blasted individually
1124 (<https://blast.ncbi.nlm.nih.gov/> ; Optimised for 'Somewhat similar sequences', default

1125 parameters) to the expected TE-consensus sequence and to the genomic region
1126 surrounding the predicted insertion point (2 kb region: 1kb upstream the predicted insertion
1127 point and 1kb downstream). When no significant hit was found, the sequence was blasted
1128 against the *Meloidogyne* reference genomes available (<https://meloidogyne.inrae.fr/>), the
1129 whole TE-consensus library, and the NR database on the NCBI blast website.
1130

1131 Acknowledgements

1132 DKLK would like to thank Joffrey Mejias and Georgios Koutsovoulos for all the pieces of
1133 advice they gave in their respective fields and for all the inspiring discussion. The authors
1134 would like to thank Erika VS Albuquerque for her help and assistance in accessing the DNA
1135 extractions from the *M. incognita* Brazilian Isolates. We would also like to thank the BIG
1136 bioinformatics platform from the PlantBios infrastructure as well as the URGI team for
1137 providing facilities and technical support.

1138 References

- 1139 Abad P, Gouzy J, Aury J-M, Castagnone-Sereno P, Danchin EGJ, Deleury E, Perfus-
1140 Barbeoch L, Anthouard V, Artiguenave F, Blok VC, et al. 2008. Genome sequence of
1141 the metazoan plant-parasitic nematode *Meloidogyne incognita*. *Nat. Biotechnol.*
1142 26:909–915.
- 1143 Agrios GN. 2005. *Plant Pathology*, 5th Edition. Burlington, USA: Elsevier Academic Press
- 1144 Barbary A, Djian-Caporalino C, Palloix A, Castagnone-Sereno P. 2015. Host genetic
1145 resistance to root-knot nematodes, *Meloidogyne* spp., in Solanaceae: from genes to
1146 the field. *Pest Manag. Sci.* 71:1591–1598.
- 1147 Bast J, Jaron KS, Schuseil D, Roze D, Schwander T. 2019. Asexual reproduction reduces
1148 transposable element load in experimental yeast populations. Coop G, Tautz D, Coop
1149 G, Charlesworth B, editors. *eLife* 8:e48548.
- 1150 Bast J, Schaefer I, Schwander T, Maraun M, Scheu S, Kraaijeveld K. 2015. No
1151 Accumulation of Transposable Elements in Asexual Arthropods. *Mol. Biol.*
1152 *Evol.*:msv261.
- 1153 Belyayev A. 2014. Bursts of transposable elements as an evolutionary driving force. *J. Evol.*
1154 *Biol.* 27:2573–2584.
- 1155 Bessereau J-L. 2006. Transposons in *C. elegans*. *WormBook* [Internet]. Available from:
1156 http://www.wormbook.org/chapters/www_transposons/transposons.html
- 1157 Blanc-Mathieu R, Perfus-Barbeoch L, Aury J-M, Rocha MD, Gouzy J, Sallet E, Martin-
1158 Jimenez C, Bailly-Bechet M, Castagnone-Sereno P, Flot J-F, et al. 2017.
1159 Hybridization and polyploidy enable genomic plasticity without sex in the most
1160 devastating plant-parasitic nematodes. *PLOS Genet.* 13:e1006777.
- 1161 Bourgeois Y, Boissinot S. 2019. On the Population Dynamics of Junk: A Review on the
1162 Population Genomics of Transposable Elements. *Genes* 10:419.
- 1163 Castagnone-Sereno P. 2006. Genetic variability and adaptive evolution in parthenogenetic
1164 root-knot nematodes. *Heredity* 96:282–289.
- 1165 Castagnone-Sereno P, Danchin EGJ. 2014. Parasitic success without sex – the nematode
1166 experience. *J. Evol. Biol.* 27:1323–1333.
- 1167 Castagnone-Sereno P, Mulet K, Danchin EGJ, Koutsovoulos GD, Karaulic M, Rocha MD,
1168 Bailly-Bechet M, Pratz L, Perfus-Barbeoch L, Abad P. 2019. Gene copy number
1169 variations as signatures of adaptive evolution in the parthenogenetic, plant-parasitic

