



UNIVERSITAT
POLITÈCNICA
DE VALÈNCIA

UNIVERSITAT POLITÈCNICA DE VALÈNCIA

Plant Molecular and Cellular Biology Joint Research
Institute (IBMCP)

Study of the evolution of the transposable element space in
plant genomes.

Master's Thesis

Master's Degree in Plant Molecular and Cellular Biotechnology

AUTHOR: Zhu, Yujie

Tutor: Mulet Salort, José Miguel

Experimental director: BOMBARELY GOMEZ, AURELIANO

ACADEMIC YEAR: 2023/2024

Contents

INTRODUCTION	1
OBJECTIVES	4
MATERIALS & METHODS	5
1. INPUT DATA ACQUISITION	5
2. GENOME METRICS RETRIEVAL	7
3. REPETITIVE ELEMENT IDENTIFICATION	7
4. COMPARATIVE ANALYSIS OF TRANSPOSABLE ELEMENTS	9
RESULTS	12
1. REPETITIVE ELEMENT IDENTIFICATION	12
2. COMPARATIVE AND ANALYSIS	15
2.1 DISTRIBUTION METRICS AND BOX PLOTS REPRESENTATIONS	15
2.1.1 RETROELEMENT COMPARATIVE AND ANALYSIS	16
2.1.2 DNA TRANSPOSON COMPARATIVE AND ANALYSIS	23
2.1.3 UNCLASSIFIED ELEMENT COMPARATIVE AND ANALYSIS.....	25
2.1.4 SMALL RNA COMPARATIVE AND ANALYSIS	27
2.1.5 SIMPLE REPEAT ELEMENT COMPARATIVE AND ANALYSIS	29
2.1.6 LOW COMPLEXITY ELEMENT COMPARATIVE AND ANALYSIS	32
2.2 HEATMAP REPRESENTATION.....	34
2.3 PRINCIPAL COMPONENT ANALYSIS.....	36
2.4 ESTIMATING DIVERSITY AND SPECIFICITY	37
DISCUSSION	42
1. THE UNIQUE TE LANDSCAPE ACROSS DIVERGENT LINEAGES	42
2. GENOME SIZE AND ITS INFLUENCE ON TE LANDSCAPE	44
3. TE ACTIVITY AND VARIABILITY IN RAPIDLY EVOLVING LINEAGES	45
CONCLUSIONS	47
REFERENCES	48

Introduction

Transposable Elements (TEs) are DNA sequences capable of changing their position within the genome (Bourque et al., 2018). Once considered "junk DNA," the availability of genome sequences and the growth of genomic databases have accelerated the study of TEs, now recognized as evolutionary features (Ramakrishnan et al., 2022).

Based on their transposition mechanisms, TEs are classified into two major classes:

Class I: RNA transposons, also known as Retrotransposons, move through a "copy-paste" mechanism. Initially, they are transcribed from DNA to RNA, and the resulting RNA intermediate is reverse-transcribed into a cDNA copy, which is then inserted into a new position in the genome using TE enzyme machinery (Bourque et al., 2018).

Class II: DNA transposons move directly via a "cut-paste" mechanism without an RNA intermediate. They insert into new genomic locations through a DNA intermediate (Muñoz-López & García-Pérez, 2010). Given that genome size, complexity, and instability are often associated with TE copy number, RNA transposons play a more significant role (Wicker et al., 2007).

These two major classes of transposons are further classified into subclasses based on their chromosomal integration mechanisms. The subclasses for RNA transposons include LTR, DIRS, and Non-LTR. Long Terminal Repeats (LTR) retrotransposons integrate into the genome through an integration enzyme-catalyzed cleavage and strand transfer process, similar to retroviruses (Brown et al., 1987). In contrast, non-LTR integration occurs via a process called target-primed reverse transcription and reverse transcription coupling (Luan et al., 1993).

Each subclass is further divided into superfamilies, with LTR's superfamilies being Copia, Gypsy, and ERV, and Non-LTR's Superfamilies being LINE, SINE, and PLE. These superfamilies are often present across diverse organisms but share common genetic organization and monophyletic origins (Bourque et al., 2018). The two primary superfamilies of LTR retrotransposons, Copia and Gypsy, are found in nearly all major eukaryotic lineages (Malik & Eickbush, 2001). Plant SINEs are derived from tRNAs, which hold a special position, as these small non-coding and non-autonomous elements of several hundred base pairs utilize the transposition mechanism of LINEs to ensure their amplification (Mhiri et al. 2022). Each superfamily is further divided into families. In principle, each TE sequence in the genome can be assigned to a family, superfamily, subclass, and class (Bourque et al., 2018).

Transposable elements play significant roles in organismal evolution and genomic dynamics, with their importance in the evolutionary process attributable to several key properties:

- **Genomic Plasticity and Diversity:** TEs induce changes and reshaping of the genome structure by mobilizing to new positions within the genome or influencing gene expression. Such genomic alterations provide genetic diversity, which is crucial for the evolutionary process.
- **Functional Innovation of Genes:** TEs possess diverse structures and functions, for instance, in plants, they can act as promoters and enhancers (Ramakrishnan et al., 2022). When TEs are inserted into genes or their regulatory regions, they can alter gene functions or expression patterns, occasionally leading to the emergence of novel functions, thereby equipping organisms with adaptability to new environments or lifestyles.
- **Genome Size and Evolution:** The vast and varying sizes of genomes are primarily due to the proliferation of TEs (Haley & Mueller, 2022). TEs are also considered significant contributors to other mechanisms like recombination rates and polyploidy (Mhiri et al., 2022). Transposition is an efficient mechanism for genome expansion, which is counterbalanced over time by DNA removal. The balance between these processes is a driving force for the evolution of eukaryotic genome sizes (Schubert & Vu, 2016).
- **Source of Genetic Variability:** TEs serve as extensive sources of genetic and hereditary polymorphism. TEs occupy a considerable portion of a species' genome, including a significant portion unique to that species (Bourque et al., 2018). The insertion and activation of TEs can trigger genetic variations, laying the groundwork for the effects of natural selection and influencing population adaptability.
- **Genomic Stability and Safeguarding:** To mitigate potential detrimental consequences from excessive transposition, evolutionary processes have developed various epigenetic mechanisms to silence and control TEs, such as small RNAs, KRAB domain-containing zinc finger proteins, DNA methylation, histone modifications, splicing inhibition, and RNA modifications (Almeida et al., 2022). TEs are influenced by multiple regulatory sources, including the TEs themselves (regulatory motifs, biology) and host plant features (epigenetic control, genome size, ploidy level, sequence elimination mechanisms, and reproductive systems). The collective interplay of these evolutionary forces might lead to a state of equilibrium (Mhiri et al., 2022).

To better understand the landscape and dynamics of TEs within genomes, sophisticated computational tools have been developed. In this study, we employed a comprehensive approach utilizing tools such as RepeatModeler, RepeatMasker, and TEsorter for TE re-annotation. These tools facilitate the identification, classification, and masking of repetitive elements, offering insights into their distribution, diversity, and potential functional implications within the genome.

Given the intricate nature and profound impact of Transposable Elements on genomic evolution, these dynamic DNA segments have emerged as focal points of research across various biological disciplines. As our understanding of TEs deepens, it becomes increasingly evident that vast territories of the TE landscape remain uncharted, brimming with potential insights. Recognizing the pressing need to delve further into this captivating realm and bridge existing knowledge gaps, the present study embarks on a rigorous exploration. To unravel the complexities and nuances of TE dynamics across plant lineages, we articulate the following hypotheses:

- *Hypothesis 1:* The most divergent lineages harbor the most unique TE landscape. Given the vast evolutionary distances and distinct ecological niches occupied by different plant lineages, it is plausible to speculate that their TE profiles have evolved in unique directions, leading to distinct TE compositions and activities.
- *Hypothesis 2:* There exists a correlation between genome size and the TE landscape across plant species. Larger genomes may offer more "space" for TE insertions, potentially influencing gene density and spatial organization. Consequently, this may affect the accumulation and distribution patterns of TEs within these genomes.
- *Hypothesis 3:* Rapidly evolving plant lineages may exhibit heightened TE activity or mutation rates. In comparison to more primitive lineages like Chlorophytes, advanced groups like Angiosperms might manifest increased TE diversity and variability, potentially reflecting adaptive responses to diverse environmental pressures.

Objectives

Transposable elements, dynamic segments of DNA within the genome, play significant roles in shaping genomic composition through processes like gene disruption and recombination rate modulation. Our hypothesis posits varying TE diversity among different plant species, confined within specific ranges tied to distinct evolutionary clades. The project aims to explore this diversity across plant lineages, spanning from algae to angiosperms. The selected 138 plant genomes, representing diverse lineages (Chlorophytes, Charophytes, Liverworts, Mosses, Hornworts, Lycophytes, Ferns, Gymnosperms, and Angiosperms), will undergo TE re-annotation using advanced tools such as RepeatModeler2, RepeatMasker, and TEsorter. This comprehensive re-annotation will offer an updated understanding of TE landscapes within each genome.

The overarching goal of the project is to unravel the diversity and specificity parameters of TE landscapes, contributing to our understanding of how TEs have influenced genomic evolution across diverse plant lineages. This research holds significance in shedding light on the functional implications of TEs and their roles in genomic processes, thus contributing to the broader fields of genomics, evolutionary biology, and plant science.

The specific goals of this project are:

- Determine the speciation and diversity of the transposable elements within 138 plant species along the plant tree of life.
- Compare the transposable landscape of these 138 species using common data visualization tools such as boxplots, heat maps, and Principal Component Analysis (PCA).
- Validate the hypothesis that the most divergent lineages possess the most unique TE landscape.
- Investigate the correlation between genome size and the TE landscape across plant species.
- Assess the potential functional implications of observed transposable element variations on specific genomic features or evolutionary traits.
- Provide insights into the evolutionary significance of transposable elements in shaping plant genome architecture and function.

Materials & Methods

1. Input Data Acquisition

To comprehensively explore the landscape of transposable elements in plant genomes, a diverse dataset was curated. Genomic data from 42 Chlorophytes were downloaded from NCBI using the "curl" command-line tool.

SPECIES	LINEAGE	SOURCE	ASSEMBLY
<i>Astrephomene gubernaculifera</i>	Chlorophyte	NCBI	Astre_guber_v1.0
<i>Auxenochlorella protothecoides (0710)</i>	Chlorophyte	NCBI	ASM73321v1
<i>Auxenochlorella protothecoides (UTEX 25)</i>	Chlorophyte	NCBI	ASM370936v1
<i>Bathycoccus prasinus</i>	Chlorophyte	NCBI	ASM222023v1
<i>Chlamydomonas eustigma</i>	Chlorophyte	NCBI	C.eustigma genome v1.0
<i>Chlamydomonas incerta</i>	Chlorophyte	NCBI	ASM1683460v1
<i>Chlamydomonas reinhardtii</i>	Chlorophyte	NCBI	Chlamydomonas_reinhardtii_v5.5
<i>Chlamydomonas schloesseri</i>	Chlorophyte	NCBI	ASM1683459v1
<i>Chlamydomonas sp. UWO 241</i>	Chlorophyte	NCBI	CUWO241_v1.0_nuclear
<i>Chlorella desiccata (nom. nud.) (UTEX 2437)</i>	Chlorophyte	NCBI	LANL_v2
<i>Chlorella desiccata (nom. nud.) (UTEX 2526)</i>	Chlorophyte	NCBI	LANL_Cdes_v1.1
<i>Chlorella ohadii</i>	Chlorophyte	NCBI	ASM2502687v1
<i>Chlorella sorokiniana</i>	Chlorophyte	NCBI	Chlorella_sorokiniana 2.0
<i>Chlorella variabilis</i>	Chlorophyte	NCBI	v 1.0
<i>Chloropicon primus (CCMP 1205)</i>	Chlorophyte	NCBI	ASM785969v1
<i>Chloropicon primus (RCC138)</i>	Chlorophyte	NCBI	ASM2320587v1
<i>Coccomyxa sp. Obi</i>	Chlorophyte	NCBI	COCOBI_1.0
<i>Coccomyxa subellipsoidea C-169</i>	Chlorophyte	NCBI	Coccomyxa_subellipsoidea v2.0
<i>Dunaliella salina</i>	Chlorophyte	NCBI	Dunsal1 v. 2
<i>Edaphochlamys debaryana</i>	Chlorophyte	NCBI	ASM1685814v1
<i>Haematococcus lacustris</i>	Chlorophyte	NCBI	Lacustris_1.0
<i>Helicosporidium sp. ATCC 50920</i>	Chlorophyte	NCBI	Helico_v1.0
<i>Micractinium conductrix</i>	Chlorophyte	NCBI	ASM224581v2
<i>Micromonas commoda</i>	Chlorophyte	NCBI	ASM9098v2
<i>Micromonas pusilla CCMP1545</i>	Chlorophyte	NCBI	Micromonas_pusilla CCMP1545 v2.0
<i>Monoraphidium minutum</i>	Chlorophyte	NCBI	Monmin1
<i>Monoraphidium neglectum</i>	Chlorophyte	NCBI	mono_v1
<i>Ostreobium quekettii</i>	Chlorophyte	NCBI	Ostreobium_ID_genome
<i>Ostreococcus lucimarinu</i>	Chlorophyte	NCBI	ASM9206v1
<i>Ostreococcus tauri (RCC1115)</i>	Chlorophyte	NCBI	Ostta1115_2
<i>Ostreococcus tauri (RCC4221)</i>	Chlorophyte	NCBI	version 140606
<i>Pedinophyceae sp. YPF-701</i>	Chlorophyte	NCBI	Pedinophyceae_YPF-701_genome
<i>Picochlorum sp. BPE23</i>	Chlorophyte	NCBI	ASM2520934v1
<i>Picochlorum sp. BPE23</i>	Chlorophyte	NCBI	ASM2520937v1
<i>Pycnococcus provasolii</i>	Chlorophyte	NCBI	Ppro_1.0
<i>Raphidocelis subcapitata</i>	Chlorophyte	NCBI	Rsub_1.0
<i>Scenedesmus sp. NREL 46B-D3</i>	Chlorophyte	NCBI	Scesp_1
<i>Tetraabaena socialis</i>	Chlorophyte	NCBI	TetSoc1
<i>Trebouxia sp. AI-2</i>	Chlorophyte	NCBI	ASM863618v1
<i>Volvox africanus</i>	Chlorophyte	NCBI	Vafri_1.0
<i>Volvox reticuliferus (NIES 3786)</i>	Chlorophyte	NCBI	Vretimale_1.0
<i>Volvox reticuliferus (NIES-3785)</i>	Chlorophyte	NCBI	Vretifemale_1.0

Genomic data from 96 plant species representing various lineages, including Charophyte (1), Charales (1) Liverworts (2), Mosses (5), Hornworts (1), Lycophytes (4), Ferns (2), Gymnosperms (5), and Angiosperms (75), were downloaded from renowned repositories, including NCBI, CNCBI, owned databases, and SGN.

