

Research Report

Genome-wide analysis of long terminal repeat retrotransposons from the cranberry *Vaccinium macrocarpon*

Nusrat Sultana^{a,b*}, Gerhard Menzel^b, Kathrin M. Seibt^b, Sònia Garcia^c, Beatrice Weber^b, Sedat Serçe^d, and Tony Heitkam^{b*}

^aDepartment of Botany, Faculty of Life and Earth Sciences, Jagannath University, Dhaka 1100, Bangladesh

^bInstitute of Botany, Technische Universität Dresden, D-01062 Dresden, Germany (present address)

^cInstitut Botànic de Barcelona (IBB-CSIC), 08038 Barcelona, Catalonia, Spain

^dDepartment of Agricultural Genetic Engineering, Ayhan Şahenk Faculty of Agricultural Sciences and Technologies, Niğde Ömer Halisdemir University, 51240, Niğde, Turkey

Running title: Retrotransposons in cranberry

Words: 7057 (main text) + 2650 (references) + 555 (figure legends)

***Corresponding authors:**

nusrat.sultana1@mailbox.tu-dresden.de, tony.heitkam@tu-dresden.de

1 ABSTRACT

2 **BACKGROUND:** Long terminal repeat (LTR) retrotransposons are widespread in plant
3 genomes and play a large role in the generation of genomic variation. Despite this, their
4 identification and characterization remains challenging, especially for non-model genomes.
5 Hence, LTR retrotransposons remain undercharacterized in *Vaccinium* genomes, although
6 they may be beneficial for current berry breeding efforts.

7 **OBJECTIVE:** Exemplarily focusing on the genome of American cranberry (*Vaccinium*
8 *macrocarpon* Aiton), we aim to generate an overview of the LTR retrotransposon landscape,
9 highlighting the abundance, transcriptional activity, sequence, and structure of the major
10 retrotransposon lineages.

11 **METHODS:** Graph-based clustering of whole genome shotgun Illumina reads was
12 performed to identify the most abundant LTR retrotransposons and to reconstruct
13 representative *in silico* full-length elements. To generate insights into the LTR
14 retrotransposon diversity in *V. macrocarpon*, we also queried the genome assembly for
15 presence of reverse transcriptases (RTs), the key domain of LTR retrotransposons.
16 Using transcriptomic data, transcriptional activity of retrotransposons corresponding to
17 the consensus was analyzed.

18 **RESULTS:** We provide an in-depth characterization of the LTR retrotransposon landscape
19 in the *V. macrocarpon* genome. Based on 475 RTs harvested from the genome assembly, we
20 detect a high retrotransposon variety, with all major lineages present. To better understand
21 their structural hallmarks, we reconstructed 26 Ty1-*copia* and 28 Ty3-*gypsy in silico*
22 consensus that capture the detected diversity. Accordingly, we frequently identify
23 association with tandemly repeated motifs, extra open reading frames, and specialized,
24 lineage-typical domains. Based on the overall high genomic abundance and transcriptional

25 activity, we suggest that retrotransposons of the Ale and Athila lineages are most promising
26 to monitor retrotransposon-derived polymorphisms across accessions.

27 **CONCLUSIONS:** We conclude that LTR retrotransposons are major components of the *V.*
28 *macrocarpon* genome. The representative consensus sequences provide an entry point for further
29 *Vaccinium* genome analyses and may be applied to derive molecular markers for enhancing
30 cranberry selection and breeding.

31 **Keywords:** cranberry, *Vaccinium macrocarpon*, repetitive DNA, LTR retrotransposon, Ty1-
32 *copia*, Ty3-gypsy

33 INTRODUCTION

34 The genus *Vaccinium* L. belongs to the family Ericaceae Juss. comprising approx. 450
35 species distributed all over the world [1]. Available molecular phylogenies suggest that the
36 genus is monophyletic, whereas intrageneric relationships are more difficult to resolve [2, 3].
37 Dependent on different systematic treatments, the genus is divided into approximately 30
38 sections and subgenera [1-4]. *Vaccinium macrocarpon* Aiton, also known as the American
39 cranberry, is native to North America [5, 6] and is one of the most commonly cultivated
40 *Vaccinium* species, along with the highbush blueberry (*V. corymbosum* L.) [5]. The small
41 berry fruits produced from wild and cultivated species of the genus are rich in nutritious
42 secondary plant metabolites, including some beneficial for human health with anticancer,
43 antioxidant, antidiabetic, and many other properties [7, 8]. The American cranberry and the
44 European cranberry (a close relative, *V. oxycoccus* L.) belong to the same section or
45 subgenus, named *Oxycoccus* [1, 4]. Although these two species differ in geographical
46 distribution, they can produce fertile offspring upon hybridization [5, 6].

47 *Vaccinium* species have ploidy levels ranging from diploid to hexaploid, with a base
48 chromosome number $x = 12$ [5]. As cranberry (*V. macrocarpon*) is a diploid and self-fertile
49 species with a relatively small genome size (approximately 470 Mb), it has become a
50 valuable reference to study the genetics and genomics of the genus [6, 9-11]. So far, whole
51 genome sequence data, a reference genome assembly, and transcriptome datasets have been
52 published for the diploid cranberry and for both the diploid and the tetraploid blueberries
53 [11-13]. In addition, several high-throughput genetics, genomics, and transcriptomics
54 datasets as well as other important breeding information of different species of the genus
55 *Vaccinium* are either published or underway, accessible through the web platform
56 <https://www.vaccinium.org>.

57 Recent technological advancements allow tracing of genome dynamics on a fine-scale level.
58 Especially transposable elements (TEs) play a major role in genome evolution and contribute
59 largely to the generation of polymorphisms between genotypes [14-18]. Hence, TEs can also
60 be applied as molecular markers for the differentiation of accessions and to inform plant
61 breeding programs [19-23]. Thus, the correct annotation and characterization of repetitive
62 sequences, particularly TEs, are a prerequisite for an enriched genome assembly as well as
63 for subsequent genome characterization and valorization [10-12]. Nevertheless, as many
64 different TE classification hierarchies have been suggested [24-29], as an integration of
65 tools, databases, and pipelines is still discussed [30], and as there is still no gold standard for
66 TE identification and classification available, TE annotation and classification for non-model
67 species is still a considerable undertaking.

68 In plants, the most abundant TEs belong to the group of long terminal repeat (LTR)
69 retrotransposons contributing up to 80 % of nuclear DNA [31-33]. They consist of two LTRs
70 flanking the main protein-encoding genes *gag* and *pol* [25, 33, 34]. Whereas *gag* (encoding
71 the capsid and nucleocapsid) is considered as a single gene, the *pol* gene generally encodes
72 several domains, including a protease (PR), an RNase H (RNH), a reverse transcriptase
73 (RT), and an integrase (INT) [25, 33, 34]. In addition to the protein domains, two conserved
74 motifs are characteristic for these full-length elements: a primer binding site (PBS) and a
75 polypurine tract (PPT) [25, 31, 33, 34]. Based on the order of the protein domains and their
76 sequence, the LTR retrotransposons are further divided into the superfamilies Ty3-*gypsy* and
77 Ty1-*copia* [24, 28].

78 In *Vaccinium* genomes, we have already gained first insights into the repetitive DNA
79 composition, especially into that of satellite DNAs [35, 36]. For dispersed repeats, first
80 estimates were also brought forward, suggesting that nearly 40 % of the cranberry genome

81 was composed of TEs [10]. Very recently, a TE annotation was also made available for the
82 genome of tetraploid blueberry (*V. corymbosum*) [13]. However, the contribution of
83 individual TE families to these genomes is still not well-studied [10-12].

84 Here, we aim to complement the broad overview studies with a deep analysis of the most
85 abundant, individual TE families in the cranberry genome. The precise knowledge of the
86 LTR retrotransposon composition, as well as the underlying sequences and structures can
87 inform annotation of berry genomes and may provide support for berry breeding and
88 selection.

89

90 **MATERIALS AND METHODS**

91 **Graph-based clustering of repeat sequences of the *V. macrocarpon* genome**

92 To representatively survey the cranberry LTR retrotransposon landscape, we used publicly
93 available whole genome paired-end Illumina reads of *Vaccinium macrocarpon* (NCBI
94 BioProject PRJNA245813) of cranberry cultivar ‘Ben Lear’ (accession number CNJ99-125-
95 1) [10]. Before clustering, we pre-treated the reads by quality filtering, adapter trimming,
96 and processing of read length: Illumina Truseq adapters were removed using Trimmomatic
97 with the parameters ILLUMINACLIP:TruSeq3-PE-2.fa:2:30:10 [37]. Fastx_trimmer and
98 FASTx quality filter from FASTX-Toolkit (http://hannonlab.cshl.edu/fastx_toolkit) were
99 used to trim all reads to 150 bp length. Shorter sequences remaining after quality filtering
100 and trimming, were removed with seqtk (<http://github.com/lh3/seqtk>). This was followed by
101 interlacing of paired-ends with the FASTX-Toolkit. We then randomly selected 2× 1M, 2×
102 5M, and 2× 10M reads representing different genome coverages ranging from 0.02× to
103 2.04×. For read clustering, we used the Galaxy web interface of the RepeatExplorer
104 pipeline2 [38, 39] with the parameters -l 39 (minimal overlap of clustering) and -o 30

105 (minimum overlap for assembly), considering the different genome coverages. The
106 generated read clusters were further assigned to the retrotransposon lineages according to
107 their similarity to reference elements in REXdb [28] and by RepeatMasker
108 (<http://www.repeatmasker.org>, [40]). Based on the cluster graph shape, protein domains and
109 repeat masker hits, each cluster was manually annotated and assigned to superclusters from
110 the RepeatExplorer run of the highest genome coverage, which was expected to generate the
111 highest number of full-length elements.

112

113 **Construction and classification of full-length LTR retrotransposons**

114 To reconstruct full-length Ty1-*copia* and Ty3-*gypsy* retrotransposon sequences, several steps
115 were followed: First, clusters and superclusters belonging to the individual retrotransposon
116 types were identified based on their graph shape and protein domain hits. Second, contigs
117 from the identified clusters were imported into Geneious Prime 2019, hereafter called
118 Geneious (<http://www.geneious.com>, [41]). Third, we used the longest contig sequence as
119 reference and mapped the remaining contigs against it using the Geneious mapper tool (with
120 interation parameter = 10). We extracted the consensus sequence from the mapped contigs.
121 Afterwards, short reads from the respective cluster were mapped against the consensus
122 contig sequence to derive the final consensus sequence. If the identified LTR retrotransposon
123 sequence was represented by only a single RepeatExplorer cluster, the consensus sequence
124 was directly used for the reconstruction of the full-length element. However, if the LTR
125 retrotransposon was split across multiple clusters, additional steps were considered for the
126 reconstruction of the full-length element: These steps included iterative multiple sequence
127 alignments of the consensus sequences of all clusters belonging to a single supercluster in
128 Geneious. The parameters for the Geneious alignment tools were: automatically determine

129 direction, alignment type = global alignment, cost matrix = 65 % similarity (5.0/-4.0), gap
130 open penalty = 12, gap extension penalty = 3, refinement iterations = 2. Then, short read
131 sequences from the respective supercluster were mapped against the reconstructed consensus
132 sequence to derive the final full-length sequences using the Geneious mapper tool (with
133 iteration parameter = 10) within Geneious.

134 The reconstructed full-length sequences were searched for the typical structural features of
135 LTR retrotransposons, namely protein domain sequences, LTR regions, PBS, PPT, and
136 internal repetitions. Different software and databases were used for these purposes: the
137 LTR_Finder web server ([42]; http://tlife.fudan.edu.cn/tlife/ltr_finder/) was applied to
138 predict common structural features along the reconstructed LTR retrotransposons, using the
139 default parameters except “predict PBS by using tRNA database”, which was variable for
140 different elements. Dotplot analyses were performed using the default parameter of the
141 EMBOSS dottup tool integrated in Geneious [41, 43]. These dotplots showed the positions
142 of LTRs as short diagonals in the upper right and lower left corners. Internal tandemly
143 repeated structures lead to regular patterns of short lines in the dotplot. Furthermore, the
144 REXdb database underlying RepeatExplorer’s protein domain finder tool (default
145 parameters) was used to identify TE protein domains along the query elements and to predict
146 their position and direction [38, 39]. Finally, the initial results from the protein domain
147 finder tool were quality-filtered using the default parameters. For the Tandem Repeat Finder
148 (TRF) tool the following parameters were used: alignment (match, mismatch, indel=2,7,7;
149 2,5,7; 2,5,5; 2,3,5), minimum alignment score to report repeats (20-60), and minimum period
150 size (10-300) [44]. The TRF output was manually inspected, compared to the dotplot
151 findings, and recorded for different features like period size, consensus sequence, number of
152 repeat units, and position within the full-length LTR retrotransposon.