- 1170 nematode *Meloidogyne incognita*. *Mol. Ecol.* 28:2559–2572.
- 1171 Castagnone-Sereno P, Wajnberg E, Bongiovanni M, Leroy F, Dalmasso A. 1994. Genetic
1172 variation in *Meloidogyne incognita* virulence against the tomato Mi resistance gene:
1173 evidence from isofemale line selection studies. *Theor. Appl. Genet.* 88:749–753.
- 1174 Faino L, Seidl MF, Shi-Kunne X, Pauper M, Berg GCM van den, Wittenberg AHJ, Thomma
1175 BPHJ. 2016. Transposons passively and actively contribute to evolution of the two-
1176 speed genome of a fungal pathogen. *Genome Res.* [Internet]. Available from:
1177 <http://genome.cshlp.org/content/early/2016/07/12/gr.204974.116>
- 1178 Finn RD, Coghill P, Eberhardt RY, Eddy SR, Mistry J, Mitchell AL, Potter SC, Punta M,
1179 Qureshi M, Sangrador-Vegas A, et al. 2016. The Pfam protein families database:
1180 towards a more sustainable future. *Nucleic Acids Res.* 44:D279–D285.
- 1181 Flutre T, Duprat E, Feuillet C, Quesneville H. 2011. Considering Transposable Element
1182 Diversification in De Novo Annotation Approaches. *PLoS ONE* 6:e16526.
- 1183 Glémin S, François CM, Galtier N. 2019. Genome Evolution in Outcrossing vs. Selfing vs.
1184 Asexual Species. In: Anisimova M, editor. *Evolutionary Genomics: Statistical and
1185 Computational Methods. Methods in Molecular Biology.* New York, NY: Springer. p.
1186 331–369. Available from: https://doi.org/10.1007/978-1-4939-9074-0_11
- 1187 Glémin S, Galtier N. 2012. Genome Evolution in Outcrossing Versus Selfing Versus Asexual
1188 Species. In: Anisimova M, editor. *Evolutionary Genomics.* Vol. 855. Totowa, NJ:
1189 Humana Press. p. 311–335. Available from: http://link.springer.com/10.1007/978-1-61779-582-4_11
- 1190
- 1191 Gross SM, Williamson VM. 2011. Tm1: A Mutator/Foldback Transposable Element Family in
1192 Root-Knot Nematodes. *PLoS ONE* 6:e24534.
- 1193 Guerreiro MPG. 2014. Interspecific hybridization as a genomic stressor inducing mobilization
1194 of transposable elements in *Drosophila*. *Mob. Genet. Elem.* 4:e34394.
- 1195 Herrero J, Muffato M, Beal K, Fitzgerald S, Gordon L, Pignatelli M, Vilella AJ, Searle SMJ,
1196 Amode R, Brent S, et al. 2016. Ensembl comparative genomics resources. Database
1197 [Internet] 2016. Available from:
1198 <https://academic.oup.com/database/article/doi/10.1093/database/bav096/2630091>
- 1199 Hill WG, Robertson A. 1966. The effect of linkage on limits to artificial selection. *Genet. Res.*
1200 8:269–294.
- 1201 Hoffmann AA, Reynolds KT, Nash MA, Weeks AR. 2008. A high incidence of
1202 parthenogenesis in agricultural pests. *Proc. R. Soc. Lond. B Biol. Sci.* 275:2473–
1203 2481.
- 1204 Jones JT, Haegeman A, Danchin EGJ, Gaur HS, Helder J, Jones MGK, Kikuchi T,
1205 Manzanilla-López R, Palomares-Rius JE, Wesemael WML, et al. 2013. Top 10 plant-
1206 parasitic nematodes in molecular plant pathology. *Mol. Plant Pathol.* 14:946–961.
- 1207 Kamath RS, Fraser AG, Dong Y, Poulin G, Durbin R, Gotta M, Kanapin A, Le Bot N, Moreno
1208 S, Sohrmann M, et al. 2003. Systematic functional analysis of the *Caenorhabditis*
1209 elegans genome using RNAi. *Nature* 421:231–237.
- 1210 Katju V, Bergthorsson U. 2013. Copy-number changes in evolution: rates, fitness effects and
1211 adaptive significance. *Front. Genet.* [Internet] 4. Available from:
1212 http://www.frontiersin.org/Evolutionary_and_Population_Genetics/10.3389/fgene.2013.00273/abstract
- 1213
- 1214 Kofler R, Gómez-Sánchez D, Schlötterer C. 2016. PoPoolationTE2: Comparative Population
1215 Genomics of Transposable Elements Using Pool-Seq. *Mol. Biol. Evol.* 33:2759–2764.
- 1216 Kondrashov AS. 1988. Deleterious mutations and the evolution of sexual reproduction.
1217 *Nature* 336:435–440.
- 1218 Koutsovoulos GD, Marques E, Arguel M-J, Duret L, Machado ACZ, Carneiro RMDG,
1219 Kozłowski DK, Bailly Bechet M, Castagnone Sereno P, Albuquerque EVS, et al.
1220 2020. Population genomics supports clonal reproduction and multiple independent
1221 gains and losses of parasitic abilities in the most devastating nematode pest. *Evol.*
1222 *Appl.* 13:442–457.
- 1223 Koutsovoulos GD, Pouillet M, Ashry AE, Kozłowski DK, Sallet E, Rocha MD, Martin-Jimenez
1224 C, Perfus-Barbeoch L, Frey J-E, Ahrens C, et al. 2019. The polyploid genome of the

