

PlantRep: a resource of plant repeats

Xizhi Luo

Chinese Academy of Agricultural Sciences Agricultural Genomes Institute at Shenzhen

Shiyu Chen

Chinese Academy of Agricultural Sciences Agricultural Genomes Institute at Shenzhen

Yu Zhang (✉ zhangyu07@caas.cn)

Chinese Academy of Agricultural Sciences Agricultural Genomes Institute at Shenzhen

<https://orcid.org/0000-0001-6547-6243>

Research Article

Keywords: database of repeat sequences, evolution of plant genome, transposable element

Posted Date: October 28th, 2021

DOI: <https://doi.org/10.21203/rs.3.rs-498874/v2>

License:   This work is licensed under a Creative Commons Attribution 4.0 International License.

[Read Full License](#)

1 **PlantRep: a resource of plant repeats**

2 Xizhi Luo¹, Shiyu Chen¹, Yu Zhang^{1,2,*}

3

4 **Key Message**

5 We re-annotated repeats of 459 plant genomes and released a new database: PlantRep
6 (<http://www.plantrep.cn/>). PlantRep sheds lights of repeat evolution and provides fundamental data
7 for deep exploration of genome.

8

9 **Key words:** database of repeat sequences, evolution of plant genome, transposable element.

10

11 ¹ Shenzhen Branch, Guangdong Laboratory for Lingnan Modern Agriculture, Genome Analysis
12 Laboratory of the Ministry of Agriculture, Agricultural Genomics Institute at Shenzhen, Chinese
13 Academy of Agricultural Sciences, Shenzhen 518124, China

14 ² School of Agriculture, Sun Yat-sen University, Shenzhen 518107, China

15 Email: zhangyu07@caas.cn

16

17

18 Introduction

19 With the rapid release of genome, exploring the comparative genomics of repeats enabled us
20 to elucidate how repeat sequences originated and amplified in different plant lineages. However,
21 large-scale evolution analysis in plants takes a considerable amount of time and computing
22 resources, hence the comparative genomics and evolution analysis of transposons have only been
23 conducted in a few selected plants¹⁻³ and the current reference repeat databases only contain a few
24 model plants⁴⁻⁶. Besides, the repeat annotations carried out independently using heterogeneous
25 pipelines cannot be used directly for comparative studies. To cope with this issue, we utilized a
26 uniformed pipeline to re-annotate repeats for 459 plant genomes and compared the repeat sequences
27 among the plant groups, including the composition, family diversity, genomic distribution and
28 evolutionary rate. The results provide a resource for the analysis and study of the repeat sequences
29 in different lineages of plants.

30 Result and Discussion

31 We re-annotated repeats from 459 released plant genomes and generated 45.72 Gb seed
32 alignments and 601,731 consensus sequences of repeats from *de novo* repeat annotation of each
33 plant genome. The repeat libraries are available in our database PlantRep. Combined with the
34 reference-based annotation, 206.04 Gb of 396,041,410 repeats were identified and categorized. We
35 released this consensus repeat annotations (PlantRep) for the plant community as an updated
36 resource for the future data-mining studies. Repeats in the PlantRep database were categorized with
37 repeat types adapted from the existing eukaryotic transposable element classes and the Dfam
38 database^{6,7} (Supplementary Table 1c). Retrotransposon includes Long Terminal Repeat (LTR),
39 Inverted Long Terminal Repeat (DIRS), Penelope (PLE), Long Interspersed Nuclear Element
40 (LINE), and Short LINE-dependent Retroposons (SINE). DNA transposons were categories into
41 Terminal Inverted Repeat (TIR), Circular dsDNA Intermediate (CirdsDNA), DNA Polymerase (DP),
42 and Circular ssDNA Intermediate (Rolling Circle, RC). Besides, low complexity, satellite and
43 simple repeat were also included.

44 To examine the diversity of repeat families, the 459 plant species were divided into 15 clades
45 based on their phylogeny: algae, bryophyte, lycophytes, fern, gymnosperms, ANA (early
46 angiosperms), magnolids, monocots, base eudicots, super rosids, fabids, malvid, super asterids,
47 lamiids, and campanulids⁸ (Supplementary Table 1b; Supplementary Fig. 2). The abundance and
48 diversity of each repeat type in each clade were characterized.