SPECIES	LINEAGE	SOURCE	ASSEMBLY
<i>Acer negundo</i>	Angiosperm (Eudicots)	NCBI	GCA_025594385.1_ASM2559438v1_genomic.fna
<i>Acer yangbiense</i>	Angiosperm (Eudicots)	NCBI	GCA_008009225.1_AYv1.1_genomic.fna
<i>Actinidia chinensis</i>	Angiosperm (Eudicots)	NCBI	GCA_003024255.1_Red5_PS1_1.69.0_genomic.fna
<i>Adansonia digitata</i>	Angiosperm (Eudicots)	NCBI	GCA_029448705.1_ASM2944870v1_genomic.fna
<i>Adiantum capillus</i>	Fern	NCBI	GCA_014529385.2_ASM1452938v2_genomic.fna
<i>Amborella trichopoda</i>	Angiosperm (Basal)	NCBI	GCF_000471905.2_AMTR1.0_genomic.fna
<i>Ananas comosus</i>	Angiosperm (Monocots)	NCBI	GCF_001540865.1_ASM154086v1_genomic.fna
<i>Annona glabra</i>	Angiosperm (Magnolids)	CNCB	GWHBCKB000000000_Anonna.genome.fna
<i>Anthoceros angustus</i>	Hornwort	NCBI	GCA_010909165.1_ASM1090916v1_genomic.fna
<i>Arabidopsis arenosa</i>	Angiosperm (Eudicots)	NCBI	GCA_905216605.1_AARE701a_genomic.fna
<i>Arabidopsis suecica</i>	Angiosperm (Eudicots)	NCBI	GCA_019202805.1_ASM1920280v1_genomic.fna
<i>Arabidopsis thaliana</i>	Angiosperm (Eudicots)	NCBI	GCF_000001735.4_TAIR10.1_genomic.fna
<i>Arabis alpina</i>	Angiosperm (Eudicots)	NCBI	GCA_900128785.1_MPIPZ.v5_genomic.fna
<i>Boechera stricta</i>	Angiosperm (Eudicots)	NCBI	GCA_018361395.1_NTU_Bstr_SAD12_2.2_genomic.fna
<i>Brachypodium distachyon</i>	Angiosperm (Monocots)	NCBI	GCF_000005505.3_Brachypodium_distachyon_v3.0_genomic.fna
<i>Brassica carinata</i>	Angiosperm (Eudicots)	NCBI	GCA_016771965.1_ASM1677196v1_genomic.fna
<i>Brassica juncea</i>	Angiosperm (Eudicots)	NCBI	GCA_018703725.1_ASM1870372v1_genomic.fna
<i>Brassica napus</i>	Angiosperm (Eudicots)	NCBI	GCF_020379485.1_Da-Ac_genomic.fna
<i>Brassica nigra</i>	Angiosperm (Eudicots)	NCBI	GCA_016432835.1_Bnig_sang_1.1_genomic.fna
<i>Brassica oleracea</i>	Angiosperm (Eudicots)	NCBI	GCF_000695525.1_BOL_genomic.fna
<i>Brassica rapa</i>	Angiosperm (Eudicots)	NCBI	GCF_000309985.2_CAAS_Brap_v3.01_genomic.fna
<i>Camelina sativa</i>	Angiosperm (Eudicots)	NCBI	GCF_000633955.1_Cs_genomic.fna
<i>Camellia sinensis</i>	Angiosperm (Eudicots)	NCBI	GCF_004153795.1_AHAU_CSS_1_genomic.fna
<i>Capsicum annuum</i>	Angiosperm (Eudicots)	NCBI	GCF_002878395.1_UCD10Xv1.1_genomic.fna
<i>Carica papaya</i>	Angiosperm (Eudicots)	NCBI	GCF_000150535.2_Papaya1.0_genomic.fna
<i>Ceratodon purpureus</i>	Moss	NCBI	GCA_014871385.1_CpurpureusR40_1_0_genomic.fna
<i>Ceratopteris richardii</i>	Fern	NCBI	GCA_020310875.1_C.richardii_v2_genomic.fna
<i>Chara braunii</i>	Algae (Charales)	NCBI	GCA_003427395.1_Cbr_1.0_genomic.fna
<i>Chlamydomonas reinhardtii</i>	Algae (Chlorophyta)	NCBI	GCF_000002595.2_Chlamydomonas_reinhardtii_v5.5_genomic.fna
<i>Cinnamomum micranthum</i>	Angiosperm (Magnolids)	NCBI	GCA_003546025.1_ASBR_Ckan_1.0_genomic.fna
<i>Citrullus lanatus</i>	Angiosperm (Eudicots)	NCBI	GCA_029034555.1_KOR_cv.242-1_genomic.fna
<i>Citrus sinensis</i>	Angiosperm (Eudicots)	NCBI	GCF_022201045.2_DVS_A1.0_genomic.fna
<i>Cucumis melo</i>	Angiosperm (Eudicots)	NCBI	GCF_025177605.1_USDA_Cmelo_AY_1.0_genomic.fna
<i>Cucumis sativus</i>	Angiosperm (Eudicots)	NCBI	GCF_000004075.3_Cucumber_9930_V3_genomic.fna
<i>Cycas panzhihuaensis</i>	Gymnosperm	NCBI	GCA_023213395.1_ASM2321339v1_genomic.fna
<i>Cydonia oblonga</i>	Angiosperm (Eudicots)	NCBI	GCA_015708375.1_ASM1570837v1_genomic.fna
<i>Cynara cardunculus</i>	Angiosperm (Eudicots)	NCBI	GCF_001531365.2_CerdV1.1_genomic.fna
<i>Cyperus esculentus</i>	Angiosperm (Monocots)	NCBI	GCA_030179965.1_Ces_v0.1_genomic.fna
<i>Datura stramonium</i>	Angiosperm (Eudicots)	CNCB	GWHBCKB000000000_Datura.genome.fna
<i>Daucus carota</i>	Angiosperm (Eudicots)	NCBI	GCF_001625215.1_ASM162521v1_genomic.fna
<i>Dendrobium officinale</i>	Angiosperm (Monocots)	NCBI	GCA_019514585.1_ASM1951458v1_genomic.fna
<i>Fragaria ananassa</i>	Angiosperm (Eudicots)	NCBI	GCA_019022445.1_NIHHS_FaW_1.0_genomic.fna
<i>Fragaria vesca</i>	Angiosperm (Eudicots)	NCBI	GCF_000184155.1_FraVesHawaii_1.0_genomic.fna
<i>Ginkgo biloba</i>	Gymnosperm	NCBI	GCA_024626585.1_ASM2462658v1_genomic.fna
<i>Glycine max</i>	Angiosperm (Eudicots)	NCBI	GCF_000004515.6_Glycine_max_v4.0_genomic.fna
<i>Glycine soja</i>	Angiosperm (Eudicots)	NCBI	GCF_004193775.1_ASM419377v2_genomic.fna
<i>Gnetum montanum</i>	Gymnosperm	NCBI	GCA_015680685.1_Gmon01_genomic.fna
<i>Gossypium arboreum</i>	Angiosperm (Eudicots)	NCBI	GCF_025698485.1_ASM2569848v2_genomic.fna
<i>Gossypium hirsutum</i>	Angiosperm (Eudicots)	NCBI	GCF_007990345.1_Gossypium_hirsutum_v2.1_genomic.fna
<i>Gossypium raimondii</i>	Angiosperm (Eudicots)	NCBI	GCF_025698545.1_ASM2569854v1_genomic.fna
<i>Ipomoea batatas</i>	Angiosperm (Eudicots)	NCBI	GCA_002525835.2_ipoBat4_genomic.fna
<i>Isoetes engelmannii</i>	Lycopods	NCBI	GCA_011763485.2_ASM1176348v2_genomic.fna
<i>Isoetes taiwanensis</i>	Lycopods	NCBI	GCA_021234155.1_ASM2123415v1_genomic.fna
<i>Isoetecium myosuroides</i>	Moss	NCBI	GCA_951799445.1_cbIsoMyos1.1_genomic.fna
<i>Lactuca sativa</i>	Angiosperm (Eudicots)	NCBI	GCF_002870075.4_Lsat_Salinas_v11_genomic.fna
<i>Lemma minuta</i>	Angiosperm (Monocots)	NCBI	GCA_024174645.1_Salk_lm5633_a03_genomic.fna
<i>Lycium ferocissimum</i>	Angiosperm (Eudicots)	NCBI	GCA_029784015.1_AGI_CSIRO_Lferr_CH_V1_genomic.fna
<i>Lycopodium clavatum</i>	Lycopods	CNCB	GWHBJY000000000_Lycopodium.genome.fna
<i>Malus domestica</i>	Angiosperm (Eudicots)	NCBI	GCF_002114115.1_ASM211411v1_genomic.fna
<i>Mangifera indica</i>	Angiosperm (Eudicots)	NCBI	GCF_011075055.1_CATAS_Mindica_2.1_genomic.fna
<i>Marchantia paleacea</i>	Liverworts	NCBI	GCA_014180765.2_ASM1418076v2_genomic.fna
<i>Marchantia polymorpha</i>	Liverworts	NCBI	GCA_003032435.1_Marchanta_polymorpha_v1_genomic.fna
<i>Medicago truncatula</i>	Angiosperm (Eudicots)	NCBI	GCF_003473485.1_MtrunA17r5.0-ANR_genomic.fna
<i>Musa acuminata</i>	Angiosperm (Monocots)	NCBI	GCF_000313855.2_ASM31385v2_genomic.fna
<i>Nelumbo nucifera</i>	Angiosperm (Eudicots)	NCBI	GCF_000365185.1_Chinese_Lotus_1.1_genomic.fna
<i>Nicotiana benthamiana</i>	Angiosperm (Eudicots)	owned	Niben261_genome_complete01.fna
<i>Nymphaea colorata</i>	Angiosperm (Basal)	NCBI	GCF_008831285.2_ASM883128v2_genomic.fna
<i>Olea europaea</i>	Angiosperm (Eudicots)	NCBI	GCA_902713445.1_OLEA9_genomic.fna
<i>Oryza sativa</i>	Angiosperm (Monocots)	NCBI	GCF_001433935.1_IRGSP-1.0_genomic.fna.gz
<i>Phalaenopsis equestris</i>	Angiosperm (Monocots)	NCBI	GCF_001263595.1_ASM126359v1_genomic.fna
<i>Phaseolus vulgaris</i>	Angiosperm (Eudicots)	NCBI	GCF_000499845.1_PhaVulg1_0_genomic.fna
<i>Phoenix dactylifera</i>	Angiosperm (Monocots)	NCBI	GCF_009389715.1_palm_55x_up_171113_Pbpolish2nd_filt_p_genomic.fna
<i>Physalis pubescens</i>	Angiosperm (Eudicots)	CNCB	GWHANUX000000000_Physalis.genome.fna
<i>Physcomitrium patens</i>	Moss	NCBI	GCF_000002425.4_Phypa_V3_genomic.fna
<i>Pistacia vera</i>	Angiosperm (Eudicots)	NCBI	GCF_008641045.1_PisVer_v2_genomic.fna
<i>Populus trichocarpa</i>	Angiosperm (Eudicots)	NCBI	GCF_000002775.5_P.trichocarpa_v4.1_genomic.fna
<i>Prunus persica</i>	Angiosperm (Eudicots)	NCBI	GCF_000346465.2_Prunus_persica_NCBIV2_genomic.fna
<i>Punica granatum</i>	Angiosperm (Eudicots)	NCBI	GCF_007655135.1_ASM765513v2_genomic.fna
<i>Salvia hispanica</i>	Angiosperm (Eudicots)	NCBI	GCF_023119035.1_UniMelb_Shisp_WGS_1.0_genomic.fna
<i>Salvia splendens</i>	Angiosperm (Eudicots)	NCBI	GCF_004379255.2_SspV2_genomic.fna
<i>Selaginella moellendorffii</i>	Lycopods	NCBI	GCF_000143415.4_v1.0_genomic.fna
<i>Setaria italica</i>	Angiosperm (Monocots)	NCBI	GCF_000263155.2_Setaria_italica_v2.0_genomic.fna
<i>Solanum lycopersicum</i>	Angiosperm (Eudicots)	NCBI	GCF_000188115.5_SL3.1_genomic.fna
<i>Solanum melongena</i>	Angiosperm (Eudicots)	SGN	Eggplant_SGN_V4.1.fna
<i>Solanum pimpinellifolium</i>	Angiosperm (Eudicots)	NCBI	GCA_014964335.1_ASM1496433v1_genomic.fna
<i>Solanum stenotomum</i>	Angiosperm (Eudicots)	NCBI	GCF_019186545.1_ASM1918654v1_genomic.fna
<i>Solanum tuberosum</i>	Angiosperm (Eudicots)	NCBI	GCF_000226075.1_SolTub_3.0_genomic.fna
<i>Sphagnum fallax</i>	Moss	NCBI	GCA_021442195.1_S.fallax_v1.1_genomic.fna
<i>Sphagnum magellanicum</i>	Moss	NCBI	GCA_021904315.1_S.magellanicum_v1.1_genomic.fna
<i>Thuja plicata</i>	Gymnosperm	NCBI	GCA_018584345.1_redcedar-v3_genomic.fna
<i>Utricularia gibba</i>	Angiosperm (Eudicots)	NCBI	GCA_002189035.1_U_gibba_v2_genomic.fna
<i>Vanilla planifolia</i>	Angiosperm (Monocots)	NCBI	GCA_023846275.1_ASM2384627v1_genomic.fna
<i>Vicia sativa</i>	Angiosperm (Eudicots)	NCBI	GCA_021764765.1_ASM2176476v1_genomic.fna
<i>Vitis vinifera</i>	Angiosperm (Eudicots)	NCBI	GCF_000003745.3_12X_genomic.fna
<i>Welwitschia mirabilis</i>	Gymnosperm	CNCB	CNA0022760_Wmirabilis_genomic.fna
<i>Zea mays</i>	Angiosperm (Monocots)	NCBI	GCF_902167145.1_Zm-B73-REFERENCE-NAM-5.0_genomic.fna

2. Genome Metrics Retrieval

To retrieve the genome metrics two programs were used.

2.1 Quast v5.10: Assessing Genome Quality

Quast, the Quality Assessment Tool for Genome Assemblies, was significant in providing a comprehensive evaluation of the quality and completeness of the assembled genomes. Quast v5.10 (Mikheenko et al., 2018) was employed to extract essential genome metrics, including Assembly Size and Number of Contigs.

```
# Activate the conda environment with Quast installed
conda activate quast_env

# Run Quast
quast -o <out_directory_name> -m 1 -g <genomic_gff> <assembly_fasta>
```

2.2 agat_sp_statistics.pl: Detailed Genome Metrics

Additionally, agat_sp_statistics.pl v1.2.0 (Dainat et al., 2023) was employed for more detailed genome metrics retrieval, including providing the Number of Genes, contributing to a comprehensive understanding of the assembled genomes.

```
# Run agat_sp_statistics.pl
agat_sp_statistics.pl -gff <genomic_gff> -o statistics.txt
```

This dual-program approach ensured a thorough assessment of the genomic landscape, laying the groundwork for subsequent analyses.

3. Repetitive Element Identification

The intricate world of repetitive elements was unraveled through a two-fold process utilizing advanced tools.

3.1 RepeatModeler2: Unveiling the Unknown Repetitive Landscape

RepeatModeler2 v2.0.4 (Flynn et al., 2020) played a significant role in identifying unknown repetitive elements within the genomic sequences. This tool employed a de novo approach, discovering repetitive motifs without relying on pre-existing libraries. RepeatModeler's primary goal was to comprehensively analyze the entire genome, automating the identification, modeling, and classification of various types of repetitive elements. Its output was a library containing identified repetitive element families. Of particular significance was its capability to unearth Long Terminal Repeats (LTRs), shedding light on crucial genomic regions characterized by repetitive dynamics.

```
# Start a Docker instance for DFAM tools
screen -L -Logfile RunDFAMToolsDATE.log /data/Software/dfam-tetools.sh

# Generate the RepeatModeler sequence database
BuildDatabase --<name> <my_genome> <my_genome.fasta>

# Run RepeatModeler2
RepeatModeler -database <my_genome> -threads 4 -LTRStruct
```

3.2 RepeatMasker: Annotating the Repetitive Landscape

Building upon the insights gained from RepeatModeler2, RepeatMasker v4.1.5 (A.F.A. Smit et al., RepeatMasker at <http://repeatmasker.org>) enriched the annotation process, unraveling the repetitive landscape encoded within the genome. This tool annotated the identified repetitive elements in the genome FASTA sequence using the DFAM database and the elements found by RepeatModeler2. RepeatMasker's focus was to mark the positions of repetitive elements in the genome sequence using a known repetitive element library (often generated by RepeatModeler).

```
# Run RepeatMasker
RepeatMasker -lib <my_genome-families.fa> -xsmall -gff -pa 4
my_genome.fasta

# Exit from the Docker instance
exit;
```

RepeatMasker not only annotated repetitive elements but also provided crucial information on their genomic locations and possible families. This tandem approach facilitated a comprehensive exploration of the repetitive elements' distribution and composition within the genomes.

3.3 TESorter: Sorting and Categorizing Repetitive Elements

Following RepeatMasker, TESorter v1.4.6 (Zhang et al., 2022) was employed to sort individual repeats, providing a refined and categorized view of the repetitive landscape.

```
# Extract the repetitive elements previously annotated by RepeatMasker
python3 /data/users/collaborators/zhuy/miniconda3/bin/RepeatMasker.py
out2seqs <my_genome.fasta.out_from_RepeatMasker> <my_genome.fasta> >
<my_genome.IndividualRepeats.fasta>

# Run TESorter
TESorter <my_genome.IndividualRepeats.fasta> -p 48
```

4. Comparative Analysis of Transposable Elements

The comprehensive re-annotation of TE landscapes extended beyond identification to estimating diversity parameters and specificity. Diversity parameters, including Shannon's entropy (Hj) and Kolmogorov complexity (KC), were calculated, gauging variations within TE landscapes. Specificity, assessed through the specialization index (δ_j index) and divergence with respect to the entire TE landscape using Kullback–Leibler divergence (Divj), offered insights into TE behavior.

```
# Estimation of Diversity Parameters
python
~/trabajo/diversity_tools/diversity_tools/shannon_index_operations.py
-i <GenomesAnalysis_count_matrix.csv> -o diversity -f <olivia_results>

# Estimation of Specificity Parameters
python
~/trabajo/diversity_tools/diversity_tools/shannon_index_operations.py
-i <GenomesAnalysis_count_matrix.csv> -o specificity -f
<olivia_results>
```

An R script was developed to create horizontal stacked bar charts depicting the diversity and specificity patterns across different plant species. The ggplot2 package v3.4.4 (Wickham et al., 2022) was utilized for its versatile features and customization options.

```
library(ggplot2)
# Create a Horizontal Stacked Bar Chart of Diversity
horizontal_stacked_bar_chart_diversity <-
ggplot(Diversity_Specificity_of_species, aes(x = interaction(SPECIES,
LINEAGE), y = DIVERSITY, fill = LINEAGE)) +
  geom_bar(stat = "identity", position = "stack") +
  labs(title = "Horizontal Stacked Bar Chart of Diversity",
       x = "Species",
       y = "Diversity",
       fill = "Lineage") +
  theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 4),
        plot.title = element_text(hjust = 0.5),
        legend.position = "bottom") +
  scale_fill_viridis_d()
theme_set(theme_classic())
print(horizontal_stacked_bar_chart_diversity)
# Create a Horizontal Stacked Bar Chart of Specificity
horizontal_stacked_bar_chart <-
ggplot(Diversity_Specificity_of_species, aes(x = interaction(SPECIES,
LINEAGE), y = SPECIFICITY, fill = LINEAGE)) +
```

```

geom_bar(stat = "identity", position = "stack") +
labs(title = "Horizontal Stacked Bar Chart of Specificity",
      x = "Species",
      y = "Specificity",
      fill = "Lineage") +
theme(axis.text.x = element_text(angle = 45, hjust = 1, size = 4),
      plot.title = element_text(hjust = 0.5),
      legend.position = "bottom") +
scale_fill_viridis_d()
theme_set(theme_classic())
print(horizontal_stacked_bar_chart_specificity)

```

Subsequently, scatter plots were generated in R to visually depict the diversity and specificity patterns across various TE superfamilies.

```

# Sort the data frame by diversity in ascending order when creating it
Diversity_Specificity_of_TEs <-
Diversity_Specificity_of_TEs[order(Diversity_Specificity_of_TEs$DIVER
SITY), ]
# Creat Scatter Plot of Diversity
scatter_plot <- ggplot(Diversity_Specificity_of_TEs, aes(x =
SUPERFAMILY, y = DIVERSITY, color = SUPERFAMILY)) +
geom_point() +
ggtitle("Scatter Plot of Diversity") +
labs(x = "Superfamily", y = "Diversity", color = "Superfamily") +
theme_classic() +
theme(axis.text.x = element_text(angle = 45, hjust = 1),
      legend.position = "top",
      plot.title = element_text(hjust = 0.5))
print(scatter_plot_divercity)
# Creat Scatter Plot of Specificity
scatter_plot_specificity <- ggplot(Diversity_Specificity_of_TEs,
aes(x = SUPERFAMILY, y = SPECIFICITY, color = SUPERFAMILY)) +
geom_point() +
ggtitle("Scatter Plot of Specificity") +
labs(x = "Superfamily", y = "Specificity", color = "Superfamily") +
theme_classic() +
theme(axis.text.x = element_text(angle = 45, hjust = 1),
      legend.position = "top",
      plot.title = element_text(hjust = 0.5))
print(scatter_plot_specificity)

```

To explore the correlation between Genome Size and Bases Masked in plant genomes, the following R script was utilized to generate a comprehensive scatter plot. This scatter

plot not only visually represents the distribution of data points but also includes an accompanying trend line. The trend line serves as a valuable tool, allowing for the identification of potential patterns or trends within the relationship between genome size and the extent of masked bases. This integrated visualization provides a nuanced understanding of the interplay between these two genomic parameters.

```
# Create scatter plot with trend line
ggplot(GenomeSize, aes(x = `GENOME SIZE (Mb)`, y = `BASES MASKED (Mb)`,
color = LINEAGE)) +
  geom_point() +
  geom_smooth(method = "lm", se = FALSE, linetype = "solid", color =
"grey") +
  labs(title = "The Correlation between Genome Size and Bases Masked",
x = "Genome Size (Mb)",
y = "Bases Masked (Mb)",
color = "Lineage") +
  theme_classic() +
  theme(plot.title = element_text(hjust = 0.5))
```

The analysis also utilized the following programs:

<https://github.com/AgustinAmata/Repeattools>

Python programs were employed to analyze TESorter results, generating visualizations such as box plots, heatmaps, and Principal Component Analysis (PCA) plots. These visualizations provided a holistic understanding of the dynamics of TEs within plant genomes.

This multi-tiered approach positioned the study to make substantial contributions to unraveling the complex interplay between TEs and plant genome evolution. The integration of advanced tools and custom analyses provided a nuanced view of the repetitive element landscape, offering insights into their roles in genomic evolution

Results

1. Repetitive Element Identification

To comprehensively understand the repetitive elements (RE) landscape within the plant genomes, a series of analyses were conducted. Initially, the program Quast v5.10 was employed to obtain the assembly size (genome size) and the number of contigs (number of scaffolds) for Chlorophytes. Subsequently, the program agat_sp_statistics.pl was employed to obtain the number of genes for Chlorophytes. Finally, RepeatModeler2 and RepeatMasker were executed, providing detailed metrics such as the bases masked and their corresponding percentages in the genomes, the summarized results are provided below:

SPECIES	LINEAGE	N_GENES	GENOME SIZE (Mb)	N_SCAFFOLDS	BASES MASKED (Mb)	PERCENTAGE
<i>Astrephomene gubernaculifera</i>	Chlorophyte	13724	103.9	207	29.5	28.44%
<i>Auxenochlorella protothecoides (0710)</i>	Chlorophyte	13724	103.9	374	29.6	28.52%
<i>Auxenochlorella protothecoides (UTEX 25)</i>	Chlorophyte	5852	21.2	217	0.7	3.29%
<i>Bathycoccus prasinos</i>	Chlorophyte	7941	15.1	21	1.0	6.90%
<i>Chlamydomonas eustigma</i>	Chlorophyte	14068	66.6	520	3.9	5.89%
<i>Chlamydomonas incerta</i>	Chlorophyte	16350	129.2	453	34.9	27.03%
<i>Chlamydomonas reinhardtii</i>	Chlorophyte	17742	111.1	53	23.2	20.91%
<i>Chlamydomonas schloesseri</i>	Chlorophyte	15571	130.2	457	39.4	30.27%
<i>Chlamydomonas sp. UWO 241</i>	Chlorophyte	16018	211.6	2458	109.4	51.71%
<i>Chlorella desiccata (nom. nud.) (UTEX 2437)</i>	Chlorophyte	8418	20.7	29	1.7	8.01%
<i>Chlorella desiccata (nom. nud.) (UTEX 2526)</i>	Chlorophyte	9308	21.6	18	1.7	8.04%
<i>Chlorella ohadii</i>	Chlorophyte	10866	57.1	486	3.5	6.13%
<i>Chlorella sorokiniana</i>	Chlorophyte	9526	59.6	159	6.5	10.90%
<i>Chlorella variabilis</i>	Chlorophyte	9780	46.2	414	5.2	11.27%
<i>Chloropicon primus (CCMP 1205)</i>	Chlorophyte	8639	17.4	20	1.4	8.24%
<i>Chloropicon primus (RCC138)</i>	Chlorophyte	8627	17.6	20	1.6	8.85%
<i>Coccomyxa sp. Obi</i>	Chlorophyte	11588	50.4	21	2.2	4.35%
<i>Coccomyxa subellipsoidea C-169</i>	Chlorophyte	9834	48.8	29	2.6	5.31%
<i>Dunaliella salina</i>	Chlorophyte	16527	343.7	5512	125.7	36.58%
<i>Edaphochlamys debaryana</i>	Chlorophyte	19227	142.1	527	26.4	18.59%
<i>Haematococcus lacustris</i>	Chlorophyte	28279	309.4	9693	159.6	51.58%
<i>Helicosporidium sp. ATCC 50920</i>	Chlorophyte	6033	12.4	5666	0.4	2.92%
<i>Micractinium conductrix</i>	Chlorophyte	9217	61.0	300	12.1	19.77%
<i>Micromonas commoda</i>	Chlorophyte	10147	21.1	19	1.1	5.22%
<i>Micromonas pusilla CCMP1545</i>	Chlorophyte	10238	22.0	21	3.5	15.42%
<i>Monoraphidium minutum</i>	Chlorophyte	15335	68.2	511	13.2	19.33%
<i>Monoraphidium neglectum</i>	Chlorophyte	16755	69.7	6720	6.0	8.55%
<i>Ostreobium quekettii</i>	Chlorophyte	10833	151.9	3134	53.8	35.42%
<i>Ostreococcus lucimarinus</i>	Chlorophyte	7603	13.2	21	1.2	9.04%
<i>Ostreococcus tauri (RCC1115)</i>	Chlorophyte	8099	14.8	70	1.7	11.46%
<i>Ostreococcus tauri (RCC4221)</i>	Chlorophyte	7869	13.0	22	0.8	5.81%
<i>Pedinophyceae sp. YPF-701</i>	Chlorophyte	7940	27.9	32	2.3	8.33%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	7227	14.9	16	1.4	9.59%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	7227	14.9	12	1.4	9.59%
<i>Pycnococcus provasolii</i>	Chlorophyte	11297	22.7	43	1.5	6.59%
<i>Raphidocelis subcapitata</i>	Chlorophyte	13383	51.2	300	6.6	12.88%
<i>Scenedesmus sp. NREL 46B-D3</i>	Chlorophyte	17398	151.9	2661	44.0	28.98%
<i>Tetrahena socialis</i>	Chlorophyte	14296	135.8	5856	17.3	12.77%
<i>Trebouxia sp. AI-2</i>	Chlorophyte	13851	52.9	217	3.0	5.71%
<i>Volvox africanus</i>	Chlorophyte	13577	129.3	448	38.1	29.46%
<i>Volvox reticuliferus (NIES 3786)</i>	Chlorophyte	13772	134.0	230	39.8	29.72%
<i>Volvox reticuliferus (NIES-3785)</i>	Chlorophyte	13898	133.1	200	39.0	29.31%

Within the presented table, a comprehensive dataset related to Repetitive Element identification is provided. The table details key parameters for each genome, including the number of genes, genome size, number of scaffolds, bases masked, and their respective percentages. Based on the table:

- *Dunaliella salina* has the largest genome size (343.7 Mb).
- *Haematococcus lacustris* presents the highest number of genes, the highest number of scaffolds (9693), the highest number of bases masked (159.6 Mb), and the maximum percentage of bases masked (51.58%).
- *Helicosporidium sp. ATCC 50920* has the smallest genome size (12.4 Mb), the lowest number of bases masked (0.4 Mb), and the minimum percentage of bases masked (2.92%).
- *Auxenochlorella protothecoides (UTEX 25)* has the lowest number of genes (5852).