- 1225 mitotic parthenogenetic root knot nematode *Meloidogyne enterolobii*.
1226 bioRxiv:586818.
- 1227 Kraaijeveld K, Zwanenburg B, Hubert B, Vieira C, De Pater S, Van Alphen JJM, Den Dunnen
1228 JT, De Kniiff P. 2012. Transposon proliferation in an asexual parasitoid. *Mol. Ecol.*
1229 21:3898–3906.
- 1230 Laricchia KM, Zdraljevic S, Cook DE, Andersen EC. 2017. Natural Variation in the
1231 Distribution and Abundance of Transposable Elements Across the *Caenorhabditis*
1232 *elegans* Species. *Mol. Biol. Evol.* 34:2187–2202.
- 1233 Lerat E, Goubert C, Guirao-Rico S, Merenciano M, Dufour A-B, Vieira C, González J. 2019.
1234 Population-specific dynamics and selection patterns of transposable element
1235 insertions in European natural populations. *Mol. Ecol.* 28:1506–1522.
- 1236 Letunic I, Bork P. 2019. Interactive Tree Of Life (iTOL) v4: recent updates and new
1237 developments. *Nucleic Acids Res.* 47:W256–W259.
- 1238 Li H, Durbin R. 2009. Fast and accurate short read alignment with Burrows–Wheeler
1239 transform. *Bioinformatics* 25:1754–1760.
- 1240 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin
1241 R. 2009. The Sequence Alignment/Map format and SAMtools. *Bioinformatics*
1242 25:2078–2079.
- 1243 Lively CM. 2010. A Review of Red Queen Models for the Persistence of Obligate Sexual
1244 Reproduction. *J. Hered.* 101:S13–S20.
- 1245 Llorens C, Futami R, Covelli L, Domínguez-Escribá L, Viu JM, Tamarit D, Aguilar-Rodríguez
1246 J, Vicente-Ripolles M, Fuster G, Bernet GP, et al. 2011. The Gypsy Database
1247 (GyDB) of mobile genetic elements: release 2.0. *Nucleic Acids Res.* 39:D70–D74.
- 1248 Lu L, Chen J, Robb SMC, Okumoto Y, Stajich JE, Wessler SR. 2017. Tracking the genome-
1249 wide outcomes of a transposable element burst over decades of amplification. *Proc.*
1250 *Natl. Acad. Sci.* 114:E10550–E10559.
- 1251 Martin M. 2011. Cutadapt removes adapter sequences from high-throughput sequencing
1252 reads. *EMBnet.journal* 17:10–12.
- 1253 McCarter JP. 2009. Molecular Approaches Toward Resistance to Plant-Parasitic
1254 Nematodes. In: Berg RH, Taylor CG, editors. *Cell Biology of Plant Nematode*
1255 *Parasitism*. Vol. 15. *Plant Cell Monographs*. Berlin, Heidelberg: Springer Berlin
1256 Heidelberg. p. 239–267. Available from:
1257 http://www.springerlink.com/index/10.1007/978-3-540-85215-5_9
- 1258 McClintock B. 1984. The significance of responses of the genome to challenge. *Science*
1259 226:792–801.
- 1260 Mhiri C, Parisod C, Daniel J, Petit M, Lim KY, Borne FD de, Kovarik A, Leitch AR,
1261 Grandbastien M-A. 2019. Parental transposable element loads influence their
1262 dynamics in young *Nicotiana* hybrids and allotetraploids. *New Phytol.* 221:1619–
1263 1633.
- 1264 Muller HJ. 1964. The Relation of Recombination to Mutational Advance. *Mutat Res* 106:2–9.
- 1265 Opperman CH, Bird DM, Williamson VM, Rokhsar DS, Burke M, Cohn J, Cromer J, Diener
1266 S, Gajan J, Graham S, et al. 2008. Sequence and genetic map of *Meloidogyne*
1267 *hapla*: A compact nematode genome for plant parasitism. *Proc Natl Acad Sci U A*
1268 105:14802–14807.
- 1269 Quesneville H, Bergman CM, Andrieu O, Autard D, Nouaud D, Ashburner M, Anxolabehere
1270 D. 2005. Combined evidence annotation of transposable elements in genome
1271 sequences. *PLoS Comput Biol* 1:166–175.
- 1272 Quinlan AR, Hall IM. 2010. BEDTools: a flexible suite of utilities for comparing genomic
1273 features. *Bioinformatics* 26:841–842.
- 1274 Rice WR. 2002. Experimental tests of the adaptive significance of sexual recombination.
1275 *Nat. Rev. Genet.* 3:241–251.
- 1276 Rodriguez F, Arkhipova IR. 2018. Transposable elements and polyploid evolution in animals.
1277 *Curr. Opin. Genet. Dev.* 49:115–123.
- 1278 Savary S, Willocquet L, Pethybridge SJ, Esker P, McRoberts N, Nelson A. 2019. The global
1279 burden of pathogens and pests on major food crops. *Nat. Ecol. Evol.* 3:430–439.