49 We examined the distribution of different types of repeat sequences in each clade
50 (Supplementary table 1c). The average percentage of repeats within the genome across all species
51 was 45.49%. The top five abundant types of repeats are LTR, TIR, LINE, Simple repeat and Rolling
52 Circle, accounting for 21.66%, 5.44%, 2.25%, 1.54%, and 0.59% of the plant genome on average
53 (Fig 1a; Supplementary table 2), respectively. Plants from different lineages display distinct
54 proportion of repeat types. In general, the proportion of repeats increased from algae, bryophyte,
55 lycophytes, fern to gymnosperms (Supplementary Fig. S2-S15). LTRs, as the largest family of plant
56 repeats, might be the major contributor for the increase of total TE population. For algae, the
57 proportions of simple repeats and LINEs exceeded LTRs, which were the highest in all the lineages,
58 suggesting a possible mechanism of controlling the LTR amplification in algae. Ferns showed a
59 higher proportion of SINE transposons. In ANA, TIRs accounted for a prominent proportion of
60 repeats. The diverse compositions of the repeat types among plant species could provide a source
61 for their unique genome evolution trajectory.

62 To trace back the evolutionary history of plant transposons, we analyzed the presence/absence
63 of types of repeats (Fig 1b; Supplementary table 3; Supplementary Fig. S16, S17). We found that
64 most of the repeat families existed in algae, indicating the common ancestor of algae and land plants
65 had already evolved the fundamental transposon layout of modern green plants (Fig 1b).

66 The diversity of transposon nucleotide can reflect the evolution rate of transposon to some
67 extent, we investigate the nucleotide diversity of transposons for each plant lineage (Fig 1c;
68 Supplementary table 4). In plants, the main nucleotide diversity peaks of the LTR concentrates is
69 at 16%, and the TIR is about 20%, but the LINE concentrates is around 25%, indicating that LINE
70 lack recently replication activity compared to LTR and TIR. The result to some extent explains the
71 genome different between plant and animal: the proportion of LINE in the plant genome is usually

72 less than 5%, but in some animal, LINE is the main repeat content of genome. As the main
73 contributor to plant genome, LTRs were selected to investigate the evolutionary history of
74 amplification (Fig 1c; Supplementary table 4). The main peak of algae is 0.2–0.23, which is smaller
75 than the peaks of all of the other species, implying that algae carried ancient LTR groups. LTR
76 amplified more recently in bryophytes compared with algae. The nucleotide diversity of LTR in
77 ferns and lycophytes showed the lowest diversity across all the land plants. The amplification of
78 LTR in the genome of gymnosperms leads to a large genome, and the results show that the degree
79 of divergence of LTR is high, indicating that it has no recent activity, which is consistent with the
80 results of Norway spruce. The nucleotide diversity of ANA and magnoliids is similar as that of
81 gymnosperms. Monocots and base eudicots displayed a more recent amplification. For dicots, the
82 nucleotide diversity was broad, indicating several rounds of amplification of LTRs. According to
83 the divergence of LTRs, we also estimated the amplification time of LTRs along the plant, with the
84 main detected amplicon that can be traced back to 1-4 Mya ago (Supplementary fig. S18). This
85 indicates that the LTRs amplify independently in each lineage, playing important roles in the
86 evolution of genome size and environmental adaptability.

87 To elucidate the contribution of repeats to genes, we calculate the frequencies of repeat
88 sequences at different sites around genes. We found that LTRs, LINEs, SINEs, and DNA transposons
89 display decrease of frequency from 1 kb upstream to 1kb downstream of transcriptional start site
90 genes (Supplementary Fig S19; Supplementary table 5), indicating that plants tend to suppress
91 transposon insertions around gene transcription start sites (TSS), the transposons located near the
92 gene might impact the expression and function of gene. The common feature of LTRs and LINEs
93 near the gene is that there is an inflection point which falls sharply from 2 kb upstream of the TSS
94 to the lowest frequency at 1 kb within genes TSS (Supplementary Fig S19-S21). It then continues
95 to rise within the gene to a distribution of 10k. The frequency of 10kb within gene is close to or
96 even higher than the frequency at 10kb upstream of the TSS. Similarly, the frequency of SINE and
97 DNA transposons from 1kb upstream of the TSS to 1kb within the gene decreases, in gene internal
98 frequency rises; but then the frequency shows a downward or stable trend within gene
99 (Supplementary Fig S19, S22), which is different from LTRs and LINEs. The result implies different
100 transposon families might adopt specific integration strategies and occupy different “niches” of

101 genome⁹. Unlike TEs, the frequency of simple repeats near TSS is opposite. Simple repeats increase
102 from 4.5 kb upstream of TTS to 0.5 kb within the gene where it reaches the highest frequency
103 (Supplementary Fig S19) , which is similar to the results of before study¹⁰. Therefore, one can
104 speculate that the high frequency of simple repeats around the gene provides a certain fault tolerance
105 rate for the stability of gene transcription.