- *Picochlorum sp. BPE23* has the lowest number of scaffolds (12).

It is evident from the data that within the Chlorophyte lineage, there is a positive correlation between genome size and the relative content of Transposable Elements. For instance, *Haematococcus lacustris*, with a genome size of 309.4M, exhibits a TE content of 51.58%. Conversely, smaller genomes within this lineage tend to have lower TE content. For example, *Ostreococcus tauri* (RCC4221), having a genome size of 13.0M, has a TE content of 5.81%.

I continued the analyses using Repeat Modeler 2 and Repeat Masker until detailed data was obtained for all the studied plants, including the number of contigs, assembly size, bases masked, and their corresponding percentages. The results are presented below:

SPECIES	LINEAGE	N. CONTIGS	GENOME SIZE (Mb)	BASES MASKED (Mb)	PERCENTAGE
<i>Acer negundo</i>	Angiosperm (Eudicots)	108	442.4	272.2	61.51%
<i>Acer yuenghiense</i>	Angiosperm (Eudicots)	280	665.9	472.0	70.88%
<i>Actinidia chinensis</i>	Angiosperm (Eudicots)	1234	553.8	232.5	41.98%
<i>Adansonia digitata</i>	Angiosperm (Eudicots)	232407	687.0	382.1	55.61%
<i>Adiantum capillus</i>	Fern	601	4822.6	4118.5	85.40%
<i>Amborella trichopoda</i>	Angiosperm (Basal)	5746	706.5	436.0	61.71%
<i>Ananas comosus</i>	Angiosperm (Monocots)	3129	382.1	218.9	57.29%
<i>Annona glabra</i>	Angiosperm (Magnoliids)	66	1027.3	593.6	57.78%
<i>Antroceros angustus</i>	Hornwort	289	119.3	57.6	48.25%
<i>Arabidopsis arenosa</i>	Angiosperm (Eudicots)	8	149.7	32.9	21.98%
<i>Arabidopsis suecica</i>	Angiosperm (Eudicots)	269	272.3	67.6	24.83%
<i>Arabidopsis thaliana</i>	Angiosperm (Eudicots)	7	119.7	20.4	17.05%
<i>Arabis alpina</i>	Angiosperm (Eudicots)	8	311.6	166.1	53.30%
<i>Boechera stricta</i>	Angiosperm (Eudicots)	22	190.5	70.7	37.12%
<i>Brachypodium distachyon</i>	Angiosperm (Monocots)	11	271.3	101.0	37.22%
<i>Brassica carinata</i>	Angiosperm (Eudicots)	1636	1087.0	653.6	60.13%
<i>Brassica juncea</i>	Angiosperm (Eudicots)	1544	933.5	501.4	53.71%
<i>Brassica napus</i>	Angiosperm (Eudicots)	3167	1001.9	575.2	57.41%
<i>Brassica nigra</i>	Angiosperm (Eudicots)	1054	534.2	308.6	57.77%
<i>Brassica oleracea</i>	Angiosperm (Eudicots)	32886	489.0	236.1	48.30%
<i>Brassica rapa</i>	Angiosperm (Eudicots)	1100	353.0	170.4	48.27%
<i>Camelina sativa</i>	Angiosperm (Eudicots)	37212	641.4	255.2	39.78%
<i>Camellia sinensis</i>	Angiosperm (Eudicots)	14028	3105.4	2267.0	73.00%
<i>Capsicum annuum</i>	Angiosperm (Eudicots)	81202	3212.5	2612.7	81.33%
<i>Carya papaya</i>	Angiosperm (Eudicots)	17766	370.4	149.8	40.45%
<i>Ceratodon purpureus</i>	Moss	151	362.5	169.1	46.66%
<i>Ceratopteris richardii</i>	Fern	6172	7462.5	6463.4	86.61%
<i>Chara braunii</i>	Algae (Charales)	11654	1751.2	1019.7	58.23%
<i>Chlamydomonas reinhardtii</i>	Algae (Chlorophyta)	53	111.1	23.9	21.47%
<i>Cinnamomum micranthum</i>	Angiosperm (Magnoliids)	2150	730.4	403.7	55.27%
<i>Citrullus lanatus</i>	Angiosperm (Eudicots)	14	361.5	216.7	59.96%
<i>Citrus sinensis</i>	Angiosperm (Eudicots)	11	299.8	151.2	50.44%
<i>Cucumis melo</i>	Angiosperm (Eudicots)	1310	438.4	280.4	63.96%
<i>Cucumis sativus</i>	Angiosperm (Eudicots)	88	226.6	82.9	36.56%
<i>Cycas panzhuangensis</i>	Gymnosperm	923	10482.7	8796.7	83.92%
<i>Cydonia oblonga</i>	Angiosperm (Eudicots)	303932	488.4	268.7	55.01%
<i>Cynara cardunculus</i>	Angiosperm (Eudicots)	8174	725.0	457.6	63.12%
<i>Cyperus esculentus</i>	Angiosperm (Monocots)	1190	296.6	124.7	42.03%
<i>Datura stramonium</i>	Angiosperm (Eudicots)	13	1974.3	1678.6	85.02%
<i>Daucus carota</i>	Angiosperm (Eudicots)	4826	421.5	184.5	43.76%
<i>Dendrobium officinale</i>	Angiosperm (Monocots)	1621	1228.7	895.7	72.90%
<i>Fragaria ananassa</i>	Angiosperm (Eudicots)	204	805.7	359.3	44.60%
<i>Fragaria vesca</i>	Angiosperm (Eudicots)	3048	214.4	63.1	29.42%
<i>Ginkgo biloba</i>	Gymnosperm	265987	2638.1	998.4	37.84%
<i>Glycine max</i>	Angiosperm (Eudicots)	284	978.9	527.0	53.84%
<i>Glycine soja</i>	Angiosperm (Eudicots)	1120	1013.8	558.1	55.05%
<i>Gnetum montanum</i>	Gymnosperm	38363	2147.7	1770.0	82.41%
<i>Gossypium arboreum</i>	Angiosperm (Eudicots)	948	1621.4	1356.4	83.66%
<i>Gossypium hirsutum</i>	Angiosperm (Eudicots)	1027	2306.1	1792.7	77.74%
<i>Gossypium raimondii</i>	Angiosperm (Eudicots)	289	751.0	524.0	69.76%
<i>Iponoea batatas</i>	Angiosperm (Eudicots)	28461	837.0	392.0	46.83%
<i>Isoetes engelmannii</i>	Lycopods	319260	641.0	281.5	43.92%
<i>Isoetes taiwanensis</i>	Lycopods	1113	1658.3	857.9	51.73%
<i>Isoetium myosotoides</i>	Moss	256	388.5	222.7	57.32%
<i>Lactuca sativa</i>	Angiosperm (Eudicots)	93	2590.4	2210.7	85.34%
<i>Lemma minuta</i>	Angiosperm (Monocots)	2381	360.4	245.4	68.08%
<i>Lycium ferocissimum</i>	Angiosperm (Eudicots)	14905	1219.4	896.3	73.51%
<i>Lycopodium clavatum</i>	Lycopods	7102	2304.7	1853.8	80.43%
<i>Malus domestica</i>	Angiosperm (Eudicots)	807	703.4	391.0	55.59%
<i>Mangifera indica</i>	Angiosperm (Eudicots)	250	392.0	197.0	50.26%
<i>Marchantia paleacea</i>	Liverworts	192	250.8	112.2	44.75%
<i>Marchantia polymorpha</i>	Liverworts	2957	225.8	55.8	24.70%
<i>Medicago truncatula</i>	Angiosperm (Eudicots)	42	430.0	221.3	51.46%
<i>Musa acuminata</i>	Angiosperm (Monocots)	7259	472.2	173.0	36.63%
<i>Nelumbo nucifera</i>	Angiosperm (Eudicots)	3603	804.6	424.4	52.74%
<i>Nicotiana benthamiana</i>	Angiosperm (Eudicots)	17640	3035.8	2192.4	72.22%
<i>Nymphaea colorata</i>	Angiosperm (Basal)	786	408.9	177.2	43.33%
<i>Olea europaea</i>	Angiosperm (Eudicots)	9753	1316.7	821.6	62.40%
<i>Oryza sativa</i>	Angiosperm (Monocots)	58	374.4	181.9	48.58%
<i>Phalaenopsis equestris</i>	Angiosperm (Monocots)	89584	1064.2	704.9	66.24%
<i>Phaseolus vulgaris</i>	Angiosperm (Eudicots)	708	521.1	271.1	52.02%
<i>Phoenix dactylifera</i>	Angiosperm (Monocots)	2391	773.2	447.9	57.92%
<i>Physalis pubescens</i>	Angiosperm (Eudicots)	327	1389.3	1161.7	83.61%
<i>Physcomitrium patens</i>	Moss	359	472.1	283.1	59.98%
<i>Pistacia vera</i>	Angiosperm (Eudicots)	1865	671.3	428.9	63.89%
<i>Populus trichocarpa</i>	Angiosperm (Eudicots)	47	392.3	182.0	46.38%
<i>Prunus persica</i>	Angiosperm (Eudicots)	192	227.6	96.1	42.24%
<i>Punica granatum</i>	Angiosperm (Eudicots)	474	320.5	160.5	50.08%
<i>Salvia hispanica</i>	Angiosperm (Eudicots)	1556	321.5	141.8	44.12%
<i>Salvia splendens</i>	Angiosperm (Eudicots)	1162	806.1	475.6	59.00%
<i>Selaginella moellendorffii</i>	Lycopods	757	212.3	108.1	50.90%
<i>Setaria italica</i>	Angiosperm (Monocots)	337	405.9	206.1	50.79%
<i>Solanum lycopersicum</i>	Angiosperm (Eudicots)	3066	838.0	517.2	62.47%
<i>Solanum melongena</i>	Angiosperm (Eudicots)	13	1164.4	807.5	69.34%
<i>Solanum pimpinellifolium</i>	Angiosperm (Eudicots)	127	808.1	568.4	70.34%
<i>Solanum stenotomum</i>	Angiosperm (Eudicots)	17084	846.4	503.1	59.44%
<i>Solanum tuberosum</i>	Angiosperm (Eudicots)	14853	705.9	433.3	61.38%
<i>Sphagnum fallax</i>	Moss	36	395.1	185.5	46.94%
<i>Sphagnum magellanicum</i>	Moss	31	439.0	217.5	49.55%
<i>Thuja plicata</i>	Gymnosperm	67899	9095.9	6448.2	70.89%
<i>Utricularia gibba</i>	Angiosperm (Eudicots)	518	100.7	29.7	29.52%
<i>Vanilla planifolia</i>	Angiosperm (Monocots)	3874	1416.4	1072.3	75.71%
<i>Vicia sativa</i>	Angiosperm (Eudicots)	18	1653.6	1346.0	81.40%
<i>Vitis vinifera</i>	Angiosperm (Eudicots)	1907	486.2	264.6	54.43%
<i>Welwitschia mirabilis</i>	Gymnosperm	22	6867.0	5231.1	76.18%
<i>Zea mays</i>	Angiosperm (Monocots)	687	2182.8	1823.4	83.54%

Based on the table:

- *Isoetes engelmannii* (Lycopods) has the highest number of contigs (319,260).
- *Cycas panzhihuaensis* (Gymnosperm) has the largest assembly size (10,482.7 Mb) and the highest number of bases masked (8,796.7 Mb).
- *Ceratopteris richardii* (Fern) has the maximum percentage of bases masked (86.61%).
- *Arabidopsis thaliana* (Angiosperm) has the lowest number of contigs (7) and the minimum percentage of bases masked (17.05%).
- *Chlamydomonas reinhardtii* (Algae) has the smallest assembly size (111.1 Mb).
- *Utricularia gibba* (Angiosperm) has the lowest number of bases masked (29.7 Mb).

In this table, plant lineages such as Angiosperms, Ferns, Hornwort, Moss, Gymnosperms, and Lycopods are represented.

It is evident that similar to Chlorophyte, there exists a positive correlation between genome size and TE content across these lineages.

Gymnosperms typically have large genomes with elevated TE content. For instance, *Cycas panzhihuaensis* has an assembly size of 10,482.7 Mb with a TE content of 83.92%. Similarly, *Gnetum montanum* possesses an assembly size of 2,147.7 Mb and a TE content of 82.41%. It's noteworthy to mention an outlier, *Ginkgo biloba*, with a genome size of 2638.1 Mb, yet a TE content of only 37.84%.

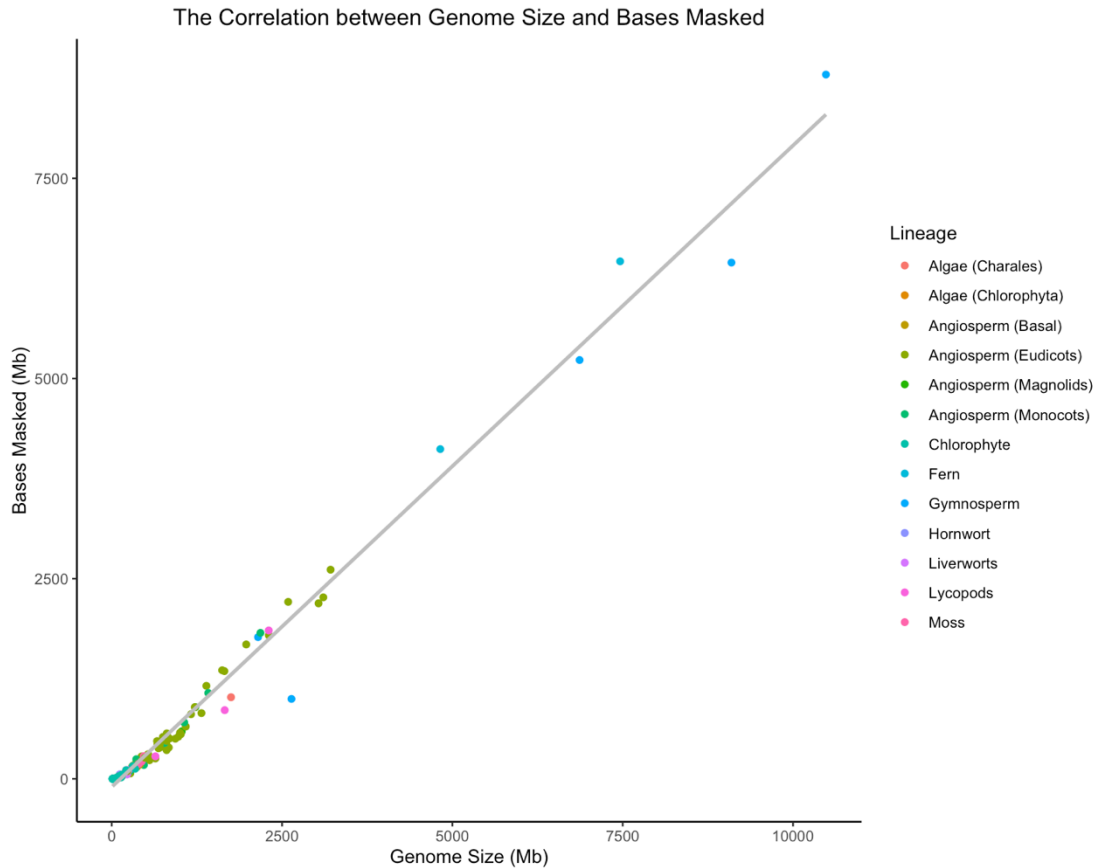
Ferns also tend to have expansive genomes and high TE content. *Ceratopteris richardii*, with an assembly size of 7,462.5 Mb, has a TE content of 86.61%. In comparison, *Adiantum capillus* has an assembly size of 4,822.6 Mb and a TE content of 86.61%.

For Angiosperms, the genome sizes vary considerably, as do the TE contents. *Capsicum annuum* has an assembly size of 3,212.5 Mb with a TE content of 81.33%, while *Nicotiana benthamiana*, with an assembly size of 3,035.8 Mb, has a TE content of 72.22%. On the other end of the spectrum, *Arabidopsis thaliana*, with an assembly size of 119.7 Mb, has a TE content of 17.05%, while *Utricularia gibba*, with an assembly size of 100.7 Mb, has a TE content of 29.52%.

In contrast to the aforementioned lineages, Mosses generally exhibit smaller genomes and lower TE content. For example, *Ceratodon purpureus* has an assembly size of 362.5 Mb and a TE content of 46.66%, while *Sphagnum fallax* possesses an assembly size of 395.1 Mb with a TE content of 46.94%.

Visualizing the Correlation between Genome Size and Bases Masked

To explore the correlation between Genome Size and Bases Masked and validate their relationship, a scatter plot was generated for effective visualization of the results.



From this scatter plot, it is evident that Gymnosperms exhibit the largest Genome Size and the highest content of Bases Masked, while Chlorophytes have the smallest Genome Size and the lowest Bases Masked content. However, regardless of lineage, there's a clear positive correlation between Genome Size and TE content. The presence of a grey reference line further emphasizes this trend, with data points clustered around the line, highlighting the consistent relationship between Genome Size and Bases Masked across all lineages.

2. Comparative and analysis

2.1 Distribution Metrics and Box Plots Representations

The analysis provides insight into the distribution and characteristics of various repetitive element classes in the genome. Within our dataset, we have captured fundamental metrics for key repetitive elements, such as LTRs and DNAs, including their quantity, length, and respective proportions.

To visualize these data comprehensively, we employed box plots, showcasing the divergence percentages of specific repetitive elements across all analyzed plant genomes. Through these box plots, insights into the variability, central tendencies, and potential outliers within the dataset are revealed, offering a clear and concise representation of comparative genomics concerning repetitive elements.

From the table, it's evident that the majority of the Retroelements across the plant lineages are LTR elements.

In lineages such as Angiosperms, Ferns, Hornwort, Moss, Gymnosperms, and Lycopods, the plant species with the highest proportion of Retroelements are *Datura stramonium* (Angiosperm) and *Zea mays* (Angiosperm). For *Datura stramonium*, 65.25% of its genome consists of Retroelements, out of which 64.24% are LTR elements. This implies that a staggering 98.5% of its retrotransposons are LTRs. Similarly, *Zea mays* has 64.31% of its genome as Retroelements, with 63.34% being LTR elements, indicating that 98.5% of its retrotransposons are LTRs.