153 The online databases used to characterize and classify each full-length sequence were The
154 Gypsy Database (GyDB) [25], Repbase [27], and REXdb database [28]. Repbase was used
155 from the web-based software censor (<https://www.girinst.org/censor/index.php>) and GyDB
156 from the online server (http://gydb.org/index.php/Main_Page). The REXdb database is
157 integrated in RepeatExplorer2 (<https://galaxy-elixir.cerit-sc.cz/>).

158 The reconstructed full-length sequences were named according to the conventions and
159 deposited in the Repbase database. Detailed information is summarized in Supplementary
160 File 1 and Tables S1.

161

162 **Assignment of retrotransposon lineages and clades**

163 A comparison against a database of well-known protein domains serves as the basis of a
164 detailed retrotransposon classification. For this, protein sequences of reverse transcriptases
165 (RTs) and, if applicable, chromodomains (CHD) were extracted from the publicly available
166 genome assembly for *V. macrocarpon* (GenBank assembly accession: GCA_000775335.1)
167 using the protein domain finder tool DANTE implemented in the RepeatExplorer Galaxy
168 web-interface [38, 39]. DANTE was used with the quality filtering (minimum identity = 0.3;
169 minimum similarity: = 0.4; minimum alignment length = 0.8; interruptions = 3). The
170 program output 2149 sequences which were then clustered with a 90 % similarity threshold
171 using CD-HIT ([45] <http://weizhongli-lab.org/cd-hit>). RT and CHD amino acid sequences
172 were aligned with MAFFT [46] within Geneious (parameters: algorithm = auto [selects
173 appropriate strategy from L-INS-I, FFT-NS-i and FFT-NS-2 according to data size], scoring
174 matrix = “200PAM/K=2”, gap open penalty = 1.53, offsetvalue = 0.123). Pairwise distance
175 matrix from the RT alignment was calculated in Geneious for boxplot visualizations, which
176 were created with R graphics [47].

177 The resulting multiple sequence alignments of RT and CHD sequences were visualized as a
178 dendrogram using the approximate maximum likelihood estimation algorithm of the
179 FastTree tool [48] and the Randomized Axelerated Maximum Likelihood RAxML method
180 [49] implemented in Geneious (default parameters). This allowed assigning the detected
181 LTR retrotransposons to published lineages and clades according to their encoded RT and
182 CHD sequences.

183

184 **Analysis of *cis*-regulatory element associated with the 5' LTR of the full-length** 185 **LTR retrotransposons**

186 We extracted the 5' LTRs of the identified Ty1-*copia* and Ty3-*gypsy* full-length elements.
187 Ambiguous nucleotides from each of the 5' LTR sequences were replaced by an N using the
188 software “Sequence Manipulation Suite” [50]
189 (http://www.bioinformatics.org/sms2/filter_dna.html). For the identification of regulatory
190 sequence motifs, 5' LTR sequences were individually searched against the Plant Care
191 database [51] (<http://bioinformatics.psb.ugent.be/webtools/plantcare/html>), a database on *cis*-
192 acting regulatory elements in plants.

193

194 **All-against-all comparison of full-length LTR retrotransposon sequences**

195 To detect the sequence similarity within lineages, and also for fine-scale classification of the
196 phylogenetically closely related sequences, all reconstructed *V. macrocarpon* full-length
197 LTR retrotransposons were subjected to an all-against-all comparison together with publicly
198 available reference sequences (Table S1). Representative reference sequences for each LTR
199 retrotransposon lineage from different plant genomes were downloaded from GyDB 2.0
200 [25], (Table S2). The 5' LTRs from these lineages were extracted using Geneious [41]. LTR

201 regions were aligned in Geneious and a distance matrix of pairwise LTR identities was
202 exported. For each lineage, an all-against-all sequence comparison was created using the
203 Python-based software FlexiDot [52]. The FlexiDot parameters were: word size of 26 (-k 26)
204 with 5 mismatches allowed (-S 5). In the dotplot, pairwise sequence identities of the 5'LTRs
205 were printed and shaded.

206

207 **Transcriptome analysis**

208 To assess retrotransposon transcription, publicly available *V. macrocarpon* paired-end
209 mRNA-seq data (Illumina GAIIx; 63.6 million 100 bp reads) from leaves and shoot tips of
210 cranberry cultivar 'Ben Lear' (accession number CNJ99-125-1) of BioProject
211 PRJNA246586 were used [10]. All reads were quality-filtered using the FASTX-toolkit
212 (http://hannonlab.cshl.edu/fastx_toolkit/) within Galaxy [53] with the parameters "quality
213 cut-off" = 10, "percent above cutoff" = 95. Adapter trimming was performed using the
214 BBDuk Geneious plugin [41, 54] with default parameters. About 36.1 million quality-
215 filtered paired-end mRNA reads were used for read mapping against the reconstructed *V.*
216 *macrocarpon* full-length elements of the Ty1-*copia* and the Ty3-*gypsy* superfamily.
217 Transcriptome reads were mapped with the 'Map to reference' tool in Geneious [41] with the
218 parameters: "mapper" = Geneious for RNA-Seq; "sensitivity" = medium-low sensitivity;
219 "span annotated mRNA introns". The graph coverage of the mapped reads was exported and
220 the transcriptome proportion was calculated the counted reads in Excel. A scatter plot
221 visualization was created to compare the genome proportion (derived from the
222 RepeatExplorer analysis) and the transcriptome proportion of the representative full-length
223 Ty1-*copia* and Ty3-*gypsy* LTR retrotransposons [47].

224

225 **Data availability**

226 Full-length retrotransposon sequences and clustered core RT and CHD sequences from this
227 analysis were deposited in the GDV database with the accession number GDV19002
228 (https://www.vaccinium.org/publication_datasets).

229

230 **RESULTS**

231 **Composition of the *V. macrocarpon* retrotransposon landscape**

232 The graph-based clustering with the RepeatExplorer2 pipeline reveals that about 91% of
233 cranberry genome is repetitive (Fig. S1, [35]). According to the automated read cluster
234 classification, only 49% of total repeat reads were assigned to known repeat types. Out of
235 those, LTR retrotransposons constitute the major portion (44%), followed by DNA
236 transposons (11%), non-LTR retrotransposons (6%), rDNAs (3%), satellite DNAs (0.06%)
237 and other repeat fractions (Fig. S1, [35]). The *in silico* analysis shows that full-length LTR
238 retrotransposons constitute a considerable portion of cranberry genome of 4.25 % *Ty1-copia*
239 and 8.41 % *Ty3-gypsy* retrotransposons (Fig. S1; Table S1).

240 For a comprehensive overview of the LTR retrotransposon population, we identified the
241 reverse transcriptase (RT) domains in the *V. macrocarpon* reference genome assembly. We
242 queried the *V. macrocarpon* assembly and identified 181 *Ty1-copia* and 294 *Ty3-gypsy* RT
243 instances. The highest number of RT sequences belonged to the Ale- (n = 110, *Ty1-copia*)
244 and Tat- (n = 184, *Ty3-gypsy*) types (Table S3).

245 Maximum-likelihood analyses allowed us to verify the classification of the reference full-
246 length elements for both superfamilies. For the *Ty1-copia* retrotransposons, we noticed a
247 close relationship between five well-defined lineages: Ikeros, Bianca, Tork, TAR, and

248 Angela. A close relationship of Ivana and SIRE sequences was also observed. Although
249 Alesia only harbors two members, it forms a single, well-separated branch. Finally, the 110
250 Ale RT sequences are polyphyletic and strongly diversified (Fig. 1A). In the superfamily
251 Ty3-*gypsy* (Fig. 2), the non-chromoviruses and chromoviruses are similarly organized. The
252 non-chromovirus retrotransposons of the Tat type were the most abundant and the most
253 diversified. The chromovirus lineage produced four branches representing the clades CRM,
254 Reina, Galadriel, and Tekay (Fig. 2A), with a clade of non-chromovirus Athila as sister.

255 The separation of the chromoviruses into clades could be verified by reconstructing a
256 dendrogram based on their signature chromodomain (CHD) sequences (Fig. S2). Although
257 full-length Reina elements were not found in the RepeatExplorer-based reconstructed
258 retrotransposon sequences (Table S4), the CHD search along the genome assembly of the
259 same *V. macrocarpon* accession, reveals that CHD protein sequences are indeed present
260 from all four chromovirus clades, i.e., CRM, Galadriel, Tekay, and Reina (Fig. S2).

261 Although branch lengths in the dendrogram provide information on the relative divergence
262 between retrotransposon groups, the boxplot visualizations allow a direct comparison of the
263 respective RT diversities (Figs. 1B, 2B). For this, pairwise identities from the RT amino acid
264 multiple sequence alignments were utilized. In Ty1-*copia* lineages, Ale RT sequences have
265 the highest diversity, as they exhibit the widest spread with a maximum of 90 %, a minimum
266 of 28 %, and a median pairwise identity of 59 %. The lowest median identity (57 %) is
267 observed for Ivana, whereas the highest conservation is observed for Bianca with a median
268 pairwise identity of 84 % (Fig. 1B). Generally, Ty1-*copia* RT sequences are more diverse
269 than those from Ty3-*gypsy*. Regarding the median RT identity, the most diverse clade in the
270 Ty3-*gypsy* superfamily is Tat (61%), whereas Reina is the least diversified (73 %). The
271 remaining clades of the Ty3-*gypsy* superfamily (Athila, CRM, Galadriel, and Tekay) have

272 similarly diverse RT sequences, with median RT identities ranging between 65 % and 74 %
273 (Fig. 2B).

274

275 **Structural diversity of Ty1-*copia* and Ty3-*gypsy* LTR retrotransposons**

276 Based on the clustered short reads, we reconstructed and characterized 26 Ty1-*copia* and 28
277 Ty3-*gypsy* full-length retrotransposons, representative for 54 repeat families in the genome
278 of *V. macrocarpon* (Table 1). The majority of these full-length elements (40 out of 54) was
279 represented by a single, circular RepeatExplorer cluster in our analysis (Fig. S3). Reads from
280 the remaining Ty1-*copia* and Ty3-*gypsy* elements were split into smaller clusters. Due to
281 overlapping read pairs, we were able to connect these smaller clusters to superclusters
282 spanning retrotransposons of full length (Table S1). Full-length consensus sequences were
283 reconstructed from eight Ty1-*copia* lineages (Ale, Alesia, Bianca, Ivana/Oryco, TAR,
284 Angela, Tork, and Sireviruses/Maximus, the last being incomplete) and three Ty3-*gypsy*
285 lineages (non-chromoviruses Tat and Athila, and chromoviruses). Typical structural features,
286 including sequence lengths, coding domains and conserved sequence motifs, were analyzed
287 in detail to understand their structural diversity (Figs. 3, 4; Fig. S4; Table 1; Table S5).

288 To illustrate the structure of the individual Ty1-*copia* lineages, we selected eight
289 representative elements (Fig. 3; Fig. S4). Among all the Ty1-*copia* lineages, we
290 reconstructed seven full-length Tork, five Ale, two Alesia, one Bianca, three Ivana/Oryco,
291 and one Sireviruses/Maximus (Table 1; Table S5): The longest consensus element belongs to
292 the latter group and is about 7 kb, it is partially incomplete, containing only the 5' LTR
293 (1545 bp). Other lineages have retrotransposon lengths ranging between 4.8-6.8 kb with
294 LTR lengths between 127-994 bp (Table 1; Table S5). While the highest genome proportion
295 among the reconstructed Ty1-*copia* elements belongs to the elements SCL22_Ale (0.57 %),

296 the lowest genome proportion is observed for SCL344_Alesia (0.01 %) (Table S5). Although
297 all full-length *in silico* Ty1-*copia* elements contain uninterrupted *gag* and *pol* protein
298 domains, some characteristic retrotransposon features were absent. For example, while most
299 of the *in silico* elements that represent Ale, Angela and Tork retrotransposons harbored the
300 PBS, PPT, and both 5' and 3' LTRs, one of these structural motifs was missing in the
301 representative *in silico* sequences of Bianca, Ivana/Oryco, TAR, and Sireviruses/Maximus
302 (Table 1; Table S5). The most common PBS-type detected in Ty1-*copia* lineages
303 corresponds to the Met-type. For Ale and TAR retrotransposons, however, a PBS of the
304 Asn/Met-type (Ale) or Glu/Ser/Met/Tyr/Trp-type (TAR) was more common (Table 1; Table
305 S5).