- 1280 Somvanshi VS, Tathode M, Shukla RN, Rao U. 2018. Nematode Genome Announcement: A
1281 Draft Genome for Rice Root-Knot Nematode, *Meloidogyne graminicola*. *J. Nematol.*
1282 50:111–116.
- 1283 Susič N, Koutsovoulos GD, Riccio C, Danchin EGJ, Blaxter ML, Lunt DH, Strajnar P, Širca
1284 S, Urek G, Stare BG. 2020. Genome sequence of the root-knot nematode
1285 *Meloidogyne luci*. *J. Nematol.* 52:1–5.
- 1286 Szitenberg A, Cha S, Opperman CH, Bird DM, Blaxter ML, Lunt DH. 2016. Genetic drift, not
1287 life history or RNAi, determine long term evolution of transposable elements.
1288 *Genome Biol. Evol.*:evw208.
- 1289 Szitenberg A, Salazar-Jaramillo L, Blok VC, Laetsch DR, Joseph S, Williamson VM, Blaxter
1290 ML, Lunt DH. 2017. Comparative Genomics of Apomictic Root-Knot Nematodes:
1291 Hybridization, Ploidy, and Dynamic Genome Change. *Genome Biol. Evol.* 9:2844–
1292 2861.
- 1293 The *C. elegans* Genome Sequencing Consortium. 1998. Genome sequence of the
1294 nematode *C. elegans*: a platform for investigating biology. *Science* 282:2012–2018.
- 1295 Trudgill DL, Blok VC. 2001. Apomictic, polyphagous root-knot nematodes: exceptionally
1296 successful and damaging biotrophic root pathogens. *Annu Rev Phytopathol* 39:53–
1297 77.
- 1298 Untergasser A, Nijveen H, Rao X, Bisseling T, Geurts R, Leunissen JAM. 2007.
1299 Primer3Plus, an enhanced web interface to Primer3. *Nucleic Acids Res.* 35:W71–
1300 W74.
- 1301 Vicient CM, Casacuberta JM. 2017. Impact of transposable elements on polyploid plant
1302 genomes. *Ann. Bot.* 120:195–207.
- 1303 Vrijenhoek RC, Parker ED. 2009. Geographical Parthenogenesis: General Purpose
1304 Genotypes and Frozen Niche Variation. In: Schön I, Martens K, Dijk P, editors. *Lost*
1305 *Sex*. Dordrecht: Springer Netherlands. p. 99–131. Available from:
1306 http://www.springerlink.com/index/10.1007/978-90-481-2770-2_6
- 1307 Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, Flavell A, Leroy P,
1308 Morgante M, Panaud O, et al. 2007. A unified classification system for eukaryotic
1309 transposable elements. *Nat. Rev. Genet.* 8:973–982.
- 1310 Wright S, Finnegan D. 2001. Genome evolution: Sex and the transposable element. *Curr.*
1311 *Biol.* 11:R296–R299.