In contrast, *Arabidopsis thaliana* (Arabidopsis) has a much lower percentage, with Retroelements making up 7.16% of its genome and LTR elements accounting for 6.02%. This means that 84.1% of its retrotransposons are LTRs.

In the Fern and Gymnosperm lineages, the number, length, and percentage of retrotransposons and LTRs are generally high. Mosses, however, exhibit a relatively lower content of retrotransposons.

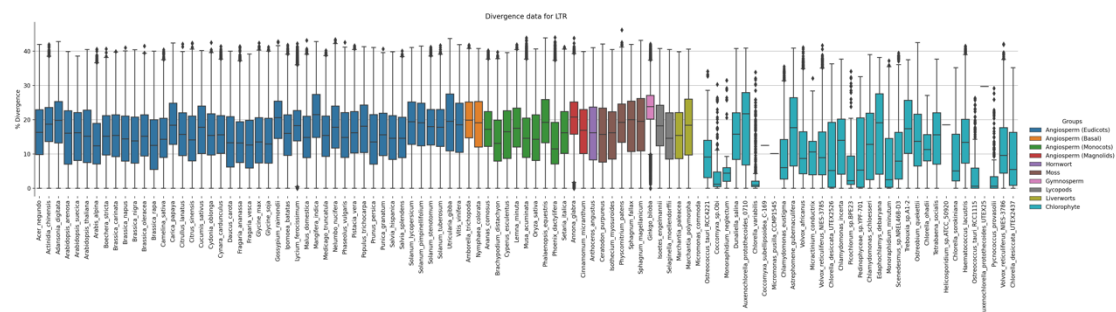
However, across all these lineages, it's notable that most plant genomes have LTR elements comprising over 80% of their Retroelements.

In contrast to the other lineages, Chlorophyte lineage genomes have notably fewer retrotransposons and LTRs. Within Chlorophyte, six genomes lack retrotransposons entirely, and approximately one-fourth lack LTR elements. In the remaining genomes, the majority have LTRs making up less than 3% of their genome.

However, there are exceptions within Chlorophyte. Notably, *Haematococcus lacustris* has retrotransposons comprising 17.53% of its genome and LTRs making up 11.35%. *Ostreobium quekettii* has retrotransposons at 16.94% and LTRs at 16.29% of its genome. A particularly unusual case is *Dunaliella salina*, where retrotransposons account for a significant 15.68% of its genome, while LTRs only constitute 2.29%.

2.1.1.1 LTR Repetitive Element Analysis

LTR element divergence data was analyzed and visually represented for each plant genome using boxplots.



This panel of boxplots displays the distribution of %Divergence values for LTR elements in all plant genomes. Each boxplot represents a single genome, allowing for a comparative assessment of divergence patterns across different species. The central line within each box denotes the median %Divergence, while the box's boundaries indicate the interquartile range (IQR). Whiskers extend to show the data's range, excluding any potential outliers.

Species-Specific Divergence: Genomes from certain plant species (e.g., *Ginkgo biloba* (Gymnosperm)), exhibited higher median %Divergence values, indicating potentially faster evolutionary rates in their LTR elements.

Distinct Low Median Divergence: Conversely, some plant genomes (e.g., *Ostreococcus tauri* RCC1115 (Chlorophyte)) presented notably lower median %Divergence values. This could imply a slower evolutionary rate or distinct evolutionary constraints acting upon the LTR elements in these species.

Uniform Divergence Patterns: A subset of genomes demonstrated narrower IQRs (e.g., *Ostreococcus tauri* RCC1115 (Chlorophyte)), suggesting a more uniform divergence pattern within those species.

Variable Divergence Patterns: In contrast to the aforementioned uniform patterns, certain genomes (e.g., *Auxenochlorella protothecoides* 0710 (Chlorophyte)) exhibited wider interquartile ranges (IQRs). This broader IQR suggests a greater variability in the divergence rates of LTR elements within these species. Such variability could arise from a myriad of factors, including differential TE activity, genomic structural variations, or varying evolutionary pressures across populations of these plants.

Outliers and Unique Evolutionary Trajectories: Beyond the central tendencies, certain genomes (e.g., *Ginkgo biloba* (Gymnosperm)) displayed outliers, signifying either exceptionally high or low %Divergence values. These outliers might be indicative of unique evolutionary trajectories, genomic events, or selective pressures influencing the evolution and divergence of LTR elements in those specific plants.

Comparative Evolutionary Dynamics: The juxtaposition of these divergence patterns across diverse plant genomes offers a glimpse into the varied evolutionary dynamics at play. The differential %Divergence values underscore the complexity of evolutionary processes and the intricate interplay between genomic structures and evolutionary pressures.

Overall, when examining the LTR elements across various plant lineages, a distinct pattern emerges. Within the Chlorophyte lineage, the median %Divergence ranges from 0 to 22, displaying significant diversity. A majority of these values predominantly lie below 15, underscoring the nuanced evolutionary patterns within the Chlorophyte lineage. The interquartile range (IQR) within this group exhibits variability, with some

species presenting narrower IQRs, while others showcase wider ones.

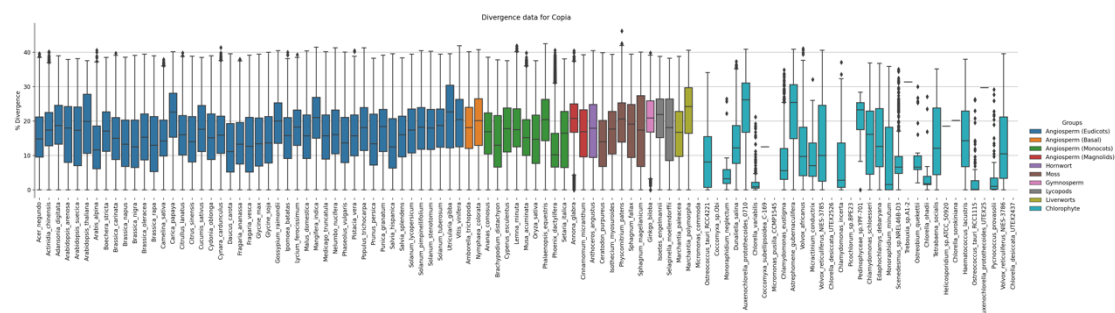
Conversely, for plant species outside the Chlorophyte lineage, the pattern is notably consistent. The median %Divergence for these species tends to be higher, clustering around 15 to 20. The IQR for these non-Chlorophyte species is relatively uniform, generally ranging between 10 and 15. Additionally, the whisker lengths for these species span from 0 to 40, indicating a broader range of variation in LTR elements across these lineages.

LTR Superfamily-Specific Divergence Analysis:

To delve deeper into the specific evolutionary patterns of LTR elements, we further segregated and analyzed the two predominant Superfamilies, Copia and Gypsy. These two superfamilies, renowned for their widespread distribution and significant impact on genome evolution, offer unique insights into the evolutionary trajectories and dynamics within the plant genomes.

Copia Repetitive Element Analysis:

Presented below are the boxplots representing the %Divergence of Copia elements across the diverse set of analyzed plant species.



For a majority of the plant lineages, spanning Angiosperms, Hornwort, Moss, Gymnosperm, Lycopods, and Liverworts, the median %Divergence of Copia elements predominantly falls within the range of 15-20. This suggests a relatively conserved evolutionary pattern across these lineages. However, notable exceptions arise in specific species such as *Utricularia gibba* (Angiosperm - Eudicots) and *Marchantia polymorpha* (Liverworts). These species exhibit a higher median %Divergence, approximately 25, indicating potentially distinct evolutionary pressures or histories for the Copia elements within these genomes.

The interquartile range (IQR) for %Divergence remains broadly consistent across most lineages. Nevertheless, outliers in some species, such as *Annona glabra* (Angiosperm - Magnoliids) and *Ginkgo biloba* (Gymnosperms), exhibit narrower IQRs. *Annona glabra* stands out further due to its particularly high median %Divergence and the presence of numerous outliers, suggesting unique evolutionary dynamics for Copia elements in this species.

Furthermore, while the whisker lengths across the analyzed lineages remain relatively

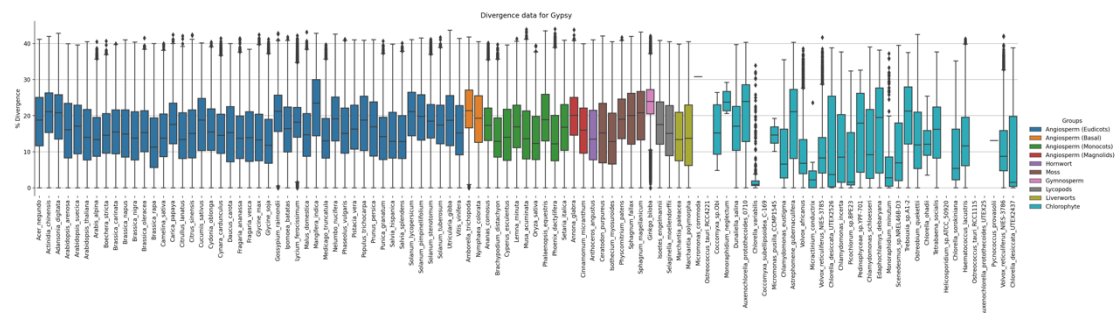
consistent, spanning from 0 to 40, exceptions like *Phoenix dactylifera* (Angiosperm - Monocots) showcase shorter whiskers. Intriguingly, this species also presents an abundance of outliers surpassing the maximum %Divergence, indicative of potential rapid or distinct evolutionary events for its Copia elements.

However, within the Chlorophyte lineage, distinct patterns emerge in the %Divergence characteristics of Copia elements. The boxplot distributions for Chlorophyte species manifest considerable variability, evidenced by a spectrum of median %Divergence values. Notably, species such as *Auxenochlorella protothecoides 0710* and *Astrephomene gubernaculifera* present elevated median %Divergence values, accompanied by extended interquartile ranges (IQRs) and whisker lengths ranging from 0 to around 40. Conversely, some species exhibit very low median values, such as *Ostreococcus tauri_RCC1115*, which exhibit a median %Divergence of zero. Moreover, the boxplots for *Ostreococcus tauri_RCC1115*, *Pycnococcus provasoli*, and *Chlorella variabilis* depict notably narrow IQRs and a pronounced presence of outliers exceeding the maximum %Divergence value.

In summary, while a majority of the analyzed plant lineages exhibit a relatively stable and conserved %Divergence for Copia elements, specific species, especially within the Chlorophyte lineage, present unique evolutionary patterns, emphasizing the diverse evolutionary dynamics of the Copia superfamily across plant genomes.

Gypsy Repetitive Element Analysis:

Presented below are the boxplots representing the %Divergence of Gypsy elements across the diverse set of analyzed plant species.



For a majority of the plant lineages, spanning Angiosperms, Hornwort, Moss, Gymnosperm, Lycopods, and Liverworts, the median %Divergence of Gypsy elements predominantly clusters within the range of 15-25. Noteworthy species in this context include *Mangifera indica* (Angiosperm - Ludicots) and *Ginkgo biloba* (Gymnosperms), both showcasing a median %Divergence around 25. Compared to the broader trend where most plant species exhibit an IQR of approximately 10, *Ginkgo biloba* stands out with a notably narrow IQR. Additionally, while the whisker lengths for the majority of species range from 0 to 40, *Ginkgo biloba*'s whiskers span approximately 12-37, accompanied by numerous outliers at both extremes.

Overall, analogous to Copia, these various lineages manifest a pattern of relatively close

and stable %Divergence for Gypsy elements.

However, within the Chlorophyte lineage, the boxplot distributions for Gypsy elements remain notably diverse. Drawing parallels to observations in the Copia distributions, species like *Auxenochlorella protothecoides* 0710 and *Astrephomene gubernaculifera* once again display elevated median %Divergence values, complemented by broad IQR and whiskers spanning from 0 to 40. Conversely, *Chlorella variabilis* continues to present a strikingly low median %Divergence. Its boxplot is characterized by a narrow IQR, shorter whiskers, and a significant representation of outliers exceeding the maximum %Divergence value.

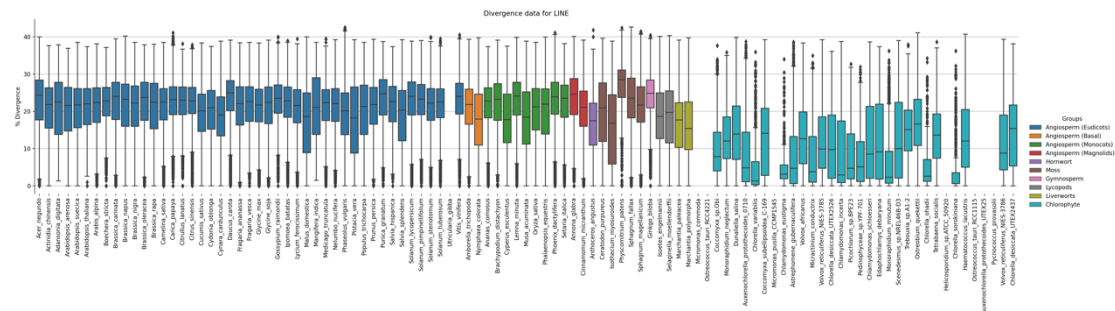
In summary, while the majority of plant lineages depict a relatively consistent %Divergence trend for Gypsy elements, specific species, particularly within the Chlorophyte lineage, highlight distinct evolutionary trajectories, underscoring the intricate evolutionary dynamics of the Gypsy superfamily across plant genomes.

2.1.1.2 Non-LTR Repetitive Element Analysis

Within the realm of Non-LTR transposons, LINEs (Long Interspersed Nuclear Elements) and SINEs (Short Interspersed Nuclear Elements) are recognized as prominent superfamily members. These elements have left indelible imprints on plant genomes due to their transposition dynamics.

LINE Repetitive Element Analysis

The following boxplots detail the %Divergence profiles of LINEs across diverse plant lineages.



For the majority of plant groups, including Angiosperms, Hornwort, Moss, Gymnosperms, Lycopods, and Liverworts, the median %Divergence typically falls within the 15 to 25 range. However, an exception is *Physcomitrium patens* (Moss), which exhibits a slightly elevated median value of around 28.

Regarding the Interquartile Range (IQR), most lineages maintain a consistent range, roughly around 10. Yet, species like *Physcomitrium patens* (Moss) and *Ginkgo biloba* (Gymnosperms) stand out with narrower IQRs, accompanied by numerous outliers falling below the minimum value.

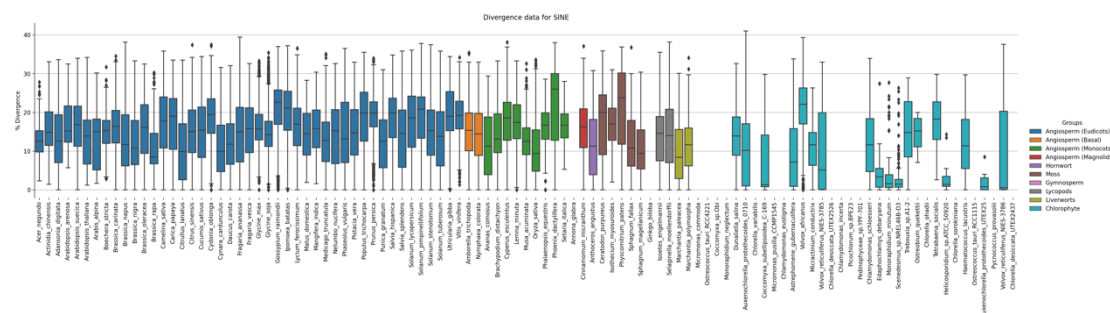
The whisker lengths vary among these lineages. While the maximum whisker value is

approximately 40 for many plants, the minimum values span between 0 and 15. Intriguingly, around half of the species across these lineages showcase outliers that fall below this minimum value.

Conversely, within the Chlorophyte lineage, there's a distinct profile for LINE elements. The median %Divergence values are notably lower, predominantly between 0 and 15. Notably, species like *Chlorella variabilis*, *Chlamydomonas eustigma*, *Monoraphidium minutum*, *Chlorella ohadii*, and *Chlorella sorokiniana* display narrow IQRs and median values below 5, accompanied by a significant number of outliers exceeding the maximum value. In contrast, the majority of other Chlorophyte plants depict wider IQRs, roughly around 15, with whisker lengths spanning from 0 to 40.

SINE Repetitive Element Analysis

The following boxplots detail the %Divergence profiles of SINEs across diverse plant lineages.



For SINE elements, the boxplot distributions present a more varied landscape compared to LINE elements. Across the Angiosperms, Hornwort, Moss, Lycopods, and Liverworts lineages, the median %Divergence for SINE elements ranges from 8 to 26. Notably, *Brassica rapa* (Angiosperm) and *Marchantia palearea* (Liverworts) exhibit the lowest median values, both at 8, while *Phoenix dactylifera* (Angiosperm) displays the highest median value at 26.

The IQR values among these lineages vary significantly. For instance, species like *Phoenix dactylifera* (Angiosperm) and *Physcomitrium patens* (Moss) have broader IQRs, ranging between 15 and 20. In contrast, species such as *Acer negundo* (Angiosperm) present narrower IQRs, approximately around 5.

Regarding whisker lengths, they vary across these lineages, with minimum values starting from 0 and maximum values up to 39. A minority of species within these lineages exhibit outliers on both ends of the distribution.

Within the Chlorophyte lineage, the SINE elements showcase a median %Divergence ranging from 2 to 23. Notably, *Volvox africanus* presents the highest median value, while *Volvox reticuliferus NIES-3786* exhibits the lowest. The IQR values for Chlorophyte species vary, with species like *Volvox reticuliferus NIES-3786* having an IQR exceeding 20. Conversely, species like *Scenedesmus sp. NREL46B-D3* and

Monoraphidium minutum have narrower IQRs, less than 5. Additionally, these two species have notably shorter whiskers, with a considerable number of outliers surpassing the maximum %Divergence value.

2.1.2 DNA Transposon Comparative and Analysis

DNA transposons are crucial agents in genomic evolution, the following table lists the number, length, and genome percentage of DNA transposons for each of the 138 plant genomes.