306 Six reference elements were selected to represent the structural variability of each Ty3-*gypsy*
307 lineage (Fig. 4). The most diversified and abundant Ty3-*gypsy* elements in the *V.*
308 *macrocarpon* genome belong to the Tat group (Table 1). A total of 13 full-length elements
309 have been characterized, with genome proportions ranging between 0.12 and 1.34 % (Table
310 1). In contrast, the other two groups only contribute 0.01-1.92 % (Athila) and 0.02-0.23 %
311 (chromoviruses) to the genome (Table 1; Table S5). Although *gag* and *pol* genes were
312 detected for all full-length elements, structurally complete elements with all motifs typical
313 for LTR retrotransposon (PBS, PPT and both 3' and 5' LTR regions) were mostly present in
314 the Galadriel clades (chromoviruses). Most elements of the other groups (chromovirus
315 Tekay, CRM and non-chromoviruses Tat and Athila) lack at least one of these structural
316 motifs. Moreover, some Ty3-*gypsy* elements have additional protein domains as compared to
317 Ty1-*copia* elements. For instance, the *pol* genes of the non-chromovirus Tat clade have dual
318 ribonuclease H domains (gRNH and aRNH), whereas the chromoviruses harbor an
319 additional chromodomain (CHD) region (Fig. 4; Table S5). A canonical CHD was only
320 found for the elements of the Tekay and the Galadriel clades, but surprisingly not within the

321 elements of the chromovirus CRM clade (Fig. 4; Fig. S4; Table S5). SCL16_Ogre is the
322 longest Ty3-*gypsy* Tat element, with a total length of 19 kb and 1251-1326 bp long LTRs. In
323 contrast, SCL80_CRM is the shortest element (5.2 kb) with LTRs ranging between 361-556
324 bp (Table S5). A Met-type PBS was only found in the CRM and Galadriel chromoviruses,
325 whereas other Ty3-*gypsy* elements were more variable regarding their PBS (Table 1).

326

327 **Cranberry retrotransposons are often associated with short tandem repeats**

328 Although tandem repeats, i. e. satellite or ribosomal DNAs, represent major repetitive
329 genome fractions that are distinct from the more dispersed transposable elements, tandemly
330 repeated motifs can also be associated with retrotransposons. In *V. macrocarpon*, we
331 detected tandem repeats in the reconstructed full-length Ty1-*copia* and Ty3-*gypsy*
332 retrotransposon sequences (Table 2; Table S6). These tandem repeats were either localized
333 within the 3' or the 5' LTRs, within the untranslated regions (UTRs) adjacent to either LTR,
334 or in some cases within the UTRs adjacent to the *gag* gene. Tandem repeats found in
335 different retrotransposons vary in sequence composition, period size, and number of
336 repetitions (Table 2; Table S6). For the Ty1-*copia* superfamily, we identified tandem repeats
337 in all lineages, embedded in either the LTRs or UTRs, with consensus period lengths ranging
338 from 2 to 49 bp with about 2 to 8 repetitions (Table 2; Table S6).

339 In Ty3-*gypsy* LTR retrotransposons, tandem repeats were detected both in non-
340 chromoviruses and chromoviruses, yet to a different extent. Non-chromovirus Tat and
341 chromoviruses harbored tandem repeats both within the LTRs and UTRs, whereas non-
342 chromovirus Athila had tandem repeats only embedded within the UTR adjacent to the *gag*
343 gene (Table 2). Among all Ty3-*gypsy* LTR retrotransposons, Tat (including Ogre) harbor
344 tandem repeats with the largest monomer size (up to 93 bp) and the highest number of

345 repetitions (up to 15 times) (Table 2; Table S6).

346

347 ***Cis*-regulatory elements are often associated with LTR sequences**

348 As LTRs frequently harbor *cis*-regulatory elements, our targeted search revealed that almost
349 all reconstructed retrotransposons include typical promoter motifs, including TATA and
350 CAAT boxes within their LTRs (Fig. S5A-B; Table 2; Table S7). However, as the LTR
351 composition varies, number and position of these promoter sequences are variable. For
352 instance, the LTRs of the Ty1-*copia* lineages Ale, Angela, TAR, Tork, Sireviruses/Maximus,
353 and all Ty3-*gypsy* lineages contained multiple promoter motifs in different positions.
354 However, the Bianca LTR only contains a CAAT box, lacking a TATA box. For Ivana, two
355 out of three *in silico* elements (SCL42_Ivana and SCL310_Ivana) contain both promoter
356 motifs in their LTRs (Fig. S5A-B; Table S7).

357 In addition to the promoter motifs, we detected other *cis*-regulatory elements for different
358 LTR retrotransposons (Table S7). These regions are mostly related to light, hormone, and
359 cell cycle responsiveness, as well as abiotic stress tolerance (Fig. S5A-B; Table S7). Among
360 the eight representative sequences of the Ty1-*copia* lineages, most regulatory motif types
361 were predicted for TAR (19), Tork (19), and Ale (10) (Table 2), whereas among the Ty3-
362 *gypsy* elements, the Tat retrotransposons (including Ogre) harbor most variability with 38
363 distinct motifs.

364

365 **Sequence homogeneity between closely related elements**

366 In an attempt to fine-scale the classification of the full-length retrotransposon elements of *V.*
367 *macrocarpon*, closely related elements were subjected to an all-against-all sequence

368 comparison with publicly available references to reveal sequence similarities among them
369 (Table S2).

370 Overall, the *V. macrocarpon* full-length elements have higher similarity to other elements
371 within their protein coding regions, contrasting the variability within their UTRs and LTRs
372 (Fig. S6A-E). However, the levels of similarity for coding and non-coding regions vary
373 depending on the lineage (Fig. S4). For instance, the RT similarities within the *Ty1-copia*
374 and within the *Ty3-gypsy* groups ranged between 57-74 % and 55-88 %, respectively. In
375 contrast, LTR similarities ranged between 17-47 % for *Ty1-copia* groups and between 11-
376 55 % for *Ty3-gypsy* groups, respectively (Fig. 6A-E; Fig. S4; Table 1). Ale and Tat
377 sequences are the most diversified, both within their RTs and their LTRs (Table 1).

378

379 **Transcriptome analysis reveals transcriptional activity of LTR retrotransposons**

380 We investigated whether the identified full-length LTR retrotransposons are transcribed into
381 RNA by analyzing the reference transcriptome database of *V. macrocarpon* from leaves and
382 shoot tips (Fig. 5; Fig. S7; Table S8). We have mapped a total of 36.1 million transcript
383 sequences against the 54 reconstructed *Ty3-gypsy* and *Ty1-copia* sequences of *V.*
384 *macrocarpon*. Subsequently, we compared the genome proportion derived from the
385 RepeatExplorer analysis against the transcriptome proportion for each individual element
386 (Fig. 5).

387 All the identified full-length elements are transcribed to a different extent, thereby
388 contributing a total of 0.1 % by *Ty1-copia* retrotransposons and 0.03 % by *Ty3-gypsy*
389 retrotransposons to the used *V. macrocarpon* transcriptome reads (Fig. 5; Fig. S7; Table S8).

390 Among all LTR retrotransposon lineages, Ale has the highest proportion of transcript
391 sequences (Fig. 5; Table S8). In fact, the two reference elements SCL5_Ale and SCL22_Ale

392 have both the highest genomic proportions of all 26 Ty1-*copia* representatives and the
393 highest percentage of mapped RNA reads of all 54 reference LTR retrotransposons in *V.*
394 *macrocarpon* (Fig. 5; Table S8). Nevertheless, the results also show that the percentage of
395 mapped transcripts is dependent on the individual full-length elements and cannot be
396 generalized across and within lineages. Nevertheless, the number of mapped transcript
397 sequences is not proportional to their genomic abundance: most lineages, including Angela,
398 Bianca, Ivana, TAR, Tork, and SIRE, have less representative transcript sequences as
399 compared to their proportional genomic abundance (Fig. 5; Table S8).

400 Despite their higher genomic abundance, the transcription level of Ty3-*gypsy*
401 retrotransposons is much lower (0.037%) than the transcription of Ty1-*copia*
402 retrotransposons (0.105%): SCL15_Ogre, SCL27_Ogre, SCL28_Tekay, and SCL31_CRM
403 generate most transcripts with similar counts (0.003-0.004 % of the total transcript
404 sequences; Fig. 5; Table S8). Noteworthy, SC31_CRM has only a low genomic abundance
405 (0.22%), but the highest number of transcripts (0.004%) as compared to other Ty3-*gypsy*
406 elements.

407 We have also investigated the distribution of mapped transcript reads along the most
408 transcribed full-length Ty1-*copia* and Ty3-*gypsy* elements to gain insight into the coverage
409 of the different structural features. In Fig. S7, the most transcribed elements from each group
410 and the distribution of mapped RNA reads are shown. The distribution of mapped transcript
411 sequences along the different structural components of the LTR retrotransposons shows
412 either a homogenous profile along the full-length element or an enrichment in the LTRs and
413 UTRs. SCL5_Ale and SCL49_Tork have homogeneously mapped transcript sequences along
414 their *gag-pol* genes and UTRs (Fig. S7). In contrast, SCL42_Ivana, SCL40_SIRE,
415 SCL49_Ogre, and SCL31_CRM have the highest number of mapped transcript reads in the

416 two LTRs. Moreover, SCL79_Bianca and SCL6_Athila have most transcripts mapped to the
417 internal UTR located between the *gag* and *pol* genes as well as to the 3' LTR (Fig. S7).

418

419 **DISCUSSION**

420 Although a draft genome assembly for *V. macrocarpon* is already available [10], we are only
421 starting to decipher its repetitive fraction. As we showed in our previous work that focused
422 on satellite DNAs in *Vaccinium* genomes [36], repetitive DNAs are still the last genomic
423 portions to be assembled completely. Nevertheless, as knowledge of repetitive DNAs is
424 needed to understand *Vaccinium* genomes in their entirety, we extend our previous work to
425 characterize the contribution of long terminal repeat (LTR) retrotransposons to the cranberry
426 genome. As we base our analysis on a random sampling of Illumina reads, we avoid biases
427 introduced by a potentially erroneous genome assembly [39]. Nevertheless, as this approach
428 may neglect less abundant and more diverged repeats [55], all quantifications are likely
429 underestimations.

430 We showed that nearly a quarter of the *V. macrocarpon* genome is made up of highly
431 abundant, repetitive TE families, with over 91% of genome considered as repetitive [35].
432 This nearly doubles the estimate retrieved for the cranberry genome assembly (40%) [10,
433 11], and likely is a result of a systematic exclusion of repetitive regions from the genome
434 assembly. Within the repetitive genome fraction, LTR retrotransposons are the most
435 abundant TEs (44 % of the repetitive fraction), followed by DNA transposons (11 %), and
436 non-LTR retrotransposons (6 %). The observed high abundance of LTR retrotransposons is
437 typical for plants and compares well to other TEs in other species such as flax [56], apple
438 [57] maize [58], and oats [59].