106 In summary, we re-annotated repeats of 459 plant species and characterized the abundance,
107 presence/absence and nucleotide diversity of the repeat types for 15 plant taxonomic groups. The
108 frequency of repeats along the gene models showed unique patterns for different repeat types. Our
109 work supplies a new resource for the future study of repeat sequences and will be helpful to plant
110 genome structure annotation.

111 **Declarations**

112 *Funding*

113 This work was supported by the National Natural Science Foundation of China (32070250), the
114 Natural Science Foundation of Guangdong Province (2020A1515011030) and the open research
115 project of “Cross-Cooperative Team” of the Germplasm Bank of Wild Species, Kunming Institute
116 of Botany, Chinese Academy of Sciences.

117 *Competing interests*

118 The authors declare that they have no competing interests.

119 *Availability of data and materials*

120 The datasets generated and/or analysed during the current study are available in the website:
121 <http://www.plantrep.cn/>.

122 *Code availability*

123 The code used during the current study available from the corresponding author on reasonable
124 request.

125 *Ethics approval and*

126 Not applicable.

127 *Consent to participate*

128 Not applicable.

129 *Consent for publication*

130 Not applicable.

131 *Authors' contributions*

132 X.L. and Y.Z. designed the study. X.L. and Y.Z. wrote the manuscript. S.C. reviewed and edited the
133 manuscript. X.L. analyzed data. Y.Z. oversaw the study.

134 **Acknowledgements**

135 We thank Sanwen Huang and Xinyan Zhang for providing suggestions, Shujun Ou for discussion
136 on intact LTR insert time calculation. Computational support was provided by the Shenzhen
137 Branch, Guangdong Laboratory for Lingnan Modern Agriculture and Genome Analysis
138 Laboratory of the Ministry of Agriculture, Agricultural Genomics Institute at Shenzhen, Chinese
139 Academy of Agricultural Sciences, Shenzhen, China. This work was supported by the National
140 Natural Science Foundation of China (32070250), the Natural Science Foundation of Guangdong
141 Province (2020A1515011030) and the open research project of “Cross-Cooperative Team” of the
142 Germplasm Bank of Wild Species, Kunming Institute of Botany, Chinese Academy of Sciences.

Reference

- 144 1 El Baidouri, M. & Panaud, O. Comparative genomic paleontology across plant kingdom reveals
145 the dynamics of TE-driven genome evolution. *Genome biology and evolution* **5**, 954-965,
146 doi:10.1093/gbe/evt025 (2013).
- 147 2 Elliott, T. A. & Gregory, T. R. Do larger genomes contain more diverse transposable elements?
148 *Bmc Evol Biol* **15**, 69, doi:10.1186/s12862-015-0339-8 (2015).
- 149 3 Schaper, E. & Anisimova, M. The evolution and function of protein tandem repeats in plants.
150 *New Phytol* **206**, 397-410, doi:10.1111/nph.13184 (2015).
- 151 4 Jurka, J. *et al.* Repbase Update, a database of eukaryotic repetitive elements. *Cytogenetic and*
152 *genome research* **110**, 462-467, doi:10.1159/000084979 (2005).
- 153 5 Bao, W., Kojima, K. K. & Kohany, O. Repbase Update, a database of repetitive elements in
154 eukaryotic genomes. *Mobile DNA* **6**, 11, doi:10.1186/s13100-015-0041-9 (2015).
- 155 6 Hubley, R. *et al.* The Dfam database of repetitive DNA families. *Nucleic Acids Res* **44**, D81-89,
156 doi:10.1093/nar/gkv1272 (2016).
- 157 7 Wicker, T. *et al.* A unified classification system for eukaryotic transposable elements. *Nat Rev*
158 *Genet* **8**, 973-982, doi:10.1038/nrg2165 (2007).
- 159 8 Group, T. A. P. An update of the Angiosperm Phylogeny Group classification for the orders and
160 families of flowering plants: APG IV. *Botanical Journal of the Linnean Society* **181**, 1-20,
161 doi:10.1111/boj.12385 (2016).
- 162 9 Zhang, X., Zhao, M., McCarty, D. R. & Lisch, D. Transposable elements employ distinct
163 integration strategies with respect to transcriptional landscapes in eukaryotic genomes. *Nucleic*
164 *Acids Res*, doi:10.1093/nar/gkaa370 (2020).
- 165 10 Huda, A., Mariño-Ramírez, L., Landsman, D. & Jordan, I. K. Repetitive DNA elements,
166 nucleosome binding and human gene expression. *Gene* **436**, 12-22,
167 doi:10.1016/j.gene.2009.01.013 (2009).