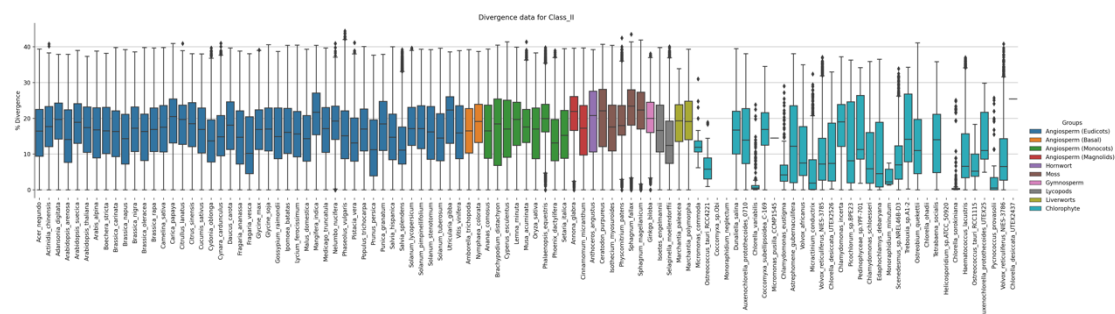
SPECIES	LINEAGE	DNA transposons (number of elements)	DNA transposons (length occupied)	DNA transposons (percentage of sequence)
<i>Acer negundo</i>	Angiosperm (Eudicots)	23076	2281943 bp	5.31%
<i>Acer sanguinonea</i>	Angiosperm (Eudicots)	34847	2214515 bp	3.33%
<i>Actinidia chinensis</i>	Angiosperm (Eudicots)	15649	5025946 bp	0.91%
<i>Adiantum alatum</i>	Angiosperm (Eudicots)	21277	9029957 bp	1.48%
<i>Adiantum capillus</i>	Fern	140142	137239168 bp	2.85%
<i>Ankorella trichopoda</i>	Angiosperm (Basal)	31356	14818905 bp	2.10%
<i>Annona comosa</i>	Angiosperm (Monocots)	2665	4462660 bp	1.17%
<i>Annona glabra</i>	Angiosperm (Magnoliids)	25657	20847775 bp	2.03%
<i>Anthoceros angustatus</i>	Homwort	3437	1965941 bp	1.65%
<i>Arabidopsis arenosa</i>	Angiosperm (Eudicots)	2783	1490506 bp	0.94%
<i>Arabidopsis suecica</i>	Angiosperm (Eudicots)	8107	4658733 bp	1.71%
<i>Arabidopsis thaliana</i>	Angiosperm (Eudicots)	2984	1920403 bp	1.60%
<i>Arabis alpina</i>	Angiosperm (Eudicots)	16338	11973226 bp	3.84%
<i>Biochara stricta</i>	Angiosperm (Eudicots)	12589	6401409 bp	3.36%
<i>Brachypodium distachyon</i>	Angiosperm (Monocots)	3723	6342990 bp	2.34%
<i>Brassica carinata</i>	Angiosperm (Eudicots)	85269	63811573 bp	5.87%
<i>Brassica juncea</i>	Angiosperm (Eudicots)	41813	28303004 bp	3.03%
<i>Brassica napus</i>	Angiosperm (Eudicots)	59041	4710812 bp	4.70%
<i>Brassica oleracea</i>	Angiosperm (Eudicots)	39036	21131203 bp	3.95%
<i>Brassica oleracea</i>	Angiosperm (Eudicots)	35262	15686904 bp	3.21%
<i>Brassica rapa</i>	Angiosperm (Eudicots)	13956	9020397 bp	2.50%
<i>Camellia sasanqua</i>	Angiosperm (Eudicots)	33282	17424140 bp	2.27%
<i>Camellia sinensis</i>	Angiosperm (Eudicots)	136002	82099863 bp	2.64%
<i>Cannicum annuum</i>	Angiosperm (Eudicots)	111672	91257570 bp	2.84%
<i>Carex papposa</i>	Angiosperm (Eudicots)	1508	650883 bp	0.17%
<i>Ceratodon purpureus</i>	Moss	21246	21881291 bp	6.04%
<i>Cerastium richardii</i>	Fern	708005	27688843 bp	3.71%
<i>Chenopodium</i>	Algae (Charales)	24602	19951973 bp	1.09%
<i>Chlamydomonas reinhardtii</i>	Algae (Chlorophyta)	92	113813 bp	0.10%
<i>Chimonium microanthum</i>	Angiosperm (Magnoliids)	29649	18022208 bp	2.47%
<i>Citrullus lanatus</i>	Angiosperm (Eudicots)	19744	8613031 bp	2.28%
<i>Citrus sinensis</i>	Angiosperm (Eudicots)	13969	7164638 bp	2.39%
<i>Citrus melo</i>	Angiosperm (Eudicots)	37829	2274973 bp	5.19%
<i>Cucumis sativus</i>	Angiosperm (Eudicots)	7159	3542145 bp	1.08%
<i>Cycas pectinata</i>	Gymnosperm	31104	17769842 bp	0.17%
<i>Cladonia oblonga</i>	Angiosperm (Eudicots)	28497	8067990 bp	1.62%
<i>Cynara cardunculus</i>	Angiosperm (Eudicots)	27138	11142394 bp	1.24%
<i>Cyperus esculentus</i>	Angiosperm (Monocots)	6570	3794589 bp	1.28%
<i>Datura stramonium</i>	Angiosperm (Eudicots)	18010	12733399 bp	0.64%
<i>Datura stramonium</i>	Angiosperm (Eudicots)	20398	8833851 bp	2.10%
<i>Dendrobium officinale</i>	Angiosperm (Monocots)	53725	2478643 bp	2.02%
<i>Fraxinus americana</i>	Angiosperm (Eudicots)	39861	2674462 bp	3.39%
<i>Fraxinus vesca</i>	Angiosperm (Eudicots)	6574	4042406 bp	1.89%
<i>Ginkgo biloba</i>	Gymnosperm	13928	5211600 bp	0.20%
<i>Glycine max</i>	Angiosperm (Eudicots)	52255	3700160 bp	3.78%
<i>Glycine soja</i>	Angiosperm (Eudicots)	46229	33166672 bp	3.27%
<i>Gnetum montanum</i>	Gymnosperm	19918	4694992 bp	0.22%
<i>Gossypium arboreum</i>	Angiosperm (Eudicots)	37786	1779843 bp	1.40%
<i>Gossypium hirsutum</i>	Angiosperm (Eudicots)	61576	40838996 bp	1.77%
<i>Gossypium raimondii</i>	Angiosperm (Eudicots)	24357	21626996 bp	2.88%
<i>Ipomoea batatas</i>	Angiosperm (Eudicots)	20242	9622488 bp	1.15%
<i>Isoetes engelmannii</i>	Lycopods	26177	9098791 bp	1.42%
<i>Isoetes taiwanensis</i>	Lycopods	63161	41043378 bp	2.48%
<i>Isotria medeoloides</i>	Moss	25252	15939378 bp	4.11%
<i>Lactuca sativa</i>	Angiosperm (Eudicots)	61984	32148664 bp	1.24%
<i>Lemna minor</i>	Angiosperm (Monocots)	5655	3485088 bp	0.94%
<i>Lycium barbarum</i>	Angiosperm (Eudicots)	32113	19780593 bp	1.62%
<i>Lycopersicon esculentum</i>	Lycopods	137085	104243086 bp	4.52%
<i>Mahoe domestica</i>	Angiosperm (Eudicots)	27620	13272521 bp	1.89%
<i>Mangifera indica</i>	Angiosperm (Eudicots)	37851	18456740 bp	4.71%
<i>Marchantia paleacea</i>	Liverworts	1413	622907 bp	0.25%
<i>Marchantia polymorpha</i>	Liverworts	3440	3009238 bp	1.36%
<i>Medicago truncatula</i>	Angiosperm (Eudicots)	22620	8031719 bp	1.87%
<i>Musa acuminata</i>	Angiosperm (Monocots)	7992	2427137 bp	0.51%
<i>Nandina domestica</i>	Angiosperm (Eudicots)	22620	8378896 bp	1.04%
<i>Nicotiana glauca</i>	Angiosperm (Eudicots)	1038444	68166679 bp	2.24%
<i>Nymphula colorata</i>	Angiosperm (Basal)	25524	9616460 bp	2.35%
<i>Olea europaea</i>	Angiosperm (Eudicots)	65866	3738667 bp	2.84%
<i>Oryza sativa</i>	Angiosperm (Monocots)	11153	14539053 bp	3.30%
<i>Phalaenopsis apicaris</i>	Angiosperm (Monocots)	57005	16838695 bp	1.88%
<i>Phaseolus vulgaris</i>	Angiosperm (Eudicots)	24958	15143968 bp	2.91%
<i>Phaseolus distachyoides</i>	Angiosperm (Monocots)	48516	27909033 bp	3.62%
<i>Physalis pubescens</i>	Angiosperm (Eudicots)	20473	9102243 bp	0.66%
<i>Physcomitrium patens</i>	Moss	3541	2374818 bp	0.20%
<i>Pisacia vera</i>	Angiosperm (Eudicots)	31656	24951778 bp	3.72%
<i>Populus trichocarpa</i>	Angiosperm (Eudicots)	17118	12993860 bp	3.31%
<i>Prunus persica</i>	Angiosperm (Eudicots)	13747	16312160 bp	7.17%
<i>Prunus pratincola</i>	Angiosperm (Eudicots)	6254	3245149 bp	1.01%
<i>Salvia hispanica</i>	Angiosperm (Eudicots)	24728	12169842 bp	3.79%
<i>Salvia splendens</i>	Angiosperm (Eudicots)	70400	5335360 bp	6.62%
<i>Sclerodermis mollis</i>	Lycopods	23300	2352378 bp	11.08%
<i>Setaria italica</i>	Angiosperm (Monocots)	26641	23395223 bp	5.52%
<i>Solanum lycopersicum</i>	Angiosperm (Eudicots)	28570	15859714 bp	1.92%
<i>Solanum melongena</i>	Angiosperm (Eudicots)	13443	6706113 bp	0.28%
<i>Solanum peltocarpum</i>	Angiosperm (Eudicots)	25814	13190001 bp	1.63%
<i>Solanum tuberosum</i>	Angiosperm (Eudicots)	25737	12975252 bp	1.53%
<i>Solanum tuberosum</i>	Angiosperm (Eudicots)	22497	8997002 bp	1.40%
<i>Sphagnum fallax</i>	Moss	27429	14202558 bp	3.99%
<i>Sphagnum magellanicum</i>	Moss	29317	15966203 bp	3.50%
<i>Thapsia galeata</i>	Gymnosperm	33411	262641406 bp	2.99%
<i>Utricularia gibba</i>	Angiosperm (Eudicots)	3442	1452433 bp	1.44%
<i>Vanilla planifolia</i>	Angiosperm (Monocots)	2199	4975384 bp	0.85%
<i>Vicia sativa</i>	Angiosperm (Eudicots)	78506	49529711 bp	3.00%
<i>Vitis vulpina</i>	Angiosperm (Eudicots)	19462	12484505 bp	2.57%
<i>Wetzelia australis</i>	Gymnosperm	67861	42157800 bp	0.61%
<i>Zea mays</i>	Angiosperm (Monocots)	59124	46562464 bp	2.13%
<i>Astragalus guberaculifer</i>	Chlorophyte	0	0	0.00%
<i>Autumnella prasinococcoides (ITEX 0718)</i>	Chlorophyte	41	72408 bp	0.07%
<i>Autumnella prasinococcoides (ITEX 23)</i>	Chlorophyte	14	3451 bp	0.02%
<i>Bathycoccus prasinus</i>	Chlorophyte	0	0	0.00%
<i>Chlamydomonas euxina</i>	Chlorophyte	90	36841 bp	0.06%
<i>Chlamydomonas reinhardtii</i>	Chlorophyte	1189	260191 bp	0.20%
<i>Chlamydomonas reinhardtii</i>	Chlorophyte	170	143750 bp	0.13%
<i>Chlamydomonas reinhardtii</i>	Chlorophyte	142	213990 bp	0.16%
<i>Chlamydomonas sp. IWO 241</i>	Chlorophyte	5794	849942 bp	0.40%
<i>Chlorella devicatus (nom. nud.) (ITEX 2437)</i>	Chlorophyte	0	0	0.00%
<i>Chlorella devicatus (nom. nud.) (ITEX 2526)</i>	Chlorophyte	27	43564 bp	0.20%
<i>Chlorella obata</i>	Chlorophyte	0	0	0.00%
<i>Chlorella sorokiniana</i>	Chlorophyte	316	1092893 bp	1.83%
<i>Chlorella sordidula</i>	Chlorophyte	235	127400 bp	0.20%
<i>Chlorella prasinus (CCMP 1205)</i>	Chlorophyte	0	0	0.00%
<i>Chlorella prasinus (CCMP 136)</i>	Chlorophyte	0	0	0.00%
<i>Cocconeis sp. C61</i>	Chlorophyte	0	0	0.00%
<i>Cocconeis subellipsoidea C-109</i>	Chlorophyte	626	380477 bp	0.78%
<i>Dunaliella salina</i>	Chlorophyte	5896	1223833 bp	0.36%
<i>Eudorckia abeyaratna</i>	Chlorophyte	178	358852 bp	0.28%
<i>Haematoxylon lacustris</i>	Chlorophyte	8464	3116883 bp	1.01%
<i>Heliosporidium sp. ATCC 30920</i>	Chlorophyte	0	0	0.00%
<i>Merismium condurci</i>	Chlorophyte	818	652321 bp	1.07%
<i>Merismium comoides</i>	Chlorophyte	0	0	0.00%
<i>Merismium pusillum (CCMP1345)</i>	Chlorophyte	0	0	0.00%
<i>Monoraphidium minutum</i>	Chlorophyte	0	0	0.00%
<i>Monoraphidium neglectum</i>	Chlorophyte	0	0	0.00%
<i>Ostreococcus anandianus</i>	Chlorophyte	3115	3160941 bp	2.08%
<i>Ostreococcus lacticornis</i>	Chlorophyte	0	0	0.00%
<i>Ostreococcus tauri (BCC1115)</i>	Chlorophyte	56	41913 bp	0.28%
<i>Ostreococcus tauri (BCC1221)</i>	Chlorophyte	34	19657 bp	0.15%
<i>Palmosiphonia sp. IFO-701</i>	Chlorophyte	0	0	0.00%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	96	175748 bp	1.18%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	96	175748 bp	1.18%
<i>Picochlorum prasinoides</i>	Chlorophyte	67	25558 bp	0.11%
<i>Raphidocelis subcapitata</i>	Chlorophyte	0	0	0.00%
<i>Scenedesmus sp. NREL 468-D3</i>	Chlorophyte	270	1371745 bp	0.09%
<i>Tetrahymena socialis</i>	Chlorophyte	292	129488 bp	0.10%
<i>Trebouxia sp. AI-2</i>	Chlorophyte	60	29089 bp	0.05%
<i>Tilapia africana</i>	Chlorophyte	646	508891 bp	0.39%
<i>Fobos reticuliferus (NIES 3786)</i>	Chlorophyte	1698	891661 bp	0.67%
<i>Fobos reticuliferus (NIES-3785)</i>	Chlorophyte	907	552913 bp	0.42%

Upon examining the distribution of DNA transposons across the analyzed plant lineages, it becomes evident that these transposons are not predominant in any lineage. However, a notable pattern emerges where the lineages of Angiosperms, Fern, Hornwort, Moss, Gymnosperms, and Lycopods consistently harbor a considerably higher abundance of DNA transposons compared to the Chlorophyte lineage. Within the Chlorophyte genomes studied, approximately one-third of the species lack DNA transposons entirely. Among the remaining Chlorophyte genomes, the majority possess a representation of DNA transposons accounting for less than 1%. Yet, mirroring the LTR element trend, *Ostreobium quekettii* stands out with the highest proportion of DNA transposons at 2.08%.

Across the other examined lineages, the proportion of DNA transposons varies considerably, ranging from 0.17% in *Cycas panzhihuaensis* (Gymnosperm) to a noteworthy 11.08% in *Selaginella moellendorffii* (Lycopod).

For most Angiosperms, the DNA transposon representation falls within the range of 1%-5%. Conversely, most Gymnosperms display a more subdued representation, with the DNA transposon proportion typically being less than 1%.

The presented boxplots delineate the %Divergence metrics for these elements.



For most plant lineages, including Angiosperms, Hornwort, Moss, Gymnosperms, Lycopods, and Liverworts, the median %Divergence of DNA predominantly falls within the range of 15-20. This consistent median %Divergence across diverse lineages may suggest a fundamental evolutionary constraint or a shared genomic stability mechanism across these plant groups. However, there are exceptions. For instance, the median values for *Mangifera indica* and *Utricularia gibba* (both Angiosperm - Eudicots) are approximately 22, while those for *Sphagnum fallax* and *Sphagnum magellanicum* (both Moss) are around 23 and 24, respectively. Such deviations could be indicative of lineage-specific evolutionary pressures or unique genomic dynamics in these species, potentially related to environmental adaptations or specific genomic rearrangements.

Across these lineages, the Interquartile Range (IQR) remains relatively consistent, approximately around 15. Yet, certain lineages exhibit narrower IQRs. Notably, species such as *Phalaenopsis equestris* (Angiosperm - Monocots), *Physcomitrium patens*, *Sphagnum fallax*, *Sphagnum magellanicum* (all Moss), *Ginkgo biloba* (Gymnosperms), and *Annona alabra* (Angiosperm - Magnolids) display narrower IQRs, with their

median values also relatively elevated. This could imply a more constrained genomic variability within these species or lineages, possibly reflecting specialized ecological niches or specific developmental pathways.

Furthermore, while the whisker lengths across most lineages range from 0 to 40, exceptions like *Salvia splendens* (Angiosperm - Luidicots) present shorter whiskers and exhibit numerous outliers surpassing the maximum divergence percentage. Such outliers may signify rapid genomic changes or unique evolutionary trajectories in these particular species, warranting further investigation into their genomic dynamics and potential adaptive significance.

In the Chlorophyte lineage, a distinct pattern emerges in the %Divergence of DNA transposons. The median values across species in this lineage vary widely, ranging from 0 to 18. Particularly, *Chlamydomonas incerta* exhibits the highest median value, whereas *Chlorella sorokiniana* and *Chlorella variabilis* approach a median value near 0, with numerous outliers beyond the maximum value. Additionally, the IQR lengths differ among species in this lineage, with *Picochlorum sp.BPE23* having an IQR of approximately 15, while *Chlorella sorokiniana*'s IQR is close to 0. The whisker lengths for these species also vary, spanning from 0 to 42, with a majority having outliers beyond the maximum value. These variations may reflect diverse genomic architectures or differential evolutionary pressures within the Chlorophyte lineage, highlighting the need for comprehensive genomic studies to elucidate the underlying mechanisms and functional implications.

2.1.3 Unclassified Element Comparative and Analysis

Unclassified elements refer to segments within the genomes that, despite being identified as repetitive elements, do not fit into the conventional categories or classifications like LTR elements or DNA transposons. Their functional significance or origin may not be fully elucidated, leading to their current unclassified status.

The following table lists the number, length, and genome percentage of unclassified elements for each of the 138 plant genomes.