439

440 **Cranberry LTR retrotransposons are diverse in sequence and structure**

441 Overall, in the *V. macrocarpon* genome, Ty3-gypsy elements are on average longer and
442 contribute a higher genome proportion than Ty1-copia elements. This was also reported in
443 other plants, such as *Avena sativa* [59], *Triticum* sp. [60], *Festuca pratensis* [61], and *Malus*
444 *domestica* [57]. Nevertheless, opposite proportions are also possible, as in *Musa acuminata*
445 and *Pyrus bretschneideri*, in which Ty1-copia retrotransposons outweighed Ty3-gypsy
446 elements [62]. Variation in sequence lengths between these two groups is thought to reflect
447 length differences in the non-coding rather than in the coding element regions, whereas the
448 genome proportion depends on their copy number rather than on the mere retrotransposon
449 length [63-65]. However, elements from Errantivirus/Athila and Ogre/Tat lineages from
450 Ty3-gypsy were found to have both the highest genome proportions and the largest sizes in
451 the *V. macrocarpon* genome. Besides the generally large lengths of Ogre/Tat
452 retrotransposons in *V. macrocarpon*, the high number of RT sequences detected in this study
453 also indicates a high copy number – both contributing to their high genomic abundance.
454 Similarly, for the Ale lineage, we found the most RT instances and a higher genome fraction
455 as compared to the other Ty1-copia lineages.

456 Characterization of *V. macrocarpon*'s RT sequences revealed that Tat and Ale are the most
457 diversified groups among the Ty3-gypsy and Ty1-copia retrotransposons, respectively. By
458 studying the pairwise RT identities, we observed that individual lineages diversify
459 differently. This is in accordance with our previous report on the related diploid species
460 *Vaccinium corymbosum* [35]. Species- and genus-specific lineage diversification and
461 heterogeneity is reported by many authors for several species. For example, within Ty1-
462 copia, Sirevirus/Maximus retrotransposons were highly diversified in sugarcane [64], in the
463 Fabaceae [65], and in maize [66], whereas Ale retrotransposons were found in highest diversity

464 in *Chenopodium quinoa* [67] and *Populus trichocarpa* [68]. Within Ty3-gypsy, the Tat
465 lineage was found to proliferate in sugarcane [64] and in the Fabaeae [65], whereas the Athila
466 clade in *Populus trichocarpa* [68], and the Tekay chromoviruses was for example more
467 abundant in *Beta vulgaris* [69]. The overall abundance and diversification of the various
468 LTR retrotransposon types also reflects the evolutionary relationships of the host species and
469 may be similar within a group of closely related species [70, 71]. The environmental
470 adaptation of a particular species may also be affected by the insertion and molecular
471 diversification of these retrotransposons [16, 72].

472 Structurally, regarding the encoded protein domains, the reverse transcriptase (RT) is the
473 most ancient domain and key to the dispersion of retrotransposons [17, 73]. The RT is
474 present in all autonomous LTR retrotransposons, whereas the modular acquisition of other
475 polyprotein domains is assumed to result from independent lineage-specific evolutionary
476 histories [17, 74]. Therefore, structural differences among the lineages are considerable. This
477 is particularly evident in the Tat and chromovirus lineages that harbor additional protein-
478 coding domains. All members of the Tat lineage in *V. macrocarpon* harbor two ribonuclease
479 H domains. Ustyantsev *et al.* (2015) [75] speculated that this duality of ribonucleases H
480 might benefit the successful strand transfer, a property common to retroviruses that suffer
481 natural selection pressure during proliferation.

482 Chromoviral integrases are typically marked by a chromodomain at their C-terminus,
483 presumably with the ability to bind specific histone variants [69 76-79]. Therefore,
484 chromoviruses may directly impact chromatin structure and function [79]. For the
485 chromovirus clade Galadriel and Tekay, a chromodomain was detected both in the *V.*
486 *macrocarpon* genome assembly and in the representative full-length elements from the read
487 cluster analysis. In contrast, for the CRM clade, chromodomain regions were extracted from

488 the genome assembly only, although both genome assembly and read cluster analysis were
489 based on the same *V. macrocarpon* accession. Compared to the CHDs of the Tekay,
490 Galadriel, and Reina clades, their identification within centromeric retrotransposons is more
491 difficult, as they typically do not show structural similarities to type I and II cellular
492 chromodomains and are highly variable even between families within a species [69, 78].
493 Nevertheless, this domain, also referred to as the CR motif or CHDCR, correlates with
494 various sequence and structural features [28]. For example, the RTs of CRM elements tend
495 to form a separate branch in the analysis of sequence relationship. Another feature is the
496 position of the CHDCR, which is defined by its extension into the 3' LTR. Although some of
497 the reference elements of the CRM clade have incomplete LTR sequences, an extended *gag-*
498 *pol* ORF into the 3' LTR is recognizable. CR chromodomains of CRM elements were found
499 to be mostly associated with the centromeric heterochromatin. These domains are likely
500 responsible for the diversification, integrity, and stability in the functional centromeric
501 heterochromatin through an RNA-mediated mechanism in different beet and grass species
502 [69, 79-81].

503 The PBS is an important structural feature for LTR retrotransposon transcription located
504 downstream of the 5' LTR, with PBS type and length characteristic for a given plant LTR
505 retrotransposon group [28]. The PBS is present in most of the *in silico* retrotransposons
506 constructed for *V. macrocarpon*. Nevertheless, some consensus lack a designated PBS,
507 likely either an assembly or identification error, or the nature of the PBS itself in the studied
508 species. In chromoviruses as well as in Ty1-*copia* retrotransposons, a PBS corresponding to
509 the initiator methionine tRNA (Met-tRNA) was most common, similar to other plant
510 genomes [28, 82]. In contrast, the PBS regions in other Ty3-*gypsy* retrotransposons were
511 more diverse. The specific PBS site is not only significant for the transcription of LTR
512 retrotransposons, but also for their post-transcriptional regulation through tRNA-derived

513 small RNAs [82, 83]. Therefore, the acquisition of distinct PBS sequences in different
514 retrotransposon lineages may point to individual evolutionary benefits.

515 The LTRs of the identified full-length Ty1-*copia* and Ty3-*gypsy* retrotransposons were
516 enriched for different *cis*-regulatory elements, including putative promoters and sequences
517 related to hormone responsiveness, cell development, and stress response. This indicates the
518 relevance of these motifs for retrotransposon transcription and activity. In addition, the *V.*
519 *macrocarpon* full-length elements harbored tandem repeats not only in their LTR region but
520 also in their UTRs. Yet, none were found within coding sequences. The presence of *cis*-
521 regulatory regions and tandem repeats in the LTR is common for LTR retrotransposons and
522 can have significant impacts on their reverse transcription and proliferation [25, 55].
523 Moreover, *cis*-regulatory sequences could act as enhancers or repressors and thereby affect
524 the expression of downstream genes as part of regulatory networks and pathways [55, 84].
525 Although still a matter of investigation, the origin of tandem repeats from different parts of
526 the retrotransposons (i.e. LTRs or UTRs) with specific functions are reported for many
527 organisms. An example are the centromeric tandem repeats that may have originated from
528 different retrotransposons in maize [85], wheat [86], and potato [87], suggesting their
529 importance for the formation of functional centromeres in these species. Moreover, the
530 existence of short tandem repeats varying in length and sequence in UTR sequences near the
531 PPT was also reported for many other plant genomes, such as legumes [63] and *Silene*
532 *latifolia* [88]. Although the specific function of these short repeats is still debated, it is
533 hypothesized that they could serve as a junction, connecting the *gag-pol* gene to the 3' end.
534 Thus, in evolutionary time-spans, they may serve as a seed for the generation of longer
535 satellite DNA arrays, with the potential to acquire a structural function for the organization
536 of chromosomes [88].

537 We found that the overall similarity among the full-length elements within each lineage
538 varied. For example, Tat retrotransposons are diversified and have little overall sequence
539 similarity within the UTR regions as was also reported for other organisms [63, 65]. The
540 highly diverse UTRs and LTRs on clade level (see Fig. S6) may indicate their ancient origin
541 and diversification in the *Vaccinium* ancestors [28]. LTRs and UTRs appear to evolve faster
542 than coding regions due to different mechanisms, like unequal recombination or illegitimate
543 recombination [63, 89-92].

544

545 **Cranberry LTR retrotransposons of the Ale-type are strongly transcribed in *V.***
546 ***macrocarpon* leaves and shoot tips**

547 As expected, all *V. macrocarpon* LTR retrotransposons investigated showed at least a basal
548 level of transcription, together making up about 0.13 % of the total transcriptome. According
549 to [10], transcription of TEs in *V. macrocarpon* is quite low (about 4.3%) and in contrast to
550 the highly transcribed protein-coding genes (about 83.59%). This is in line with our study, in
551 which only the most abundant LTR retrotransposons were queried. A low transcription of
552 retrotransposons compared to conventional protein-coding genes is a common finding in
553 most plant species, including grasses [93], flax [56], and sugarcane [64]. The low
554 transcription of full-length LTR retrotransposons indicates out that these sequences are
555 rather strictly regulated in cranberry [10].

556 Transcriptional activity of full-length LTR retrotransposons is not correlated with their
557 genome proportion in *V. macrocarpon*. Higher levels of expression than those found here for
558 less abundant elements (such as Ale elements) have been reported in sugarcane [64], maize
559 [58], poplar [94], *Arabidopsis arenosa* [95], and sunflower [96]. Overall, Ty1-*copia*
560 sequences appeared to be more transcribed than Ty3-*gypsy* elements in *V. macrocarpon*

561 genome being the highest transcription profiles (~0.08%) those of the Ale full-length
562 elements. In contrast, Athila elements have a higher genomic percentage (2.44%), but have
563 fewer transcript reads assigned (0.004 %), even not covering the full reference element
564 length. In fact, it is considered that Ale and Sireviruses are some of the most
565 transcriptionally active Ty1-*copia* lineages in plants [58, 64, 66, 95-97]. In cranberry, for Ale
566 retrotransposons, this also seems to hold true, whereas for Sireviruses, fewer genomic and
567 transcriptomic proportions were detected. Nevertheless, both Ale and Sireviruses are
568 preferentially accumulated in the euchromatic region of the genome in different species [58,
569 95, 96] and hence could affect gene regulatory pathways.

570 In contrast, retrotransposon heterochromatization and fragmentation may decrease element
571 transcription. We know from other plant genomes that some lineages are more prone to
572 truncation and heterochromatic burial, e.g. Athila/Errantiviruses or even the related
573 pararetroviruses in sugar beet [98, 99]. Similar tendencies may impact the TE landscape in
574 *V. macrocarpon*, yielding lineages that are less transcribed than others, e.g.
575 Athila/Errantiviruses and Sireviruses. Transcriptional activity of TEs is generally dependent
576 on several factors and may be correlated with development [100], (epi)genetic regulation
577 [15, 16, 64, 101], and environmental adaptation [15, 95, 100]. Therefore, TE transcription
578 can depend on the genomic neighborhood, the presence of *cis*-regulatory motifs, tissue types,
579 developmental states, and environmental effects [100, 101]. Here we observed considerable
580 differences regarding the genome and transcriptome abundances of repetitive sequences in *V.*
581 *macrocarpon*, which imply that the mechanisms of transcriptional regulation vary depending
582 on the LTR retrotransposon.

583 Regarding their applicability, the retrotransposons identified here may provide suitable
584 targets for the development of molecular markers in assisted breeding programs [19-23]. We

585 suggest targeting retrotransposons that are most abundant and for which high transcription
586 implies closeness to genes. Hence, we would suggest Ale retrotransposons within the Ty1-
587 *copia* superfamily, as well as Athila retrotransposons within the Ty3-*gypsy*. Both are
588 abundant in genomic and transcriptomic databases, and thus may provide many primer
589 binding sites and likely are situated in the more openly packed euchromatin.

590

591 **Conclusion**

592 Using short reads and the sequence assembly, we have provided an exhaustive overview of
593 the LTR retrotransposon landscape in the cranberry genome. We have detected all major
594 LTR retrotransposon lineages, with some showing association to short tandemly repeated
595 motifs. Considering their high genomic abundance and transcriptional activity, we suggest
596 that Ale and Athila TEs likely represent the most useful targets to survey TE-derived
597 polymorphisms across genotypes.

598

599 **Conflict of interest**

600 The authors have no conflict of interest to report.

601

602 **Acknowledgments**

603 This work was supported by the Scientific and Technological Research Council of Turkey
604 (TUBITAK) TUBITAK-2215 PhD scholarship and the Scientific Research Projects Unit
605 (BAP) of Niğde Ömer Halisdemir University (FEB 2017/18 DOKTEP). This work also
606 acknowledges the financial support by the Jagannath University, Bangladesh, in the form of
607 study grants as well as the Georg Forster fellowship of the Alexander von Humboldt
608 Foundation awarded to NS.