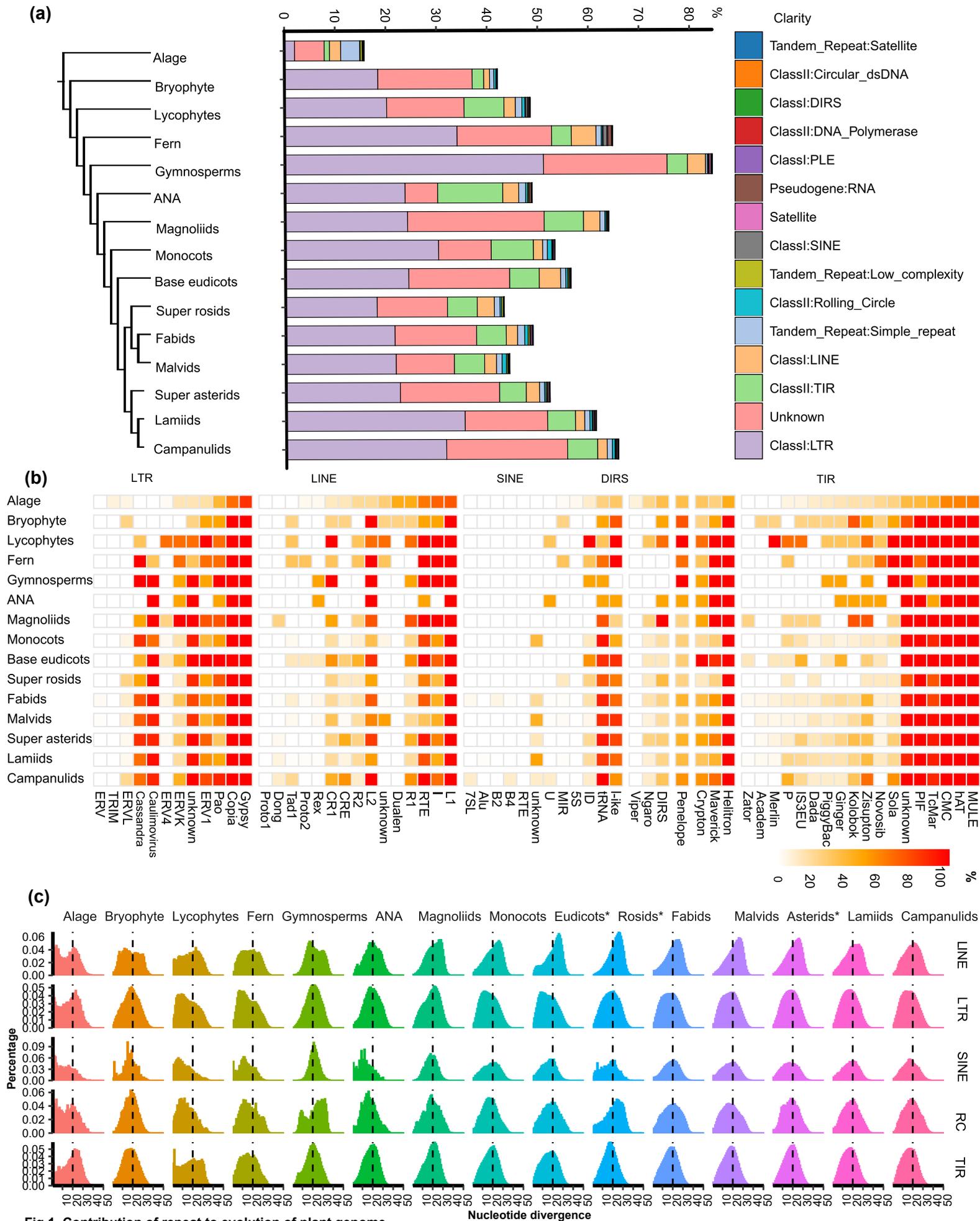


Fig 1. Contribution of repeat to evolution of plant genome.

(a) The average percentage of different types of repeats within the genome of 15 groups from green plant kingdom. The left panel displayed the phylogenetic trees of the plant lineages. (b) The percentage of species carrying certain repeat family in each group. (c) Nucleotide diversity of transposons within each plant lineages. The x-axial label means the nucleotide diversity percentage of repeat. The vertical dashed line represents the divergence rate of 20%. Eudicots*, base eudicots; Rosids*, super rosids; Asterids*, super asterids.

Supplementary Files

This is a list of supplementary files associated with this preprint. Click to download.

- [PlantReparesourceofplantrepeatsSupplementaryFigure.pdf](#)
- [PlantReparesourceofplantrepeatsSupplementaryMaterial.pdf](#)
- [SupplementaryTable1.xlsx](#)
- [SupplementaryTable2.xlsx](#)
- [SupplementaryTable3.xlsx](#)
- [SupplementaryTable4.xls](#)
- [SupplementaryTable5.xls](#)