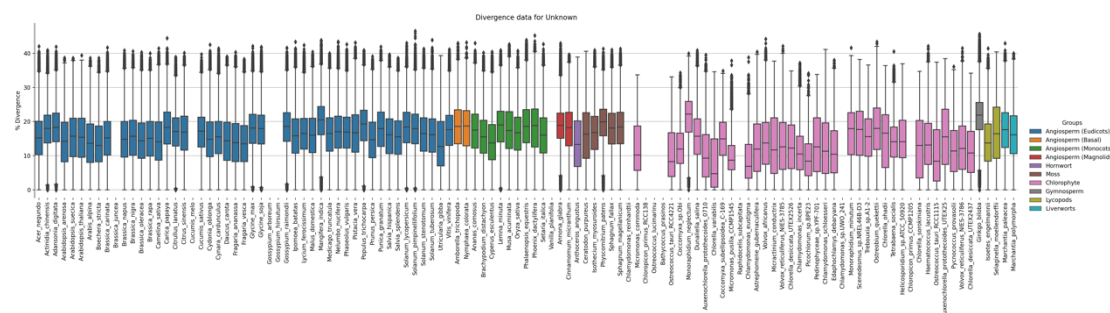
SPECIES	LINEAGE	Unclassified (number of elements)	Unclassified (length occupied)	Unclassified (percentage of sequence)
<i>Acer negundo</i>	Angiosperm (Eudicots)	544760	130037607 bp	29.39%
<i>Acer yangbiense</i>	Angiosperm (Eudicots)	577505	135738787 bp	20.38%
<i>Actinidia chinensis</i>	Angiosperm (Eudicots)	744075	138377348 bp	24.98%
<i>Adansonia digitata</i>	Angiosperm (Eudicots)	1144391	26414665 bp	38.49%
<i>Adiantum capillus</i>	Fern	3749055	131801763 bp	27.33%
<i>Amborella trichopoda</i>	Angiosperm (Basal)	945500	232627118 bp	32.93%
<i>Ananas comosus</i>	Angiosperm (Monocots)	452670	116046036 bp	30.37%
<i>Annona glabra</i>	Angiosperm (Magnoliids)	1439494	371633160 bp	36.18%
<i>Anthoeceros angustus</i>	Hornwort	86679	33324172 bp	27.92%
<i>Arabidopsis arenosa</i>	Angiosperm (Eudicots)	53727	15920018 bp	10.64%
<i>Arabidopsis suecica</i>	Angiosperm (Eudicots)	88525	32856068 bp	12.07%
<i>Arabidopsis thaliana</i>	Angiosperm (Eudicots)	21921	7459066 bp	6.23%
<i>Arabis alpina</i>	Angiosperm (Eudicots)	170695	54991419 bp	17.65%
<i>Boehea stricta</i>	Angiosperm (Eudicots)	89134	24168195 bp	12.68%
<i>Brachypodium distachyon</i>	Angiosperm (Monocots)	175368	37044778 bp	13.65%
<i>Brassica carinata</i>	Angiosperm (Eudicots)	720593	230304662 bp	21.19%
<i>Brassica juncea</i>	Angiosperm (Eudicots)	573264	217141236 bp	23.26%
<i>Brassica napus</i>	Angiosperm (Eudicots)	637744	239012010 bp	23.86%
<i>Brassica nigra</i>	Angiosperm (Eudicots)	290455	106523763 bp	19.94%
<i>Brassica oleracea</i>	Angiosperm (Eudicots)	438284	122432596 bp	25.04%
<i>Brassica rapa</i>	Angiosperm (Eudicots)	212763	83552437 bp	23.67%
<i>Camellia sativa</i>	Angiosperm (Eudicots)	330593	103209815 bp	16.09%
<i>Camellia sinensis</i>	Angiosperm (Eudicots)	222404	592622309 bp	19.08%
<i>Capsicum annuum</i>	Angiosperm (Eudicots)	2357972	782159937 bp	24.35%
<i>Carica papaya</i>	Angiosperm (Eudicots)	177433	69989504 bp	18.89%
<i>Ceratodon purpureus</i>	Moss	209537	108539954 bp	29.94%
<i>Ceratopteris richardii</i>	Fern	4742183	1521578784 bp	20.39%
<i>Chara braunii</i>	Algae (Charales)	2239563	526321037 bp	30.05%
<i>Chlamydomonas reinhardtii</i>	Algae (Chlorophyta)	27258	6751800 bp	6.68%
<i>Cinnamomum micranthum</i>	Angiosperm (Magnoliids)	849081	180647370 bp	24.73%
<i>Citrus latifolia</i>	Angiosperm (Eudicots)	473072	116557476 bp	32.25%
<i>Citrus sinensis</i>	Angiosperm (Eudicots)	207264	60778753 bp	20.27%
<i>Cucumis melo</i>	Angiosperm (Eudicots)	418573	142851878 bp	32.58%
<i>Cucumis sativus</i>	Angiosperm (Eudicots)	178666	44910196 bp	19.82%
<i>Cycas panhellenensis</i>	Gymnosperm	743503	274544065 bp	26.17%
<i>Cydonia oblonga</i>	Angiosperm (Eudicots)	744719	157526361 bp	32.25%
<i>Cynara cardunculus</i>	Angiosperm (Eudicots)	891700	191402492 bp	26.40%
<i>Cyperus esculentus</i>	Angiosperm (Monocots)	326916	73495076 bp	24.78%
<i>Danara stramonium</i>	Angiosperm (Eudicots)	997980	352427016 bp	17.85%
<i>Daucus carota</i>	Angiosperm (Eudicots)	528592	89948690 bp	21.34%
<i>Dendrobium officinale</i>	Angiosperm (Monocots)	1052446	328746339 bp	26.36%
<i>Fragaria ananassa</i>	Angiosperm (Eudicots)	418155	119347769 bp	14.81%
<i>Fragaria vesca</i>	Angiosperm (Eudicots)	124784	33812268 bp	15.77%
<i>Ginkgo biloba</i>	Gymnosperm	2095638	56808060 bp	21.53%
<i>Glycine max</i>	Angiosperm (Eudicots)	659719	197258115 bp	20.15%
<i>Glycine soja</i>	Angiosperm (Eudicots)	657508	190551905 bp	18.80%
<i>Gnetum montanum</i>	Gymnosperm	151794	53739456 bp	25.02%
<i>Grossypium arboreum</i>	Angiosperm (Eudicots)	1085306	419635130 bp	25.88%
<i>Grossypium hirsutum</i>	Angiosperm (Eudicots)	2229684	806786366 bp	34.99%
<i>Grossypium raimondii</i>	Angiosperm (Eudicots)	827870	226686515 bp	30.18%
<i>Ipomoea batatas</i>	Angiosperm (Eudicots)	1341028	234230149 bp	27.98%
<i>Isoetes engelmannii</i>	Lycopods	894777	171125709 bp	26.70%
<i>Isoetes subaenariensis</i>	Lycopods	1188150	380920875 bp	22.97%
<i>Isoetes macrospora</i>	Moss	582520	172208054 bp	44.33%
<i>Lactuca sativa</i>	Angiosperm (Eudicots)	1854922	721243084 bp	27.84%
<i>Lemna minuta</i>	Angiosperm (Monocots)	262617	80694251 bp	24.88%
<i>Lycium ferocissimum</i>	Angiosperm (Eudicots)	1376314	388427203 bp	31.86%
<i>Lycopodium clavatum</i>	Lycopods	1419861	465833119 bp	20.21%
<i>Mala domestica</i>	Angiosperm (Eudicots)	490811	107174828 bp	15.24%
<i>Mangifera indica</i>	Angiosperm (Eudicots)	404622	100428467 bp	25.62%
<i>Marchantia paleacea</i>	Liverworts	91949	32983442 bp	13.15%
<i>Marchantia polymorpha</i>	Liverworts	62461	26266591 bp	11.63%
<i>Medicago truncatula</i>	Angiosperm (Eudicots)	554946	115755563 bp	26.92%
<i>Musa acuminata</i>	Angiosperm (Monocots)	206778	62323735 bp	13.20%
<i>Nelumbo macleodii</i>	Angiosperm (Eudicots)	944578	183280865 bp	22.78%
<i>Nicotiana glauca</i>	Angiosperm (Eudicots)	2439044	729796446 bp	24.04%
<i>Nymphaea colorata</i>	Angiosperm (Basal)	352541	93276435 bp	22.81%
<i>Olea europaea</i>	Angiosperm (Eudicots)	2451453	449864431 bp	34.17%
<i>Orzyza sativa</i>	Angiosperm (Monocots)	361697	79787375 bp	21.31%
<i>Phalaenopsis equestris</i>	Angiosperm (Monocots)	1123275	387859297 bp	36.45%
<i>Phaseolus vulgaris</i>	Angiosperm (Eudicots)	387921	106343957 bp	20.41%
<i>Phoenix dactylopera</i>	Angiosperm (Monocots)	385815	101231084 bp	13.09%
<i>Physalis pubescens</i>	Angiosperm (Eudicots)	702077	355823893 bp	25.61%
<i>Physcomitrium patens</i>	Moss	105865	34376740 bp	7.28%
<i>Pistacia vera</i>	Angiosperm (Eudicots)	450625	114620514 bp	17.07%
<i>Populus trichocarpa</i>	Angiosperm (Eudicots)	399467	103447188 bp	26.37%
<i>Prunus persica</i>	Angiosperm (Eudicots)	151840	410061084 bp	18.02%
<i>Punica granatum</i>	Angiosperm (Eudicots)	269529	89840239 bp	28.03%
<i>Salvia hispanica</i>	Angiosperm (Eudicots)	327757	72564732 bp	22.57%
<i>Salvia splendens</i>	Angiosperm (Eudicots)	627415	155463347 bp	19.29%
<i>Selaginella moellendorffii</i>	Lycopods	79062	35867924 bp	16.89%
<i>Setaria italica</i>	Angiosperm (Monocots)	239840	67729328 bp	16.69%
<i>Solanum lycopersicum</i>	Angiosperm (Eudicots)	710452	219993742 bp	26.79%
<i>Solanum melongena</i>	Angiosperm (Eudicots)	962324	319296233 bp	27.42%
<i>Solanum pimpinellifolium</i>	Angiosperm (Eudicots)	725276	238023301 bp	29.45%
<i>Solanum stenotomum</i>	Angiosperm (Eudicots)	787416	209450434 bp	24.75%
<i>Solanum tuberosum</i>	Angiosperm (Eudicots)	751728	197190543 bp	27.93%
<i>Sphagnum fallax</i>	Moss	528940	313199787 bp	33.71%
<i>Sphagnum magellanicum</i>	Moss	566915	14498732 bp	33.03%
<i>Thaui plicata</i>	Gymnosperm	5112945	1987481838 bp	21.85%
<i>Utricularia gibba</i>	Angiosperm (Eudicots)	46760	17420498 bp	17.30%
<i>Vanilla planifolia</i>	Angiosperm (Monocots)	750493	488967057 bp	34.52%
<i>Vicia sativa</i>	Angiosperm (Eudicots)	875356	294064471 bp	17.78%
<i>Vitis vinifera</i>	Angiosperm (Eudicots)	368225	113188427 bp	23.28%
<i>Webbia biala</i>	Gymnosperm	6734176	2104830097 bp	30.51%
<i>Zea mays</i>	Angiosperm (Monocots)	837546	259709792 bp	16.88%
<i>Astraphomena gubernaculifera</i>	Chlorophyte	40070	11349855 bp	10.93%
<i>Auxenochlorella protohemicoides (0710)</i>	Chlorophyte	48692	12594245 bp	12.13%
<i>Auxenochlorella protohemicoides (UTEX 25)</i>	Chlorophyte	648	319987 bp	1.51%
<i>Bathycoccus prasinos</i>	Chlorophyte	1201	354784 bp	2.35%
<i>Chlamydomonas axicoma</i>	Chlorophyte	7174	1751837 bp	2.63%
<i>Chlamydomonas incerta</i>	Chlorophyte	80379	11320685 bp	8.76%
<i>Chlamydomonas reinhardtii</i>	Chlorophyte	33671	7146352 bp	6.43%
<i>Chlamydomonas schlesseri</i>	Chlorophyte	49746	12415401 bp	9.54%
<i>Chlamydomonas sp. UWO 241</i>	Chlorophyte	286445	51069248 bp	24.13%
<i>Chlorella desiccata (nom. nud.) (UTEX 2437)</i>	Chlorophyte	1410	608452 bp	2.93%
<i>Chlorella desiccata (nom. nud.) (UTEX 2526)</i>	Chlorophyte	1935	659116 bp	3.06%
<i>Chlorella obadii</i>	Chlorophyte	3131	8044659 bp	14.11%
<i>Chlorella sorokiniana</i>	Chlorophyte	2347	1240924 bp	2.08%
<i>Chlorella variabilis</i>	Chlorophyte	3707	1023134 bp	2.22%
<i>Chlorococcum primum (CCMP 1205)</i>	Chlorophyte	2969	737946 bp	4.24%
<i>Chlorococcum primum (RCC138)</i>	Chlorophyte	3451	850942 bp	4.84%
<i>Coccomyxa sp. Ch1</i>	Chlorophyte	5931	884838 bp	1.76%
<i>Coccomyxa subellipsoidea C-169</i>	Chlorophyte	9132	1062252 bp	2.18%
<i>Dundalella salina</i>	Chlorophyte	323978	61259470 bp	17.82%
<i>Edaphochlamys debaryana</i>	Chlorophyte	36483	11210618 bp	7.89%
<i>Haematococcus lacustris</i>	Chlorophyte	264151	75750219 bp	24.49%
<i>Heliosporidium sp. ATCC 50920</i>	Chlorophyte	1131	223095 bp	1.80%
<i>Micractinium concolor</i>	Chlorophyte	12812	4256088 bp	6.98%
<i>Micromonas commoda</i>	Chlorophyte	1604	347194 bp	1.64%
<i>Micromonas pusilla CCMP1545</i>	Chlorophyte	10660	1519793 bp	6.92%
<i>Monoraphidium minutum</i>	Chlorophyte	8628	2049762 bp	3.01%
<i>Monoraphidium neglectum</i>	Chlorophyte	10102	1442908 bp	2.07%
<i>Ostreobium queketti</i>	Chlorophyte	77186	23765604 bp	15.65%
<i>Ostreococcus lucimarinus</i>	Chlorophyte	1033	570596 bp	4.32%
<i>Ostreococcus tauri (RCC1115)</i>	Chlorophyte	1185	4047316 bp	2.74%
<i>Ostreococcus tauri (RCC4221)</i>	Chlorophyte	860	235602 bp	1.81%
<i>Pedonophyceae sp. YPF-701</i>	Chlorophyte	2127	778667 bp	2.79%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	957	588300 bp	3.94%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	957	588300 bp	3.94%
<i>Picochlorum prasinoviride</i>	Chlorophyte	2302	510174 bp	2.25%
<i>Raphidocelis subcapitata</i>	Chlorophyte	3252	772664 bp	1.51%
<i>Scenedesmus sp. NREL 46B-D3</i>	Chlorophyte	107676	17972991 bp	11.83%
<i>Tetraena socialis</i>	Chlorophyte	69821	10668108 bp	7.86%
<i>Trebouxia sp. AT-2</i>	Chlorophyte	8064	1417979 bp	2.68%
<i>Tetraselmis africana</i>	Chlorophyte	43997	12486256 bp	9.65%
<i>Tetraselmis africana (NIES-3786)</i>	Chlorophyte	47983	17130208 bp	12.79%
<i>Tetraselmis africana (NIES-3785)</i>	Chlorophyte	45359	16178529 bp	12.16%

Across the spectrum of plant species analyzed, unclassified elements are universally present. However, a pronounced pattern emerges wherein the lineages of Angiosperms, Fern, Hornwort, Moss, Gymnosperms, and Lycopods consistently exhibit a notably higher abundance of unclassified elements in terms of quantity, length, and proportion compared to the Chlorophyte lineage.

Within the Chlorophyte genomes, the proportion of unclassified elements spans a range from 1.51% in *Auxenochlorella protothecoides* (UTEX 25) to a substantial 24.49% in *Haematococcus lacustris*. The majority of these Chlorophyte species typically present a representation of unclassified elements falling between 1-10%.

In contrast, the other examined lineages showcase a broader spectrum of unclassified element proportions, ranging from 6.23% in *Arabidopsis thaliana* (Angiosperm) to a striking 44.33% in *Isoetecium myosuroides* (Moss). Most plant lineages outside Chlorophytes display a prevalence of unclassified elements ranging between 10-35%.

The presented boxplots delineate the %Divergence metrics for these elements.



For the Unknown elements within the Angiosperms, Hornwort, Gymnosperms, Lycopods, and Liverworts lineages, the %Divergence exhibits overall stability, with median values ranging between 15 and 20%. However, Gymnosperm stands out with a notably higher median value, approximately around 23%. Across these lineages, the IQR values hover around 10, and whisker lengths span from 0 to approximately 35. It's noteworthy that, except for Hornwort, the other lineages showcase a significant number of outliers surpassing the maximum %Divergence value.

In contrast, within the Chlorophyte lineage, the Unknown elements present a distinctly lower median %Divergence range, predominantly falling between 5 and 15%. A notable exception is observed with *Monoraphidium neglectum*, which boasts the highest median value, approximately 22%. The IQR values among Chlorophyte species vary, with narrower ranges observed in species like *Micromonas pusilla* CCMP1545, approximately around 5, and broader ranges in species like *Auxenochlorella protothecoides* (UTEX25), which is approximately 15. The whisker lengths among these species are diverse, and roughly half of them exhibit outliers that exceed the maximum %Divergence value.

2.1.4 Small RNA Comparative and Analysis

Across the diverse landscape of plant species analyzed, the representation of Small RNA remains notably sparse.

Within the Chlorophyte lineage, a striking pattern emerges where nearly half of the species surveyed lack discernible Small RNA. However, certain exceptions stand out. Notably, *Volvox africanus* displays a Small RNA proportion of 6.31%, followed closely by *Astrephomene gubernaculifera* at 5.72% and *Auxenochlorella protothecoides (0710)* at 5.62%. The remaining species within this lineage typically exhibit Small RNA proportions ranging between 0-0.68%.

In contrast, among the lineages of Angiosperms, Fern, Hornwort, Moss, Gymnosperms, and Lycopods, *Cucumis melo* (Angiosperm) records the highest Small RNA representation at 4.53%. Intriguingly, the Fern species *Ceratopteris richardii* stands as a unique outlier, devoid of detectable Small RNA. For the majority of species within these lineages, the proportion of Small RNA remains predominantly below the 1% threshold.

2.1.5 Simple Repeat Element Comparative and Analysis

Simple repeats, often referred to as microsatellites or short tandem repeats (STRs), are short sequences of DNA motifs that are repeated consecutively multiple times. These repetitive sequences are ubiquitous across genomes and play roles in various genomic processes, including gene regulation, chromosomal organization, and genome stability. Their abundance and variability can offer insights into the evolutionary dynamics and genetic makeup of organisms.

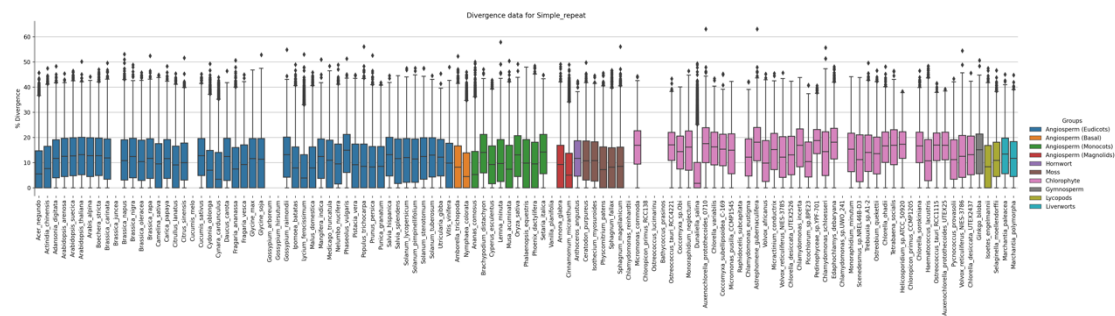
The following table lists the number, length, and genome percentage of simple repeat elements for each of the 138 plant genomes.