609

610 **Abbreviations**

611 LTR - long terminal repeat; TEs - transposable elements; PR - protease; gRNH - RNase H;
612 RT- reverse transcriptase; INT - integrase, PBS - primer binding site; PPT - polypurine tract;
613 CHD – chromodomain; UTR – untranslated region; DNA - deoxyribonucleic acid; RNA -
614 ribonucleic acid.

615 **References**

- 616 [1] Trehane J. Blueberries, cranberries, and other vacciniiums. Timber Press; 2004.
- 617 [2] Kron KA, Powell EA, Luteyn JL. Phylogenetic relationships within the blueberry tribe
618 (Vaccinieae, Ericaceae) based on sequence data from matK and nuclear ribosomal ITS
619 regions, with comments on the placement of Satyria. *Am J Bot.* 2002;89(2):327-36.
- 620 [3] Powell EA, Kron KA. Hawaiian blueberries and their relatives—a phylogenetic analysis
621 of *Vaccinium* sections *Macropelma*, *Myrtillus*, and *Hemimyrtillus* (Ericaceae). *Syst*
622 *Bot.* 2002;27(4):768-79.
- 623 [4] Vander Kloet SP, Dickinson TA. A subgeneric classification of the genus *Vaccinium*
624 and the metamorphosis of *V.* section *Bracteata* Nakai: more terrestrial and less
625 epiphytic in habit, more continental and less insular in distribution. *J Plant Res.*
626 200;122(3):253-68.
- 627 [5] Hancock JF, Lyrene P, Finn CE, Vorsa N, Lobos GA. Blueberries and cranberries. In:
628 Hancock JF, editor. *Temperate fruit crop breeding*. Springer; 2008, p. 115-50.
- 629 [6] Česonienė L, Daubaras R, Paulauskas A, Zukauskienė J, Zych M. Morphological and
630 genetic diversity of European cranberry (*Vaccinium oxycoccos* L., Ericaceae) clones in
631 Lithuanian reserves. *Acta Soc Bot Pol.* 2013;82(3):211-217.
- 632 [7] Moyer RA, Hummer KE, Finn CE, Frei B, Wrolstad RE. Anthocyanins, phenolics, and
633 antioxidant capacity in diverse small fruits: *Vaccinium*, *Rubus*, and *Ribes*. *J Agric*
634 *Food Chem.* 2002;50(3):519-25.
- 635 [8] Nickavar B, Amin G. Anthocyanins from *Vaccinium arctostaphylos* berries. *Pharm*
636 *Biol.* 2004;42(4-5):289-91.
- 637 [9] Zdepski, A., Debnath, S.C., Howell, A., Polashock, J., Oudemans, P., Vorsa, N. and

- 638 Michael, T.P. Cranberry. In Genetics, genomics and breeding of berries; Folta, K.,
639 Kole, C., Ed.; Boca Raton, F.L., p. 200 ISBN 9781578087075 - CAT# N10335, 2011.
- 640 [10] Polashock J, Zelzion E, Fajardo D, Zalapa J, Georgi L, Bhattacharya D, Vorsa N. The
641 American cranberry: first insights into the whole genome of a species adapted to bog
642 habitat. *BMC Plant Biol.* 2014;14(1):1-8.
- 643 [11] Diaz-Garcia L, Garcia-Ortega LF, González-Rodríguez M, Delaye L, Iorizzo M,
644 Zalapa J. Chromosome-level genome assembly of the American cranberry (*Vaccinium*
645 *macrocarpon* Ait.) and its wild relative *Vaccinium microcarpum*. *Front Plant Sci.*
646 2021;12:137.
- 647 [12] Gupta V, Estrada AD, Blakley I, Reid R, Patel K, Meyer MD, Andersen SU, Brown
648 AF, Lila MA, Loraine AE. RNA-Seq analysis and annotation of a draft blueberry
649 genome assembly identifies candidate genes involved in fruit ripening, biosynthesis
650 of bioactive compounds, and stage-specific alternative splicing. *Gigascience.*
651 2015;4(1):s13742-015.
- 652 [13] Colle M, Leisner CP, Wai CM, Ou S, Bird KA, Wang J, Wisecaver JH, Yocca AE,
653 Alger EI, Tang H, Xiong Z. Haplotype-phased genome and evolution of phytonutrient
654 pathways of tetraploid blueberry. *GigaScience.* 2019;8(3):giz012.
- 655 [14] Biémont C, Vieira C. Junk DNA as an evolutionary force. *Nature.*
656 2006;443(7111):521-4.
- 657 [15] Oliver KR, McComb JA, Greene WK. Transposable elements: powerful contributors
658 to angiosperm evolution and diversity. *Genome Biol Evol.* 2013;5(10):1886-901.
- 659 [16] Wendel JF, Jackson SA, Meyers BC, Wing RA. Evolution of plant genome
660 architecture. *Genome Biol.* 2016;17(1):1-4.

- 661 [17] Ustyantsev K, Blinov A, Smyshlyaev G. Convergence of retrotransposons in
662 oomycetes and plants. *Mob DNA*. 2017;8(1):1-1.
- 663 [18] Niu XM, Xu YC, Li ZW, Bian YT, Hou XH, Chen JF, Zou YP, Jiang J, Wu Q, Ge S,
664 Balasubramanian S. Transposable elements drive rapid phenotypic variation in
665 *Capsella rubella*. *Proc Natl Acad Sci*. 2019;116(14):6908-13.
- 666 [19] Kalendar R, Schulman AH. Transposon-based tagging: IRAP, REMAP, and iPBS. In
667 *Molecular Plant Taxonomy 2014* (pp. 233-255). Humana Press, Totowa, NJ.
- 668 [20] Jiang S, Zong Y, Yue X, Postman J, Teng Y, Cai D. Prediction of retrotransposons and
669 assessment of genetic variability based on developed retrotransposon-based insertion
670 polymorphism (RBIP) markers in *Pyrus* L. *Mol Genet Genom*. 2015;290(1):225-37.
- 671 [21] Zong Y, Kang H, Fang Q, Chen X, Zhou M, Ni J, Zhang Y, Wang L, Zhu Y, Guo W.
672 Phylogenetic relationship and genetic background of blueberry (*Vaccinium* spp.) based
673 on retrotransposon-based SSAP molecular markers. *Scientia Horticulturae*.
674 2019;247:116-22.
- 675 [22] Seibt KM, Wenke T, Wollrab C, Junghans H, Muders K, Dehmer KJ, Diekmann K,
676 Schmidt T. Development and application of SINE-based markers for genotyping of
677 potato varieties. *Theor Appl Genet*. 2012;125(1):185-96.
- 678 [23] Reiche B, Koegler A, Morgenstern K, Brueckner M, Weber B, Heitkam T, Seibt KM,
679 Troeber U, Meyer M, Wolf H, Schmidt T. Application of retrotransposon-based Inter-
680 SINE Amplified Polymorphism (ISAP) markers for the differentiation of common
681 poplar genotypes. *Can J For Res*. 2021 (ja).
- 682 [24] Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, Flavell A, Leroy
683 P, Morgante M, Panaud O, Paux E. A unified classification system for eukaryotic

- 684 transposable elements. *Nat Rev Genet.* 2007;8(12):973-82.
- 685 [25] Llorens C, Futami R, Covelli L, Domínguez-Escribá L, Viu JM, Tamarit D, Aguilar-
686 Rodríguez J, Vicente-Ripolles M, Fuster G, Bernet GP, Maumus F. The Gypsy
687 Database (GyDB) of mobile genetic elements: release 2.0. *Nucleic Acids Res.*
688 2010;39(suppl_1):D70-4.
- 689 [26] Piégu B, Bire S, Arensburger P, Bigot Y. A survey of transposable element
690 classification systems—a call for a fundamental update to meet the challenge of their
691 diversity and complexity. *Mol Phylogenet Evol.* 2015;86:90-109.
- 692 [27] Bao W, Kojima KK, Kohany O. Repbase Update, a database of repetitive elements in
693 eukaryotic genomes. *Mob DNA.* 2015;6(1):1-6.
- 694 [28] Neumann P, Novák P, Hošťáková N, Macas J. Systematic survey of plant LTR-
695 retrotransposons elucidates phylogenetic relationships of their polyprotein domains
696 and provides a reference for element classification. *Mob DNA.* 2019;10(1):1-7.
- 697 [29] Storer J, Hubley R, Rosen J, Wheeler TJ, Smit AF. The Dfam community resource of
698 transposable element families, sequence models, and genome annotations. *Mob DNA.*
699 2021;12(1):1-4.
- 700 [30] Elliott TA, Heitkam T, Hubley R, Quesneville H, Suh A, Wheeler TJ. TE Hub: A
701 community-oriented space for sharing and connecting tools, data, resources, and
702 methods for transposable element annotation. *Mob DNA.* 2021;12(1):1-5.
- 703 [31] Wicker T, Keller B. Genome-wide comparative analysis of copia retrotransposons in
704 Triticeae, rice, and *Arabidopsis* reveals conserved ancient evolutionary lineages and
705 distinct dynamics of individual copia families. *Genome Res.* 2007;17(7):1072-81.
- 706 [32] Schnable PS, Ware D, Fulton RS, Stein JC, Wei F, Pasternak S, Liang C, Zhang J,

- 707 Fulton L, Graves TA, Minx P. The B73 maize genome: complexity, diversity, and
708 dynamics. *Science*. 2009;326(5956):1112-5.
- 709 [33] Galindo-González L, Mhiri C, Deyholos MK, Grandbastien MA. LTR-
710 retrotransposons in plants: Engines of evolution. *Gene*. 2017;626:14-25.
- 711 [34] Marín I, Lloréns C. Ty3/Gypsy retrotransposons: description of new *Arabidopsis*
712 *thaliana* elements and evolutionary perspectives derived from comparative genomic
713 data. *Mol Biol Evol*. 2000;17(7):1040-9.
- 714 [35] Sultana N, Serçe S, Menzel G, Heitkam T, Schmidt T. Comparative analysis of
715 repetitive sequences reveals genome differences between two common cultivated
716 *Vaccinium* Species (*V. corymbosum* and *V. macrocarpon*). *J Mol Biol Biotech*.
717 2017;2:24.
- 718 [36] Sultana N, Menzel G, Heitkam T, Kojima KK, Bao W, Serçe S. Bioinformatic and
719 Molecular Analysis of Satellite Repeat Diversity in *Vaccinium* Genomes. *Genes*.
720 2020;11(5):527.
- 721 [37] Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina
722 sequence data. *Bioinformatics*. 2014;30(15):2114-20.
- 723 [38] Novak P, Neumann P, Pech J, Steinhaisl J, Macas J. RepeatExplorer: a Galaxy-based
724 web server for genome-wide characterization of eukaryotic repetitive elements from
725 next-generation sequence reads. *Bioinformatics*. 2013;29(6):792-3.
- 726 [39] Novák P, Neumann P, Macas J. Global analysis of repetitive DNA from unassembled
727 sequence reads using RepeatExplorer2. *Nat Protoc*. 2020;15(11):3745-76.
- 728 [40] Smit AF, Hubley R, Green P. RepeatMasker Open-4.0. 2015.
- 729 [41] Kearsse M, Moir R, Wilson A, Stones-Havas S, Cheung M, Sturrock S, Buxton S,

- 730 Cooper A, Markowitz S, Duran C, Thierer T. Geneious Basic: an integrated and
731 extendable desktop software platform for the organization and analysis of sequence
732 data. *Bioinformatics*. 2012;28(12):1647-9.
- 733 [42] Xu Z, Wang H. LTR_FINDER: an efficient tool for the prediction of full-length LTR
734 retrotransposons. *Nucleic Acids Res*. 2007;35(suppl_2):W265-8.
- 735 [43] Rice P, Longden I, Bleasby A. EMBOSS: the European molecular biology open
736 software suite. *Trends Genet*. 2000;16(6):276-7.
- 737 [44] Benson G. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic*
738 *Acids Res*.1999;27(2):573-80.
- 739 [45] Li W, Godzik A. Cd-hit: a fast program for clustering and comparing large sets of
740 protein or nucleotide sequences. *Bioinformatics*. 2006;22(13):1658-9.
- 741 [46] Katoh K, Standley DM. MAFFT multiple sequence alignment software version 7:
742 improvements in performance and usability. *Mol Biol Evol*. 2013;30(4):772-80.
- 743 [47] Team RC. R: A language and environment for statistical computing.
- 744 [48] Price MN, Dehal PS, Arkin AP. FastTree 2—approximately maximum-likelihood trees
745 for large alignments. *PloS one*. 2010;5(3):e9490.
- 746 [49] Stamatakis A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of
747 large phylogenies. *Bioinformatics*. 2014;30(9):1312-3.
- 748 [50] Stothard P. The sequence manipulation suite: JavaScript programs for analyzing and
749 formatting protein and DNA sequences. *Biotechniques*. 2000;28(6):1102-4.
- 750 [51] Lescot M, Déhais P, Thijs G, Marchal K, Moreau Y, Van de Peer Y, Rouzé P,
751 Rombauts S. PlantCARE, a database of plant cis-acting regulatory elements and a
752 portal to tools for in silico analysis of promoter sequences. *Nucleic Acids Res*.

- 753 2002;30(1):325-7.
- 754 [52] Seibt KM, Schmidt T, Heitkam T. FlexiDot: highly customizable, ambiguity-aware
755 dotplots for visual sequence analyses. *Bioinformatics*. 2018;34(20):3575-7.
- 756 [53] Gordon A, Hannon GJ. Fastx-toolkit. FASTQ/A short-reads preprocessing tools
757 (unpublished) http://hannonlab.cshl.edu/fastx_toolkit. 2010;20;5.
- 758 [54] Bushnell B. BMap: a fast, accurate, splice-aware aligner. Lawrence Berkeley
759 National Lab.(LBNL), Berkeley, CA (United States); 2014 Mar 17.
- 760 [55] Orozco-Arias S, Isaza G, Guyot R. Retrotransposons in plant genomes: structure,
761 identification, and classification through bioinformatics and machine learning. *Int J*
762 *Mol Sci*. 2019;20(15):3837.
- 763 [56] González LG, Deyholos MK. Identification, characterization and distribution of
764 transposable elements in the flax (*Linum usitatissimum* L.) genome. *BMC Genomics*.
765 2012;13(1):1-7.
- 766 [57] Sun HY, Dai HY, Zhao GL, Ma Y, Ou CQ, Li H, Li LG, Zhang ZH. Genome-wide
767 characterization of long terminal repeat-retrotransposons in apple reveals the
768 differences in heterogeneity and copy number between Ty1-*copia* and Ty3-*gypsy*
769 retrotransposons. *J Integr Plant Biol*. 2008;50(9):1130-9.
- 770 [58] Meyers BC, Tingey SV, Morgante M. Abundance, distribution, and transcriptional
771 activity of repetitive elements in the maize genome. *Genome Res*. 2001;11(10):1660-
772 76.
- 773 [59] Liu Q, Li X, Zhou X, Li M, Zhang F, Schwarzacher T, Heslop-Harrison JS. The
774 repetitive DNA landscape in *Avena* (Poaceae): chromosome and genome evolution
775 defined by major repeat classes in whole-genome sequence reads. *BMC Plant Biol*.

- 776 2019;19(1):1-7.
- 777 [60] Wicker T, Gundlach H, Spannagl M, Uauy C, Borrill P, Ramírez-González RH, De
778 Oliveira R, Mayer KF, Paux E, Choulet F. Impact of transposable elements on genome
779 structure and evolution in bread wheat. *Genome Biol.* 2018;19(1):1-8.
- 780 [61] Ebrahimzadegan R, Houben A, Mirzaghaderi G. Repetitive DNA landscape in
781 essential A and supernumerary B chromosomes of *Festuca pratensis* Huds. *Sci Rep.*
782 2019;9(1):1-1.
- 783 [62] Vitte C, Fustier MA, Alix K, Tenailon MI. The bright side of transposons in crop
784 evolution. *Brief Funct Genom.* 2014;13: 276–95.
- 785 [63] Macas J, Neumann P. Ogre elements—a distinct group of plant Ty3/gypsy-like
786 retrotransposons. *Gene.* 2007;390(1-2):108-16.
- 787 [64] Domingues DS, Cruz GM, Metcalfe CJ, Nogueira FT, Vicentini R, de S Alves C, Van
788 Sluys MA. Analysis of plant LTR-retrotransposons at the fine-scale family level
789 reveals individual molecular patterns. *BMC Genomics.* 2012;13(1):1-4.
- 790 [65] Macas J, Novak P, Pellicer J, Čížková J, Koblížková A, Neumann P, Fukova I, Doležel
791 J, Kelly LJ, Leitch IJ. In depth characterization of repetitive DNA in 23 plant genomes
792 reveals sources of genome size variation in the legume tribe Fabaeae. *PloS one.*
793 2015;10(11):e0143424.
- 794 [66] Bousios A, Kourmpetis YA, Pavlidis P, Minga E, Tsaftaris A, Darzentas N. The
795 turbulent life of Sirevirus retrotransposons and the evolution of the maize genome:
796 more than ten thousand elements tell the story. *Plant J* 2012;69(3):475-88.
- 797 [67] Kolano B, Bednara E, Weiss-Schneeweiss H. Isolation and characterization of reverse
798 transcriptase fragments of LTR retrotransposons from the genome of *Chenopodium*

- 799 *quinoa* (Amaranthaceae). Plant Cell Rep. 2013;32(10):1575-88.
- 800 [68] Natali L, Cossu RM, Mascagni F, Giordani T, Cavallini A. A survey of Gypsy and
801 Copia LTR-retrotransposon superfamilies and lineages and their distinct dynamics in
802 the *Populus trichocarpa* (L.) genome. Tree Genet Genomes. 2015;11(5):1-3.
- 803 [69] Weber B, Heitkam T, Holtgräwe D, Weisshaar B, Minoche AE, Dohm JC,
804 Himmelbauer H, Schmidt T. Highly diverse chromoviruses of *Beta vulgaris* are
805 classified by chromodomains and chromosomal integration. Mob DNA 2013;4(1):1-6.
- 806 [70] Dodsworth S, Chase MW, Kelly LJ, Leitch IJ, Macas J, Novák P, Piednoël M, Weiss-
807 Schneeweiss H, Leitch AR. Genomic repeat abundances contain phylogenetic signal.
808 Syst Biol. 2015;64(1):112-26.
- 809 [71] Vitales D, Garcia S, Dodsworth S. Reconstructing phylogenetic relationships based on
810 repeat sequence similarities. Mol Phylogenetics Evol. 2020;147:106766.
- 811 [72] Benabdelmouna A, Darmency H. Copia-like retrotransposons in the genus *Setaria*:
812 Sequence heterogeneity, species distribution and chromosomal organization. Plant Syst
813 Evol. 2003;237(3):127-36.
- 814 [73] Xiong Y, Eickbush TH. Origin and evolution of retroelements based upon their reverse
815 transcriptase sequences. EMBO J. 1990;9(10):3353-62.
- 816 [74] Capy P, Maisonhaute C. Acquisition/loss of modules: the construction set of
817 transposable elements. Russ J Genet. 2002;38(6):594-601.
- 818 [75] Ustyantsev K, Novikova O, Blinov A, Smyshlyaev G. Convergent evolution of
819 ribonuclease H in LTR retrotransposons and retroviruses. Mol Biol Evol.
820 2015;32(5):1197-207.
- 821 [76] Gorinšek B, Gubenšek F, Kordiš D. Evolutionary genomics of chromoviruses in

- 822 eukaryotes. *Mol Biol Evol.* 2004;21(5):781-98.
- 823 [77] Kordiš D. A genomic perspective on the chromodomain-containing retrotransposons:
824 Chromoviruses. *Gene.* 2005;347(2):161-73.
- 825 [78] Novikova O. Chromodomains and LTR retrotransposons in plants. *Commun Integr*
826 *Biol.* 2009;2(2):158-62.
- 827 [79] Gao X, Hou Y, Ebina H, Levin HL, Voytas DF. Chromodomains direct integration of
828 retrotransposons to heterochromatin. *Genome Res.* 2008;18(3):359-69.
- 829 [80] Nagaki K, Song J, Stupar RM, Parokonny AS, Yuan Q, Ouyang S, Liu J, Hsiao J,
830 Jones KM, Dawe RK, Buell CR. Molecular and cytological analyses of large tracks of
831 centromeric DNA reveal the structure and evolutionary dynamics of maize
832 centromeres. *Genetics.* 2003;163(2):759-70.
- 833 [81] Neumann P, Navrátilová A, Koblížková A, Kejnovský E, Hřibová E, Hobza R,
834 Widmer A, Doležel J, Macas J. Plant centromeric retrotransposons: a structural and
835 cytogenetic perspective. *Mob DNA.* 2011;2(1):1-6.
- 836 [82] Martinez G. tRNAs as primers and inhibitors of retrotransposons. *Mob Genet Element.*
837 2017;7(5):1-6.
- 838 [83] Cullen H, Schorn AJ. Endogenous Retroviruses Walk a Fine Line between Priming
839 and Silencing. *Viruses.* 2020;12(8):792.
- 840 [84] Lu Z, Marand AP, Ricci WA, Ethridge CL, Zhang X, Schmitz RJ. The prevalence,
841 evolution and chromatin signatures of plant regulatory elements. *Nat Plants.*
842 2019;5(12):1250-9.
- 843 [85] Sharma A, Wolfgruber TK, Presting GG. Tandem repeats derived from centromeric
844 retrotransposons. *BMC Genomics.* 2013;14(1):1-1.

- 845 [86] Cheng ZJ, Murata M. A centromeric tandem repeat family originating from a part of
846 Ty3/gypsy-retroelement in wheat and its relatives. *Genetics*. 2003;164(2):665-72.
- 847 [87] Tek AL, Song J, Macas J, Jiang J. Sobo, a recently amplified satellite repeat of potato,
848 and its implications for the origin of tandemly repeated sequences. *Genetics*.
849 2005;170(3):1231-8.
- 850 [88] Kejnovsky E, Kubat Z, Macas J, Hobza R, Mracek J, Vyskot B. Retand: a novel family
851 of gypsy-like retrotransposons harboring an amplified tandem repeat. *Mol Genet*
852 *Genom*. 2006;276(3):254-63.
- 853 [89] Devos KM, Brown JK, Bennetzen JL. Genome size reduction through illegitimate
854 recombination counteracts genome expansion in *Arabidopsis*. *Genome Res*.
855 2002;12(7):1075-9.
- 856 [90] Vitte C, Panaud O. Formation of solo-LTRs through unequal homologous
857 recombination counterbalances amplifications of LTR retrotransposons in rice *Oryza*
858 *sativa* L. *Mol Biol Evol*. 2003;20(4):528-40.
- 859 [91] Ma J, Bennetzen JL. Rapid recent growth and divergence of rice nuclear genomes.
860 *Proc Natl Acad Sci*. 2004 Aug 24;101(34):12404-10.
- 861 [92] Sanchez DH, Gaubert H, Drost HG, Zabet NR, Paszkowski J. High-frequency
862 recombination between members of an LTR retrotransposon family during
863 transposition bursts. *Nat Commun*. 2017;8(1):1-7.
- 864 [93] Vicient CM, Jääskeläinen MJ, Kalendar R, Schulman AH. Active retrotransposons are
865 a common feature of grass genomes. *Plant Physiol*. 2001;125(3):1283-92.
- 866 [94] Cossu RM, Buti M, Giordani T, Natali L, Cavallini A. A computational study of the
867 dynamics of LTR retrotransposons in the *Populus trichocarpa* genome. *Tree Genet*

- 868 Genomes. 2012;8(1):61-75.
- 869 [95] Wos G, Choudhury RR, Kolář F, Parisod C. Transcriptional activity of transposable
870 elements along an elevational gradient in *Arabidopsis arenosa*. Mob DNA.
871 2021;12(1):1-2.
- 872 [96] Qiu F, Ungerer MC. Genomic abundance and transcriptional activity of diverse gypsy
873 and copia long terminal repeat retrotransposons in three wild sunflower species. BMC
874 Plant Biol. 2018;18(1):1-8.
- 875 [97] Bousios A, Darzentas N, Tsaftaris A, Pearce SR. Highly conserved motifs in non-
876 coding regions of Sirevirus retrotransposons: the key for their pattern of distribution
877 within and across plants?. BMC Genomics. 2010;11(1):1-4.
- 878 [98] Wollrab C, Heitkam T, Holtgräwe D, Weisshaar B, Minoche AE, Dohm JC,
879 Himmelbauer H, Schmidt T. Evolutionary reshuffling in the Errantivirus lineage Elbe
880 within the *Beta vulgaris* genome. Plant J. 2012;72(4):636-51.
- 881 [99] Schmidt N, Seibt KM, Weber B, Schwarzacher T, Schmidt T, Heitkam T. Broken,
882 silent, and in hiding: Tamed endogenous pararetroviruses escape elimination from the
883 genome of sugar beet (*Beta vulgaris*). Ann Bot. 2021; mcab042.
- 884 [100] Meagher TR, Vassiliadis C. Phenotypic impacts of repetitive DNA in flowering
885 plants. New Phytol. 2005;168(1):71-80.
- 886 [101] Ito H. Small RNAs and transposon silencing in plants. Dev. Growth Differ.
887 2012;54(1):100-7.
- 888

889 Table 1. General structural features and repeat compositions of full-length LTR retrotransposons. Classification of the in silico retrotransposons
 890 was performed according to Neumann *et al.* 2019 [25].