SPECIES	LINEAGE	Simple repeats (number of elements)	Simple repeats (length occupied)	Simple repeats (percentage of sequence)
<i>Acer negundo</i>	Angiosperm (Eudicots)	216119	7980535 bp	1.80%
<i>Acer yungbiense</i>	Angiosperm (Eudicots)	239501	9299191 bp	1.40%
<i>Acridula chinensis</i>	Angiosperm (Eudicots)	156921	588283 bp	1.06%
<i>Adiantum digitatum</i>	Angiosperm (Eudicots)	226504	8931992 bp	1.30%
<i>Adiantum capillus</i>	Fern	1150482	62486243 bp	1.30%
<i>Amborella trichopoda</i>	Angiosperm (Basal)	168689	10254994 bp	1.45%
<i>Ananas comosus</i>	Angiosperm (Monocots)	170938	7886524 bp	2.06%
<i>Annona glabra</i>	Angiosperm (Magnoliids)	183931	6681055 bp	0.65%
<i>Anhoecerus angustatus</i>	Horowitz	18880	1004233 bp	0.84%
<i>Arabidopsis arenosa</i>	Angiosperm (Eudicots)	44312	1740104 bp	1.16%
<i>Arabidopsis suecica</i>	Angiosperm (Eudicots)	78273	3122193 bp	1.15%
<i>Arabidopsis thaliana</i>	Angiosperm (Eudicots)	34625	1364872 bp	1.14%
<i>Arabis alpina</i>	Angiosperm (Eudicots)	46948	1929752 bp	0.62%
<i>Boehmeria stricta</i>	Angiosperm (Eudicots)	51062	2123258 bp	1.11%
<i>Brachypodium distachyon</i>	Angiosperm (Monocots)	44534	1868055 bp	0.69%
<i>Brassica carinata</i>	Angiosperm (Eudicots)	199749	11572338 bp	1.06%
<i>Brassica juncea</i>	Angiosperm (Eudicots)	166310	6920826 bp	0.74%
<i>Brassica napus</i>	Angiosperm (Eudicots)	183943	7918357 bp	0.79%
<i>Brassica nigra</i>	Angiosperm (Eudicots)	85814	4050821 bp	0.76%
<i>Brassica oleracea</i>	Angiosperm (Eudicots)	91722	3838788 bp	0.79%
<i>Brassica rapa</i>	Angiosperm (Eudicots)	74227	3079625 bp	0.87%
<i>Camelina sativa</i>	Angiosperm (Eudicots)	186422	7029626 bp	1.10%
<i>Camellia sinensis</i>	Angiosperm (Eudicots)	630211	25595728 bp	0.82%
<i>Capsicum annuum</i>	Angiosperm (Eudicots)	279674	12699613 bp	0.40%
<i>Carica papaya</i>	Angiosperm (Eudicots)	102294	4132927 bp	1.12%
<i>Ceratodon purpureus</i>	Moss	55652	2711214 bp	0.75%
<i>Ceratopteris richardii</i>	Fern	1025552	417886560 bp	5.60%
<i>Chara braunii</i>	Algae (Charales)	785394	41815766 bp	2.39%
<i>Chlamydomonas reinhardtii</i>	Algae (Chlorophyta)	134153	8856305 bp	7.97%
<i>Cinnamomum micranthum</i>	Angiosperm (Magnoliids)	239595	8282682 bp	1.13%
<i>Citrus limon</i>	Angiosperm (Eudicots)	145137	5442063 bp	1.51%
<i>Citrus sinensis</i>	Angiosperm (Eudicots)	85262	3255718 bp	1.09%
<i>Cucumis melo</i>	Angiosperm (Eudicots)	130413	5631800 bp	1.28%
<i>Cucumis sativus</i>	Angiosperm (Eudicots)	91557	3703318 bp	1.63%
<i>Cycas panzhuanaensis</i>	Gymnosperm	964483	168403281 bp	1.61%
<i>Cydonia oblonga</i>	Angiosperm (Eudicots)	132844	4997358 bp	1.02%
<i>Cypripedium calceolatum</i>	Angiosperm (Eudicots)	153792	655300 bp	0.90%
<i>Cyperus excelsus</i>	Angiosperm (Monocots)	123005	5825357 bp	1.96%
<i>Datura stramonium</i>	Angiosperm (Eudicots)	164338	8307676 bp	0.42%
<i>Daucus carota</i>	Angiosperm (Eudicots)	73200	3309077 bp	0.78%
<i>Dendrobium officinale</i>	Angiosperm (Monocots)	227153	14108215 bp	1.15%
<i>Fragaria ananassa</i>	Angiosperm (Eudicots)	187697	7576361 bp	0.94%
<i>Fragaria vesca</i>	Angiosperm (Eudicots)	55480	2191992 bp	1.02%
<i>Ginkgo biloba</i>	Gymnosperm	134293	6085465 bp	0.23%
<i>Glycine max</i>	Angiosperm (Eudicots)	239175	10744244 bp	1.10%
<i>Glycine soja</i>	Angiosperm (Eudicots)	244417	11057772 bp	1.09%
<i>Gnetum montanum</i>	Gymnosperm	101453	4563683 bp	0.21%
<i>Gossypium arboreum</i>	Angiosperm (Eudicots)	161236	7161254 bp	0.44%
<i>Gossypium hirsutum</i>	Angiosperm (Eudicots)	257783	11778002 bp	0.51%
<i>Gossypium raimondii</i>	Angiosperm (Eudicots)	119992	5265445 bp	0.70%
<i>Iponoea batatas</i>	Angiosperm (Eudicots)	285929	11855713 bp	1.42%
<i>Isoetes engelmannii</i>	Lycopods	100994	4776477 bp	0.75%
<i>Isoetes taiwanensis</i>	Lycopods	135930	8977573 bp	0.54%
<i>Isobacterium mossourades</i>	Moss	39166	1765016 bp	0.45%
<i>Lactuca sativa</i>	Angiosperm (Eudicots)	455619	25837336 bp	1.00%
<i>Lemna minor</i>	Angiosperm (Monocots)	62815	3948166 bp	1.10%
<i>Lycium ferocissimum</i>	Angiosperm (Eudicots)	205579	8843951 bp	0.73%
<i>Lycopersicon esculentum</i>	Lycopods	290222	16998745 bp	0.74%
<i>Mahoe domestica</i>	Angiosperm (Eudicots)	119727	4811980 bp	0.68%
<i>Mangifera indica</i>	Angiosperm (Eudicots)	90264	4658202 bp	1.19%
<i>Marchantia paleacea</i>	Liverworts	36023	1463820 bp	0.58%
<i>Marchantia polymorpha</i>	Liverworts	51777	2017629 bp	0.89%
<i>Medicago truncatula</i>	Angiosperm (Eudicots)	103286	4299562 bp	1.00%
<i>Musa acuminata</i>	Angiosperm (Monocots)	100246	4217653 bp	0.89%
<i>Nelumbo nucifera</i>	Angiosperm (Eudicots)	168601	7235483 bp	0.90%
<i>Nicotiana glauca</i>	Angiosperm (Eudicots)	290146	14545708 bp	0.48%
<i>Nymphaea colorata</i>	Angiosperm (Basal)	118284	5016483 bp	1.23%
<i>Olea europaea</i>	Angiosperm (Eudicots)	240977	9091452 bp	0.69%
<i>Oryza sativa</i>	Angiosperm (Monocots)	96355	4409375 bp	1.18%
<i>Phalenopsis equestris</i>	Angiosperm (Monocots)	180283	11034902 bp	1.04%
<i>Phaseolus vulgaris</i>	Angiosperm (Eudicots)	110406	5484400 bp	1.66%
<i>Phaseolus vulgaris</i>	Angiosperm (Eudicots)	148440	6927681 bp	0.90%
<i>Physalis pubescens</i>	Angiosperm (Eudicots)	111989	5085878 bp	0.37%
<i>Physcomitrium patens</i>	Moss	173701	6861714 bp	1.45%
<i>Pistacia vera</i>	Angiosperm (Eudicots)	153767	6014269 bp	0.90%
<i>Populus trichocarpa</i>	Angiosperm (Eudicots)	113002	4309088 bp	1.10%
<i>Prunus persica</i>	Angiosperm (Eudicots)	74760	2781607 bp	1.22%
<i>Panicum granatum</i>	Angiosperm (Eudicots)	92832	3396816 bp	1.06%
<i>Sabia hispanica</i>	Angiosperm (Eudicots)	68783	2975443 bp	0.93%
<i>Salvia splendens</i>	Angiosperm (Eudicots)	134023	5889180 bp	0.73%
<i>Selaginella moellendorffii</i>	Lycopods	35754	1665109 bp	0.78%
<i>Senecio jalticus</i>	Angiosperm (Monocots)	60686	2875102 bp	0.71%
<i>Solanum lycopersicum</i>	Angiosperm (Eudicots)	106107	6230521 bp	0.75%
<i>Solanum melongena</i>	Angiosperm (Eudicots)	122660	6336681 bp	0.54%
<i>Solanum pimpinellifolium</i>	Angiosperm (Eudicots)	112890	6479954 bp	0.80%
<i>Solanum stenotomum</i>	Angiosperm (Eudicots)	106195	5236532 bp	0.62%
<i>Solanum tuberosum</i>	Angiosperm (Eudicots)	94486	6114333 bp	0.87%
<i>Sphagnum fallax</i>	Moss	111726	4817712 bp	1.04%
<i>Sphagnum magellanicum</i>	Moss	114686	4377903 bp	1.00%
<i>Thuja plicata</i>	Gymnosperm	690845	4926713 bp	0.54%
<i>Utricularia gibba</i>	Angiosperm (Eudicots)	18972	741826 bp	0.74%
<i>Vanilla planifolia</i>	Angiosperm (Monocots)	402897	25924738 bp	18.30%
<i>Vicia sativa</i>	Angiosperm (Eudicots)	198155	11652949 bp	0.70%
<i>Vitis vinifera</i>	Angiosperm (Eudicots)	135155	5232790 bp	1.08%
<i>Weibschia mirabilis</i>	Gymnosperm	712709	59627004 bp	0.87%
<i>Zea mays</i>	Angiosperm (Monocots)	124313	6583005 bp	0.30%
<i>Astraphome gubernalifera</i>	Chlorophyte	83244	4906531 bp	4.72%
<i>Auxochloris prostratoides (0710)</i>	Chlorophyte	80659	4339408 bp	4.18%
<i>Auxochloris prostratoides (UTEX 25)</i>	Chlorophyte	6602	309359 bp	1.46%
<i>Bathycoccus prasinus</i>	Chlorophyte	14792	361766 bp	3.73%
<i>Chlamydomonas eustigma</i>	Chlorophyte	10416	461491 bp	0.69%
<i>Chlamydomonas incerta</i>	Chlorophyte	145558	9162993 bp	7.09%
<i>Chlamydomonas reinhardtii</i>	Chlorophyte	128360	7568749 bp	6.81%
<i>Chlamydomonas schlotheimia</i>	Chlorophyte	165410	9645577 bp	7.41%
<i>Chlamydomonas sp. UWO 241</i>	Chlorophyte	325473	19742756 bp	9.33%
<i>Chlorella desiccata (nom. nud.) (UTEX 2437)</i>	Chlorophyte	6670	254164 bp	1.23%
<i>Chlorella desiccata (nom. nud.) (UTEX 2526)</i>	Chlorophyte	6856	257987 bp	1.20%
<i>Chlorella ohadii</i>	Chlorophyte	47643	2464557 bp	4.32%
<i>Chlorella sorokiniana</i>	Chlorophyte	54937	3057484 bp	5.13%
<i>Chlorella variabilis</i>	Chlorophyte	58420	3038107 bp	6.58%
<i>Chlorococcus prasinus (CCMP 1205)</i>	Chlorophyte	15298	581738 bp	3.34%
<i>Chlorococcus prasinus (RCC138)</i>	Chlorophyte	15307	585947 bp	3.33%
<i>Coccomyxa sp. Obi</i>	Chlorophyte	7042	330049 bp	0.65%
<i>Coccomyxa subellipsoidea C-169</i>	Chlorophyte	7164	342922 bp	0.72%
<i>Dunaliella salina</i>	Chlorophyte	138242	8687227 bp	2.53%
<i>Edaphoclamys debariana</i>	Chlorophyte	113393	6534007 bp	4.60%
<i>Haematococcus lacustris</i>	Chlorophyte	221539	24135783 bp	7.80%
<i>Heliosporidium sp. ATCC 50920</i>	Chlorophyte	2569	114252 bp	0.92%
<i>Micractinium conductrix</i>	Chlorophyte	91518	4987179 bp	8.17%
<i>Micromonas commoda</i>	Chlorophyte	13332	652938 bp	3.09%
<i>Micromonas pusilla (CCMP1345)</i>	Chlorophyte	35056	1807403 bp	8.23%
<i>Monoraphidium minimum</i>	Chlorophyte	130603	7981439 bp	11.71%
<i>Monoraphidium neglectum</i>	Chlorophyte	69185	4174548 bp	5.99%
<i>Ostreobium queketti</i>	Chlorophyte	16001	789417 bp	0.52%
<i>Ostreococcus lucimarinus</i>	Chlorophyte	9907	613178 bp	4.64%
<i>Ostreococcus tauri (RCC1115)</i>	Chlorophyte	8176	488356 bp	3.31%
<i>Ostreococcus tauri (RCC1221)</i>	Chlorophyte	8012	470460 bp	3.61%
<i>Palmonyxa sp. YFP-701</i>	Chlorophyte	14713	815279 bp	2.92%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	2262	82192 bp	0.55%
<i>Picochlorum sp. BPE23</i>	Chlorophyte	2262	82192 bp	0.55%
<i>Pycnococcus provasolii</i>	Chlorophyte	19498	821222 bp	3.62%
<i>Raphidocelis subcapitata</i>	Chlorophyte	86247	5308884 bp	10.38%
<i>Scenedesmus sp. NREL 46B-D3</i>	Chlorophyte	137213	9654960 bp	6.36%
<i>Tetrahymena socialis</i>	Chlorophyte	77584	4335814 bp	3.19%
<i>Trebouxia sp. 41-2</i>	Chlorophyte	11383	633824 bp	1.20%
<i>Volvox africanus</i>	Chlorophyte	75462	4157060 bp	3.21%
<i>Volvox reitelbaueri (NIES 3786)</i>	Chlorophyte	67181	3742659 bp	2.79%
<i>Volvox reitelbaueri (NIES-3785)</i>	Chlorophyte	65641	3343448 bp	2.51%

Across the diverse array of plant species examined, the presence of Simple Repeats is universal. Within the lineages of Angiosperms, Fern, Hornwort, Moss, Gymnosperms, and Lycopods, the proportion of the genome occupied by Simple Repeats remains predominantly below the 2% mark. However, some exceptions capture attention. Notably, *Vanilla planifolia* (Angiosperm) stands out with a substantial representation of Simple Repeats, accounting for 18.3% of its genome. In contrast, *Ginkgo biloba* (Gymnosperm) presents the most modest representation at a mere 0.23%.

Interestingly, the scenario shifts when examining the Chlorophyte lineage. Here, in contrast to the other lineages, the proportion of Simple Repeats tends to be more pronounced. Within this lineage, several species exhibit notable percentages of Simple Repeats, such as *Monoraphidium minutum* at 11.71%, *Raphidocelis subcapitata* at 10.38%, *Chlamydomonas sp. UWO 241* at 9.33%, and *Micromonas pusilla CCMP1545* at 8.23%. The least representation within this lineage is observed in *Ostreobium quekettii* at 0.52%.

The presented boxplots delineate the %Divergence metrics for these elements.



Across the lineages of Angiosperms, Hornwort, Gymnosperms, Lycopods, and Liverworts, a notably stable pattern emerges in the boxplot. The overall median %Divergence for Simple Repeats in these lineages tends to be low, ranging approximately from 5 to 15. *Ginkgo biloba*, representing the Gymnosperms, stands out with the highest median value among these lineages. Furthermore, the IQR values for these plants remain consistently within the range of 15-20, indicating a relatively uniform distribution of %Divergence.

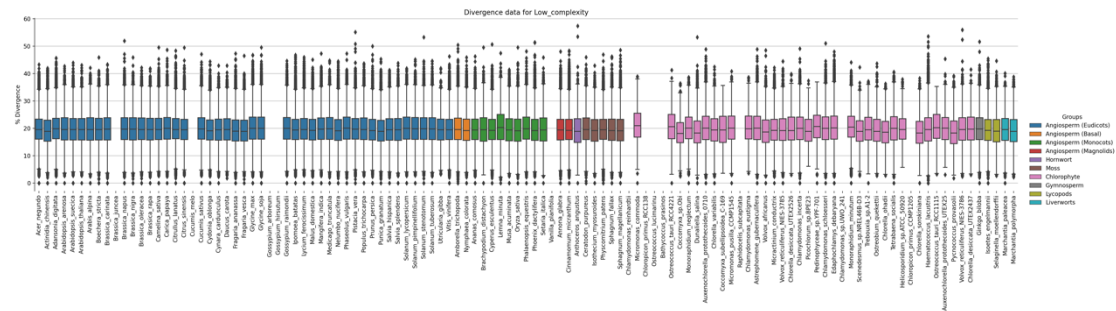
Contrastingly, the Chlorophyte lineage exhibits heightened variability in Simple Repeats. The median %Divergence values for these species oscillate between 10 and 20, with IQR typically falling within the 10-15 range, reflecting a broader spectrum of %Divergence values compared to other lineages.

Across all lineages, the whisker lengths display a common lower bound of 0. However, the upper bounds span a wider range, typically between 35 and 45. Many plants across these lineages manifest numerous outliers beyond the upper limit. An exception is *Dunaliella salina* (Chlorophyte), characterized by whisker lengths ranging from 0 to 25. This species presents a notably low median value of around 3, coupled with a narrow IQR of approximately 10, and remarkably exhibits the highest count of outliers.

Across the spectrum of plant lineages examined, the presence of Low Complexity Elements is pervasive, with their representation being notably consistent across different species. Despite their ubiquitous presence, the proportion occupied by LCEs in these genomes remains relatively modest. The highest proportion observed is in *Edaphochlamys debaryana* (Chlorophyte), where LCEs account for 1.18% of the sequence.

In stark contrast, all the rest of the plant species, spanning various lineages, exhibit LCE proportions below the 1% threshold. Notably, *Adiantum_capillus* (Fern) showcases the most minimal representation, with LCEs accounting for a mere 0.03% of its genome. It is noteworthy to mention that several species, including *Ginkgo_biloba* and *Gnetum_montanum* (both Gymnosperms), as well as *Scenedesmus sp. NREL 46B-D3* (Chlorophyte), demonstrate a consistent low representation of LCEs, each registering at 0.04%.

The presented boxplots delineate the %Divergence metrics for these elements.



Across all examined lineages, a striking consistency emerges in the %Divergence values associated with Low Complexity elements. Regardless of the lineage, the median %Divergence hovers around the 20 mark, showcasing an interesting uniformity without any deviations. This uniform median suggests that, despite the evolutionary differences and varied genomic contexts across these lineages, there might be consistent underlying processes or constraints governing the evolution or retention of low-complexity elements.

Additionally, the IQR values, ranging approximately from 16 to 23, further highlight this stability. Such a narrow IQR indicates that the majority of plants within these lineages share a similar %Divergence distribution, with minimal variability. The whisker lengths, spanning from 5 to 35, suggest that while there is some range in the data, most plants tend to fall within this specified range. However, the presence of outliers at both ends of the whiskers across the majority of these lineages might indicate specific evolutionary pressures or events affecting a subset of plants, leading to extreme %Divergence values.

2.2 Heatmap Representation



The heatmap delineates the abundance distribution of transposon elements across seven distinct plant lineages. The horizontal axis maps approximately fifty different transposon elements, while the vertical axis spans seven plant lineages: Angiosperms, Liverworts, Hornwort, Moss, Gymnosperms, Lycopods, and Chlorophyte.

Angiosperms: Within the Angiosperms, particularly among the Eudicots, elements such as Copia, L1, Helitron, CMC-EnSpm, and hAT-Tip100 consistently exhibit elevated abundance levels. In contrast, within the Basal Angiosperms, the element 'mixture' demonstrates exceptionally high abundance. Among the Monocots, while elements like Penelope and ERV1 manifest notably high abundance in certain species, others maintain a relatively consistent low overall TE abundance. Within the Magnolids, the element Merlin stands out with its remarkably high abundance.

Liverworts: Notably, Tad1 and hAT-hobo display heightened abundance levels. Such an increase might signify lineage-specific amplification events, which can contribute to genetic diversification or response to environmental pressures within Liverworts.

Moss: Elements like PiggyBac, TcMar-Stowaway, and Sola2 exhibit notably high abundances. Their elevated presence suggests lineage-specific amplification dynamics, potentially influencing the genomic architecture and evolutionary trajectories within Moss.

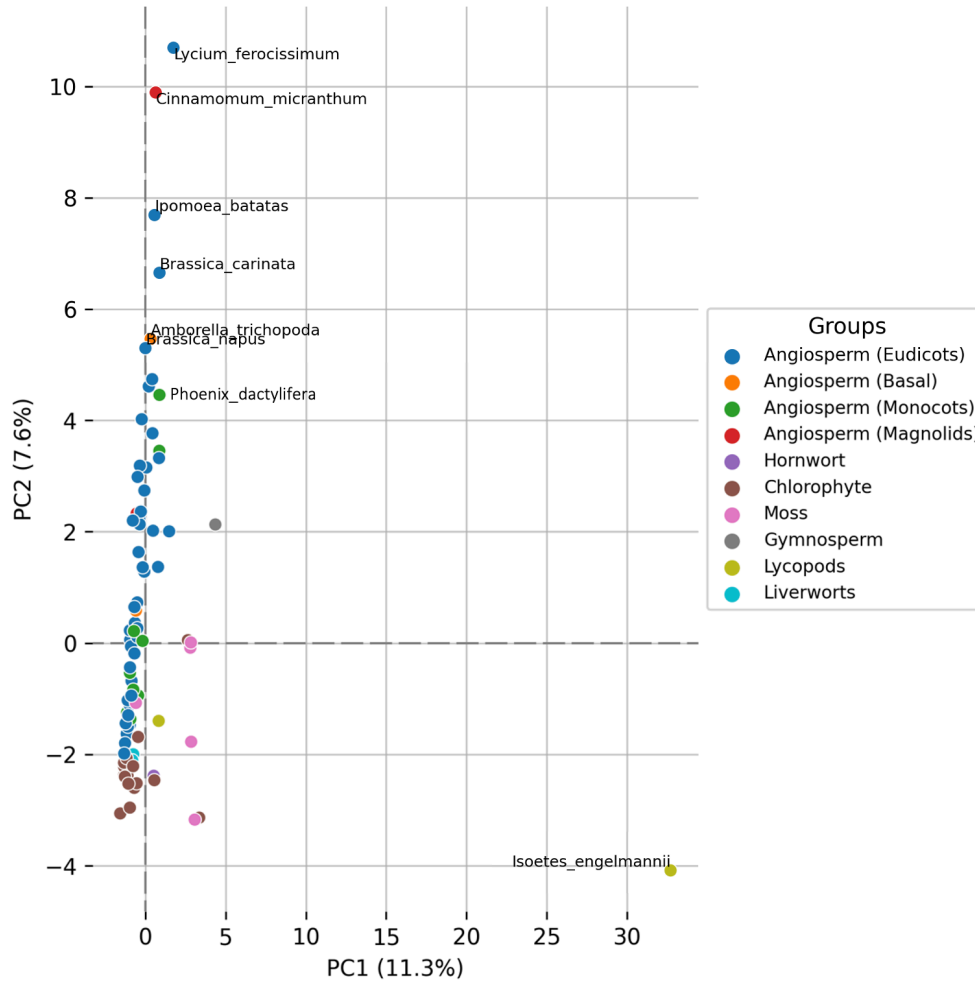
Gymnosperms: Within this lineage, Ngaro, Copia, and L1 display elevated abundance levels. LI-Tx1, in particular, showcases a dominant prevalence. The distinct transposon abundance profile in Gymnosperms, differing from patterns observed in Angiosperms, underscores unique lineage-specific dynamics, possibly reflecting evolutionary pressures or adaptive responses specific to Gymnosperms.

Lycopods: R1-LOA, MULE-NOF, Crypton-A, R2, Deceiver, and RTE-X demonstrate exceptionally high abundance. Such pronounced levels might suggest potential outliers or anomalies within this lineage, warranting further investigation into their functional implications and evolutionary significance.

Hornwort: In this lineage, Tc1 Mariner and L1-dep manifest notably elevated abundance. This observation hints at potential lineage-specific amplification events, highlighting the dynamic nature of transposon activity within Hornworts.

Chlorophyte: Relative to other lineages, Chlorophyte showcases reduced transposon element abundance. Nonetheless, even within this lineage, certain elements manifest significant prevalence, which might play roles in genomic stability or adaptive mechanisms specific to Chlorophyte species.

2.3 Principal Component Analysis



The PCA plot offers a multidimensional perspective on the distribution of transposon elements across diverse plant lineages. The plot is structured with the horizontal axis (PC1, representing 11.3% variance) and the vertical axis (PC2, representing 7.6% variance).