Class / Superfamily	Lineage (group) superclade clade subclade	Total full- length elements	Genome proportion ¹ (%)	Total element size (kb)	LTR size (bp)	PBS (tRNA types) ²	Range of pairwise similarity (%)	
							LTR	RT
Ty1- <i>copia</i>	Ale	5	0.02-0.6	5.2-4.9	401-127	Asn, Met	17-41	58-72
	Alesia	2	0.01-0.1	5-5.4	289-562	n.d.	54	63
	Bianca	1	0.1	7	313-326	n.d.		
	Ivana	3	0.02-0.2	4.9-5.5	417-438	Met	41.5-47.4	62-69
	TAR	6	0.1-0.3	5.5-6.5	585-994	Glu, Ser, Met, Tyr, Trp	21-41	63-74
	Angela	1	0.21	6.8	796	Met		
	Tork	7	0.1-0.3	4.8-5.6	329-773	Met	17-43	57-70
	SIRE	1	0.2	7	1,545	Met		
	Total	26	4.3					
Ty3- <i>gypsy</i>	non-chromovirus OTA Tat Ogre	7	0.1-1.3	11.5-19.9	761-4560	Asn, Arg, Thr	10.9-28.9	55.4- 73.7
	non-chromovirus OTA Tat TatV	6	0.1-1.0	10.5-11.7	646-1646	Lys, Arg	20.3-35.6	62.3- 74.7
	non-chromovirus OTA Athila	5	0.01-1.9	5.7-9.8	854-1403	Asp	8.7-49.6	69-81.8
	chromovirus Tekay	4	0.03-0.2	5.4-8.6	1149-1925	Gln	49.4	87.5 67.1-
	chromovirus CRM	4	0.02-0.2	5.3-5.5	361-556	Met	47.6-54.3	75.9
	chromovirus Galadriel	2	0.1	5.3-6.4	465-1061	Met	54.5	67.6
Total	28	8.4						

¹Genome proportion = range of genome proportion for each specific reconstructed full-length elements
²n. d. = not detected.

891 Table 2: Summary of tandem repeats and *cis*-regulatory elements associated with the reconstructed full-length LTR-retrotransposon sequences.

Lineage /clade /subclade	Location of tandem repeats	Consensus period size	Number of repetitions	Total number of different <i>cis</i>-regulatory motifs identified in 5' LTR region
Ale	5' and 3' LTR, 5' UTR	2-48	1.9-7.9	10
Alesia	5 LTR, 5' UTR	30	2.2-4.9	5
Bianca	5'UTR	25	2.8	2
Ivana	5'LTR	27	2.9	8
TAR	5'LTR	10-17	1.9-4.4	19
Angela	5'LTR	18	5.4	9
Tork	5'and 3' LTR, 5' UTR	11-49	2-6.8	19
SIRE	5' LTR	15	2	12
Ogre	5' and 3' LTR, UTR adjacent to the <i>gag</i> gene	6-91	2-14.2	38
TatV	5' and 3' LTR, UTR	15-93	2.1-15.1	24
Athila	3' LTR, UTR	21-42	2-2.5	26
Tekay	5' and 3' LTR, UTR	2-29	1.9-18.5	10
CRM	5' and 3' LTR, UTR	17-48	1.9-4.5	8
Galadriel	n.d.	n.d.	n.d.	13

n.d. = no tandem repeat detected

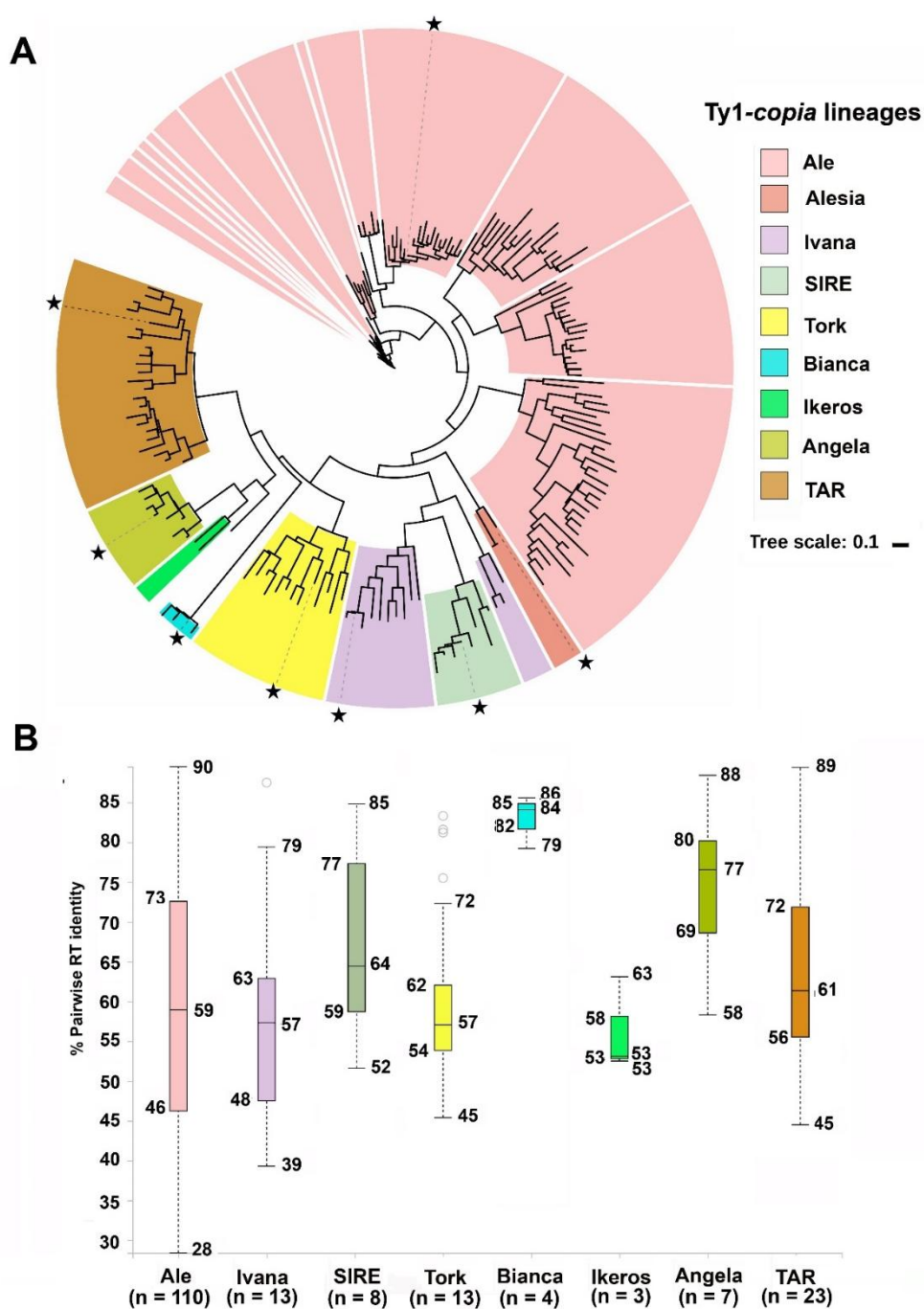


Fig. 1. Diversity of the different Ty1-copia lineages in the *V. macrocarpon* genome. A) The dendrogram was derived from 181 Ty1-copia RT amino acid sequences calculated with the approximate Maximum-Likelihood method using the FastTree tool (Price *et al.* 2010) [45]. Lineages are color-coded (see internal legend). The stars indicate the position of the representative reconstructed full-length elements from *V. macrocarpon* genome shown in Fig. 3A. B) The boxplots illustrate the pairwise sequence identities of amino acid sequences of all Ty1-copia RTs. The lineage Alesia was excluded from the boxplots as only a single copy was detected in the *V. macrocarpon* genome. Colors in the boxplots correspond to panel (A). For each lineage, the total number of RT sequences (n) is provided, which includes the reference sequence.

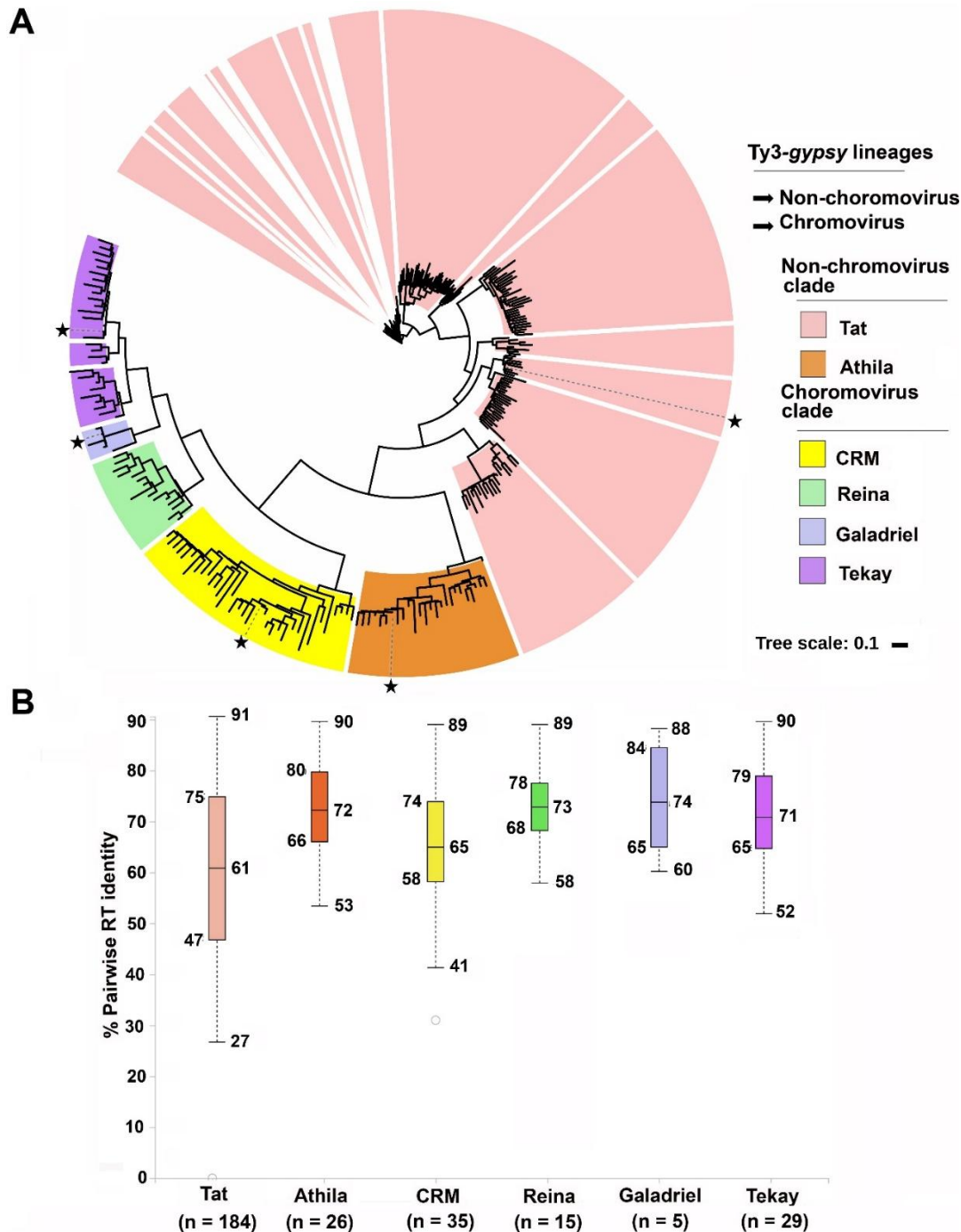


Fig. 2. Diversity of the different Ty3-gypsy lineages in the *V. macrocarpon* genome. A) The dendrogram was derived from 294 Ty3-gypsy RT amino acid sequences calculated with the approximate Maximum-Likelihood method with the FastTree tool [45]. Clades are color-coded (see internal legend). The stars indicate the position of the representative for the reconstructed full-length elements from *V. macrocarpon* genome showed in Fig. 4A. B) The boxplots illustrate the pairwise sequence identities of amino acid sequences of all Ty3-gypsy RTs. Colors in the boxplots correspond to panel (A). For each clade, the total number of RT sequences (n) is provided, which includes the reference sequence

Representative Ty1-*copia* elements

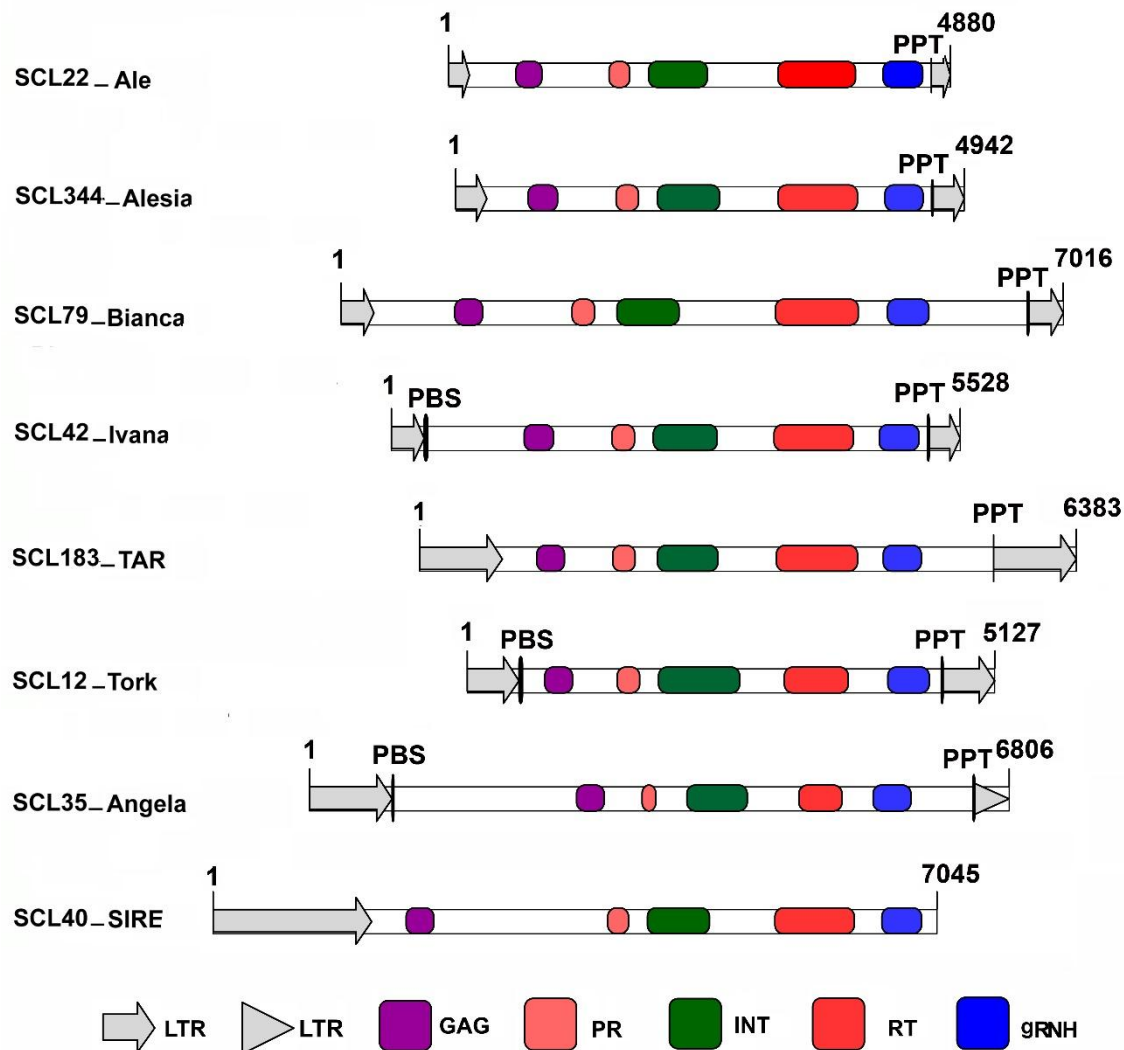


Fig. 3. Structural features of representative *in silico* elements of Ty1-*copia* LTR retrotransposons in the *V. macrocarpon* genome. Long terminal repeats (LTRs) are represented as grey open arrows (intact LTR) and grey arrowheads (truncated LTR). Only the SCL40_SIRE consensus is incomplete and lacks the 3' LTR. Black vertical lines adjacent to the LTRs represent primer binding site (PBS) and polypurine tracts (PPT), while thickness is according to length of this region (Table S2). Structural features of *gag* domain and the four genes of the *pol* domain are depicted as color-coded boxes: GAG = *gag* domain, PR = protease, RT = reverse transcriptase, gRNH = *gypsy*-type RNase H, INT = integrase.

Representative Ty3-gypsy elements

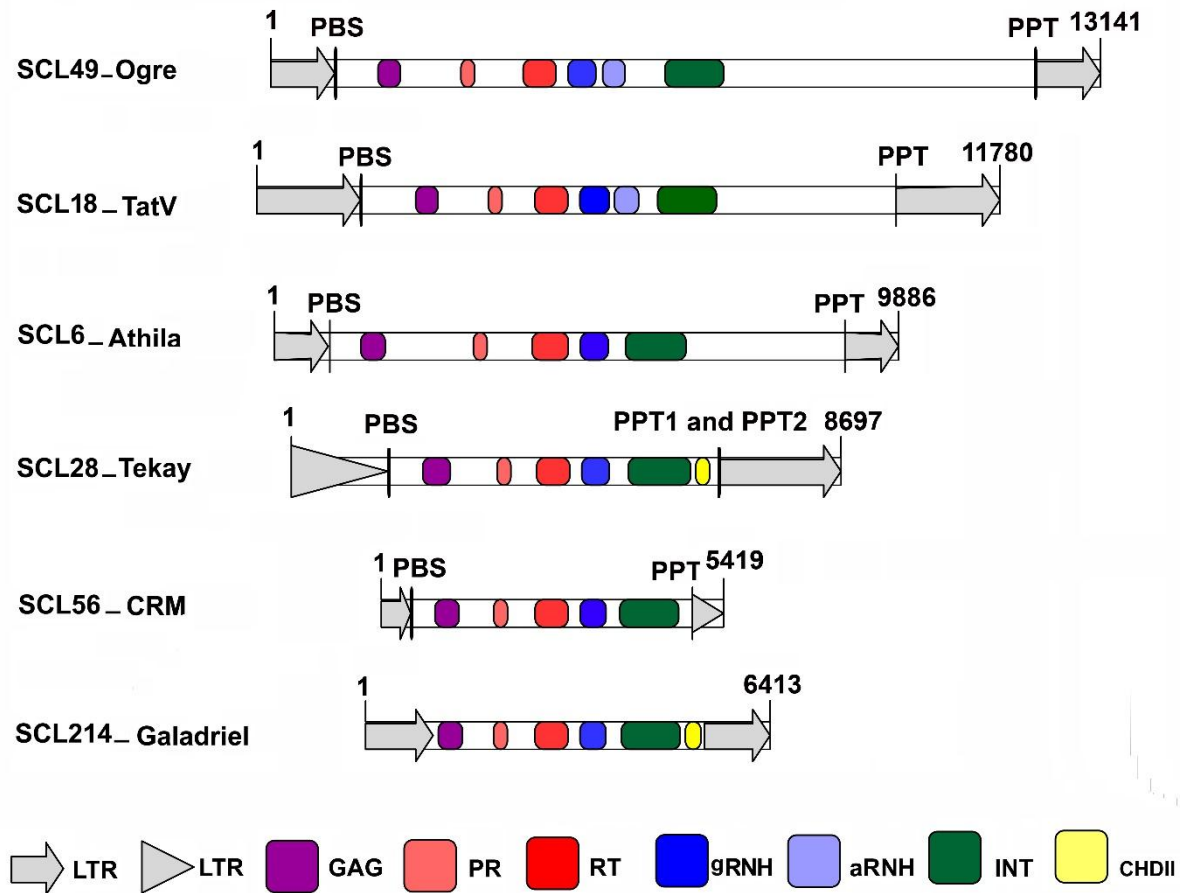


Fig. 4. Structural features of full-length representative consensus elements of Ty3-gypsy LTR retrotransposons in the *V. macrocarpon* genome. Long terminal repeats (LTRs) are represented as grey open arrows (intact LTR) and grey arrowheads (truncated LTR). Black vertical lines adjacent to the LTRs represent primer binding site (PBS) and polypurine tracts (PPT), while thickness is according to length of this region (Table S2). Protein domains encoded in *gag* and *pol* are shown within the LTRs border. Structural features of the *gag* and *pol* genes are depicted as color-coded boxes: GAG = *gag* gene, PR = protease, RT = reverse transcriptase, gRNH = gypsy RNase H, aRNH = archaeal RNase H, INT = integrase, CHDII = chromodomain II.

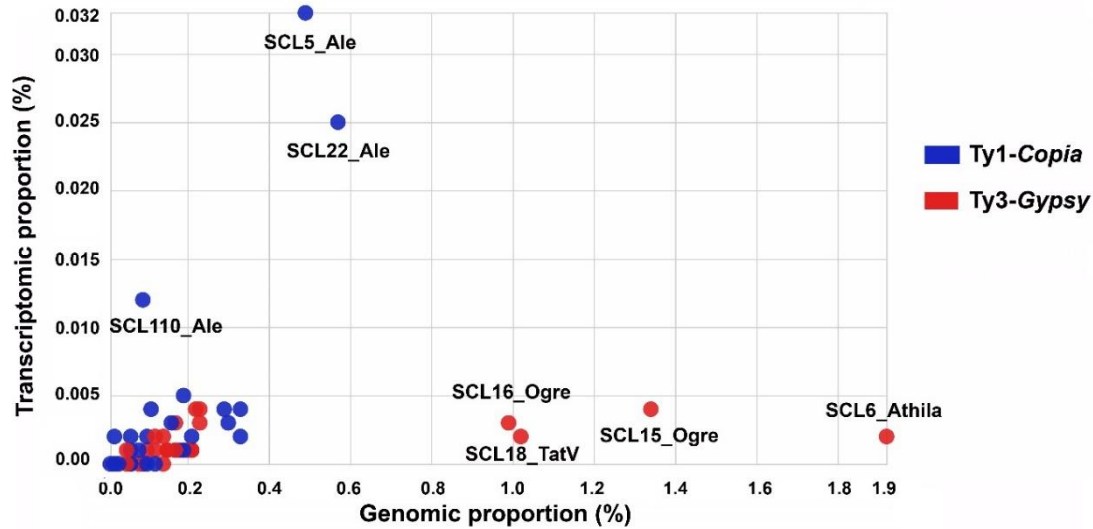


Fig. 5. Comparison of genome and transcriptome proportions for the 26 Ty1-*copia* and 28 Ty3-*gypsy* full-length consensus retrotransposons from *V. macrocarpon*. The genome proportions of the reconstructed full-length Ty1-*copia* and Ty3-*gypsy* elements are plotted against their transcriptome proportions. Genome proportions were calculated from the RepeatExplorer output (Table S2). Transcriptome proportions were calculated from read mapping of the publicly available Illumina cDNA sequence reads of *V. macrocarpon* (accession number PRJNA246586). Major groups of the 54 elements are color-coded (see internal legend). Only the names of elements of the highest genome and/or transcriptome proportions are annotated. For details, see Table S2.