Distribution Patterns across Lineages:

Angiosperms: A majority of Angiosperms, particularly the Eudicots, are clustered around the central region of the plot, with values spanning from -2 to 6 along PC2. Notably, three Angiosperm (Eudicots) outliers extend to values of 7, 8, and 11 along PC2, hinting at unique genomic or evolutionary pressures. Basal Angiosperms occupy values around 1, 2, and 5 on PC2. The Monocots cluster closely around the central region, ranging from -1 to 5 along PC2. In contrast, a Magnolid species stands out with a distinctive position at the value of 10 along PC2.

Liverworts: Their tight clustering around (0, -2) underscores a consistent transposon profile, possibly reflecting evolutionary constraints or shared genomic characteristics within this lineage.

Moss: Moss species' spread between (2.5, -3) and (2.5, 1) suggests varied transposon dynamics. The positioning may reflect lineage-specific transposon amplification or suppression events, warranting further investigation.

Gymnosperms: Their positioning around (5, 2) indicates a unique transposon landscape distinct from Angiosperms. This might signify lineage-specific transposon evolution or genomic stabilization mechanisms in Gymnosperms.

Lycopods: The dual positioning of Lycopods is particularly intriguing. While one conforms to the expected range near (1, -1), the outlier at (33, -4) is of notable interest. This anomaly aligns with earlier heatmap observations, suggesting specific transposon elements or events driving this divergence.

Hornwort: Hornworts' predominant presence near (0, -3) suggests a shared transposon profile or evolutionary history. This clustering indicates potential conserved transposon dynamics or genomic stability within Hornworts.

Chlorophyte: Chlorophyte species, primarily clustering around (0, -3) to (0, -2), display a consistent transposon landscape. This uniformity may indicate genomic stability or shared evolutionary constraints across Chlorophyte species.

2.4 Estimating Diversity and Specificity

In the study, we embarked on an extensive re-annotation of the TE landscapes, going beyond mere identification to delve into the estimation of diversity parameters and specificity. By computing Shannon's Entropy (H_j), we were able to accurately gauge variations and intricacies within the TE landscapes. Furthermore, the assessment of specificity, determined through the specialization index (δ_j index) and the divergence relative to the entire TE landscape using Kullback–Leibler divergence (Div_j), provided profound insights into the behavior of these TEs.

The following two tables list the diversity and specificity across diverse plant species and of distinct transposable elements (TEs).

The first table presents an analysis of diversity and specificity across over a hundred different plant species. Notably, the highest diversity is observed in *Phoenix dactylifera* (Angiosperm (Monocots)) with a value of 2.61. This species also exhibits the lowest specificity at 4.29. Conversely, *Chlorella ohadii* (Chlorophyte) demonstrates the highest specificity of 6.25 but has the lowest diversity at 0.66.

A closer examination reveals a pattern within various plant lineages. Several species within the Angiosperm lineage, such as *Arabidopsis suecica* (diversity: 2.29, specificity: 4.62) and *Brassica carinata* (diversity: 2.41, specificity: 4.49), show relatively high diversity with moderately low specificity. In contrast, species like *Anthoceros angustus* from Hornworts (diversity: 2.00, specificity: 4.91) and *Marchantia polymorpha* from Liverworts (diversity: 2.32, specificity: 4.59) present higher specificity with decent diversity. Additionally, within the Lycopods, *Selaginella moellendorffii* demonstrates a diversity of 2.57 and a specificity of 4.34.

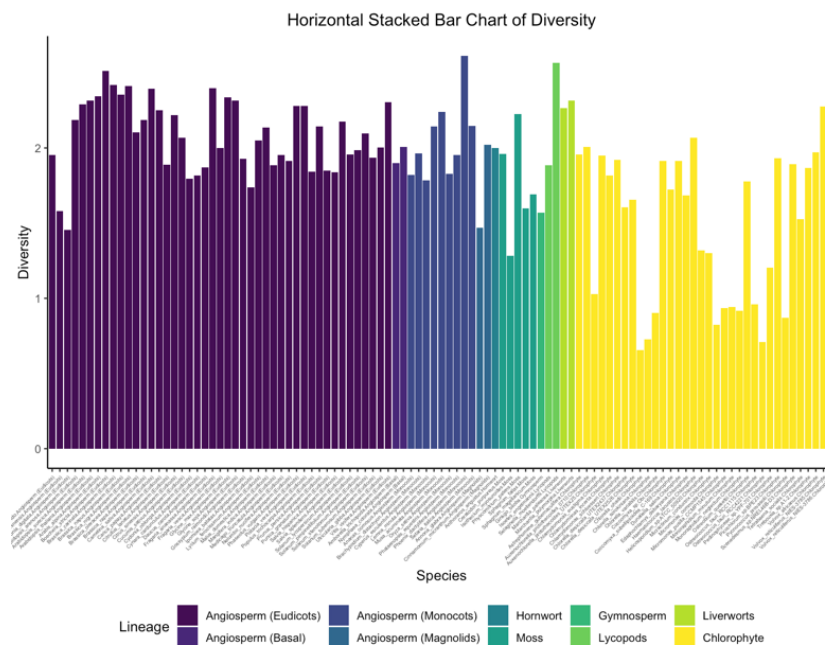
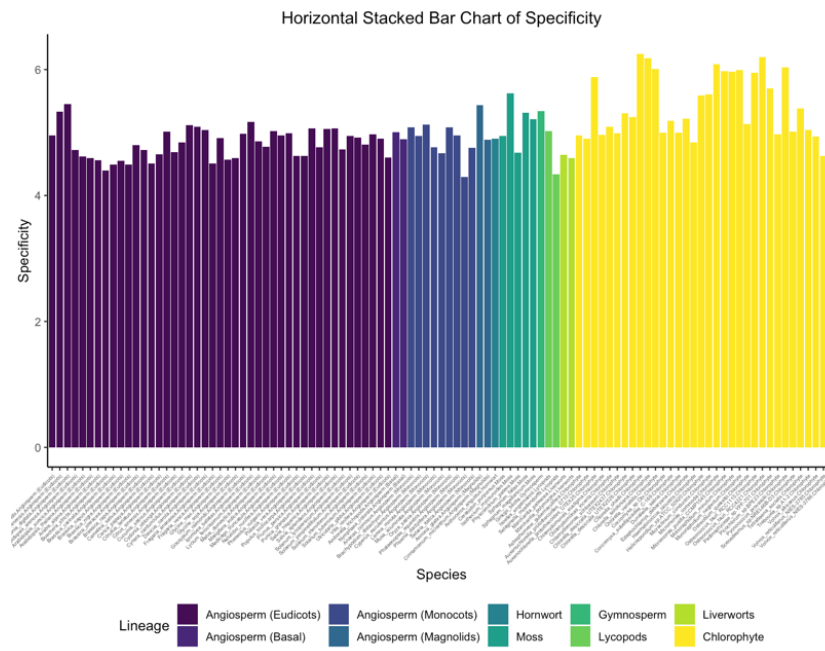
However, certain lineages display lower diversity and higher specificity. For instance, *Sphagnum fallax* from Moss has a diversity of 1.60 and a specificity of 5.31. Similarly, *Ginkgo biloba* from Gymnosperms shows a diversity of 1.60 and a specificity of 5.34. A unique observation is made with *Ostreococcus tauri* RCC4221 (Chlorophyte) having a diversity of 0.71 but a remarkably high specificity of 6.20. Notably, *Volvox reticuliferus* NIES-3785 from Chlorophyte exhibits a diversity of 2.28 and a specificity of 5.31.

The second table offers insights into the diversity and specificity of 120 different TE superfamily elements. A striking observation is the presence of 20 TEs with a diversity value of 0, including Kolobok, Ginger-1, tRNA-L1, tRNA-CR1, TATE, ERV-Foamy, PIF-ISL2EU, subtelo, hAT-hAT19, Alu, R2, TcMar-Tc4, Deceiver, Hydra, hAT-hATx, CMC-Transib, tRNA-Core-RTEERV4, TcMar-Ant1, Crypton-H, ARTEFACT. These TEs are predominantly found in specific plant species, leading to their elevated specificity values, peaking at 6.69.

In contrast, the TE with the highest diversity is Simple repeats at 6.34, being ubiquitously present across all analyzed plant species. This widespread distribution translates to its notably low specificity of 0.34. Other TEs like Unknown (diversity: 5.88, specificity: 0.81), Copia (diversity: 5.64, specificity: 1.05), Gypsy (diversity: 5.35, specificity: 1.34), and Low complexity (diversity: 6.23, specificity: 0.46) also exhibit high diversity but with reduced specificity values, suggesting their prevalent presence across various plant genomes.

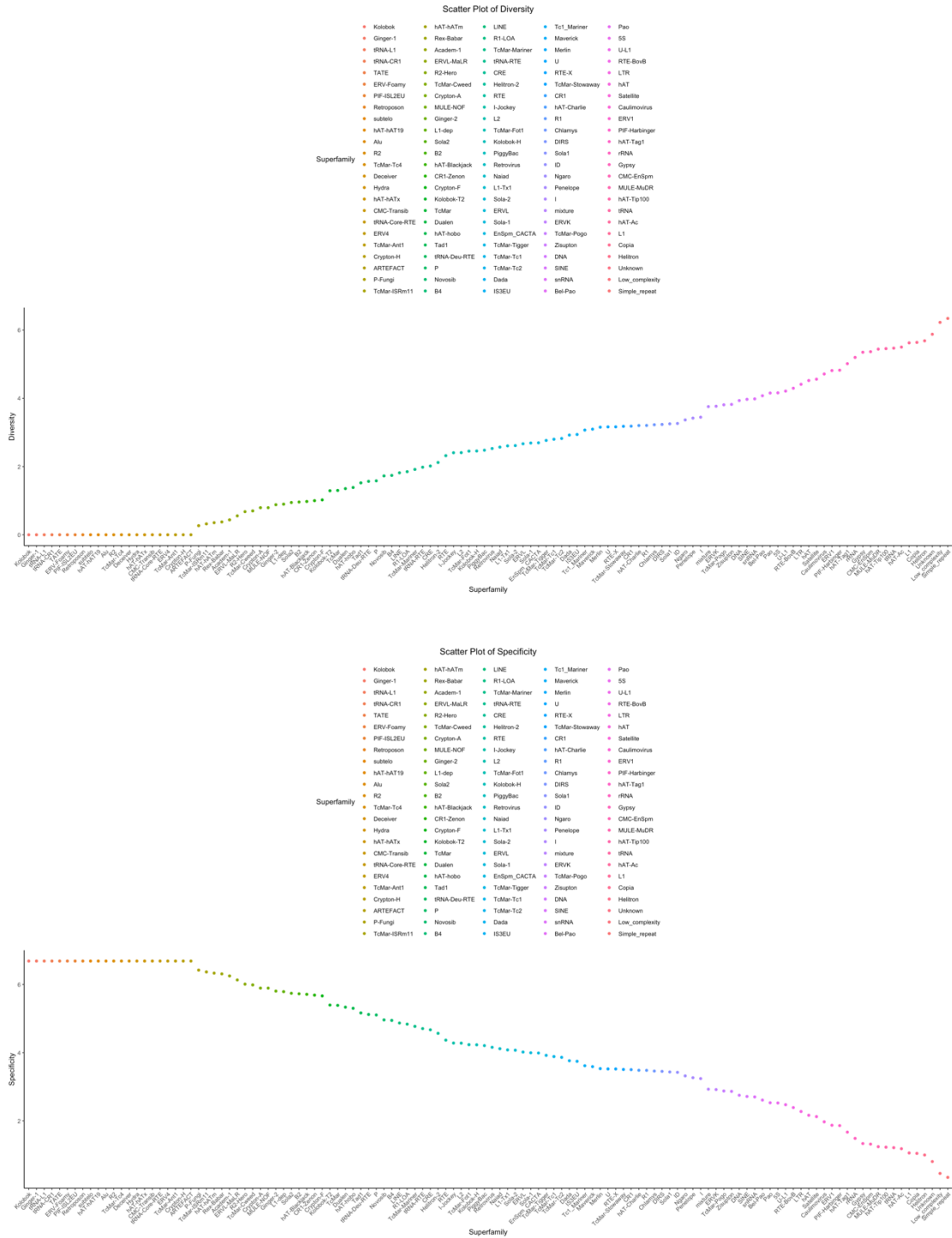
Visualizing Diversity and Specificity in Plant Species and TE Superfamilies

To provide a more intuitive understanding, horizontal stacked bar charts were utilized to visualize the diversity and specificity of plant species.



The horizontal stacked bar charts reveal distinct patterns among different plant lineages. Notably, Angiosperms exhibit significantly higher overall diversity compared to Chlorophytes. Additionally, Liverworts, Lycopods, and Hornworts tend to showcase relatively high diversity with lower specificity. In contrast, Gymnosperms demonstrate moderate levels of both diversity and specificity. Moss, on the other hand, exhibits lower diversity but higher specificity. These observations collectively highlight the nuanced diversity and specificity profiles across different plant lineages.

Subsequently, scatter plots were generated to visually represent the diversity and specificity patterns across various TE superfamilies.



As evident from the depicted figures, several TE superfamilies, including Kolobok, Ginger-1, tRNA-L1, tRNA-CR1, TATE, ERV-Foamy, PIF-ISL2EU, subtelo, hAT-hAT19, Alu, R2, TcMar-Tc4, Deceiver, Hydra, hAT-hATx, CMC-Transib, tRNA-Core-RTE, ERV4, TcMar-Ant1, Crypton-H, and ARTEFACT, exhibit a diversity value of 0, concurrently demonstrating the highest specificity.

On the contrary, simple repeats, low complexity, unknown, helitron, copia, L1, and others showcase remarkably high diversity, coupled with notably low specificity, implying a widespread and less constrained presence across the genome.

Discussion

1. The Unique TE Landscape Across Divergent Lineages

The intricate interplay between transposable elements and their host genomes has long been a subject of interest in evolutionary biology. The analysis combining both heatmap and PCA visualizations, provides a nuanced understanding of TE landscape variations across seven distinct plant lineages.

The heatmap delineates clear lineage-specific abundance distributions for various transposon elements. For instance, the pervasive abundance of Copia, L1, Helitron, CMC-EnSpm, and hAT-Tip100 in the majority of Angiosperms underscores the potential for active transposition events within this lineage. The last table further highlights the notably high diversity values associated with these elements. Such patterns may arise from a combination of environmental pressures and inherent genomic factors, leading to genomic plasticity and adaptability in Angiosperms.

As the heatmap indicates an increase in the content of specific transposable elements, there is a potential corresponding tendency for elevated values along the PC2 among various angiosperm species. Notably, *Lycium ferocissimum* exemplifies this trend, displaying the highest PC2 value, which aligns with its substantial content of several transposon elements, including L1 and Copia. Likewise, in the case of *Ipomoea batatas*, heightened levels of Merlin, I, and Alu elements are apparent, correlating with a noteworthy PC2 value of approximately 8. Moreover, within monocots, *Phoenix dactylifera* distinguishes itself with the highest PC2 value, showcasing a significant abundance of Penelope and ERV1 transposable elements. This observation suggests a plausible association between the increase in specific TE content in the heatmap and the resulting variations in PC2 values across diverse angiosperm taxa.

Interestingly, both Moss and Liverworts exhibit notable lineage-specific amplification dynamics. Elements like PiggyBac, TcMar-Stowaway, and Sola2 in Moss, and Tad1 and hAT-hobo in Liverworts, respectively, highlight the dynamic nature of TE evolution within these lineages. Such dynamics might be pivotal in shaping the genomic architecture and influencing adaptive responses to environmental cues.

The distinct clustering of Gymnosperms around specific PCA coordinates emphasizes their unique TE landscape, diverging from the patterns observed in Angiosperms. Such divergence might be indicative of lineage-specific evolutionary pressures or adaptive responses, underscoring the need for detailed genomic studies in Gymnosperms.

A particularly noteworthy observation in the analysis pertains to the Lycopods lineage. In the PCA plot, an outlier was identified corresponding to the Lycopods species, *Isoetes engelmannii*, with coordinates (33, -4). Upon examination, this species

exhibited a Diversity value of 1.88 and a Specificity value of 5.02, potentially reflecting specialized evolutionary dynamics or distinct genomic features within this species. The heatmap results were subsequently checked, revealing that *Isoetes engelmannii* exhibited a marked elevation in the abundances of several transposon elements, including R1-LOA, MULE-NOF, Crypton-A, R2, Deceiver, and RTE-X. Examination of the data reveals that the Diversity values for R2 and Deceiver are 0, indicative of exceptionally high specificity, implying that these transposon elements exist exclusively in *Isoetes engelmannii*. Given the significant abundance of these transposons in *Isoetes engelmannii*, it is plausible to infer that the pronounced anomaly in the PCA plot might be attributed to these exceptionally high transposon levels. The presence and abundance of these elements in *Isoetes engelmannii* warrant further investigation to elucidate their potential functional implications and evolutionary significance within the Lycopods lineage.

Furthermore, Chlorophytes showcase unique features. In the PCA plot, Chlorophytes primarily cluster around -3 to -2 along PC2. Meanwhile, in the heatmap, Chlorophytes exhibit significantly diminished TE content.

The findings align with the hypothesis that the most divergent lineages indeed harbor the most unique TE landscapes. The variations in TE abundance, distribution, and dynamics across different plant lineages reflect their respective evolutionary trajectories, adaptive potentials, and genomic architectures.

Contextualizing with Existing Literature:

Comparing our results with the findings from the study titled "PlantLTRdb: An interactive database for 195 plant species LTR-retrotransposons" (Mokhtar et al. 2023) underscores the breadth and depth of TE research in plants. The PlantLTRdb database, with its comprehensive coverage of 195 plant species, has highlighted that LTR-RT activity stands as a pivotal factor driving genome evolution. Specifically, the activation and proliferation of LTR-RTs can lead to significant expansions in genome size, with certain plant genomes having over 70% of their content attributed to these elements. Furthermore, the aforementioned study delineates that there exists a strong positive correlation between the total length of LTR-RTs and genome size across these plant species.

A critical point of intersection between our study and the findings from PlantLTRdb lies in the functional impacts of LTR-RTs on genes. The study emphasizes that LTR-RTs, when situated near or within genes, have the potential to directly modify gene function. Additionally, the research conducted a detailed statistical analysis of LTR-RT lengths and their respective superfamily, such as Copia and Gypsy elements. This information provides a more nuanced understanding of the structural variations and complexities within LTR-RTs across plant genomes.

Another significant dimension highlighted by the study is the insertion age of the plant species studied, which reflects the evolutionary rate associated with the uniqueness of their genomic content. This temporal aspect adds a layer of complexity to our understanding of TE dynamics, suggesting that the genomic impacts of LTR-RTs may vary depending on the evolutionary histories of individual plant lineages.

Moreover, the research states that environmental stresses can serve as triggers for TE activation. This perspective resonates with our discussion on the potential adaptive responses of plants to their environments through TE dynamics. It's conceivable that certain TE activations within specific plant lineages could be a direct response to environmental pressures, thereby contributing to the observed genomic plasticity and adaptability in these lineages.

While PlantLTRdb offers a holistic, database-driven perspective on TE dynamics, our research provides a granular analysis of lineage-specific dynamics. Together, these insights paint a comprehensive picture of TE evolution across plant taxa, emphasizing the intricate roles TEs play in plant evolution, genomic diversification, and environmental adaptation.

In conclusion, the synthesis of heatmap, PCA analyses, and insights from PlantLTRdb enhances our understanding of TE dynamics in various plant lineages. Mokhtar et al.'s study enriches this discourse, emphasizing the multifaceted impacts of LTR-RTs on genome evolution, gene function, environmental adaptation, and evolutionary rates reflected in insertion ages. This underscores the necessity for a comprehensive approach in unraveling complex interactions between TEs and plant genomes.

2. Genome Size and Its Influence on TE Landscape

The hypothesis 2 posited a correlation between genome size and the TE landscape across plant species, suggesting that larger genomes may provide more "space" for TE insertions, influencing gene density and spatial organization, ultimately affecting TE accumulation and distribution patterns.

The scatter plot reinforces this hypothesis, revealing distinct patterns across various plant lineages. Gymnosperms, characterized by extensive genomes, particularly exemplified by *Cycas panzihuaensis*, display the largest Genome Size and the highest content of Bases Masked, exceeding 80%. Conversely, Chlorophytes, with smaller genomes, exemplified by *Helicosporidium sp. ATCC 50920*, exhibits the smallest Genome Size and the lowest Bases Masked content.

Despite lineage variations, the scatter plot consistently illustrates a positive correlation between Genome Size and Bases Masked. The grey reference line serves as a visual guide, emphasizing the clustering of data points around the line. This consistent

relationship highlights that, irrespective of lineage, larger genomes tend to correlate with higher TE content, aligning with our hypothesis.

In conclusion, the scatter plot effectively supports the hypothesis, providing visual evidence for the positive correlation between Genome Size and Bases Masked content across diverse plant lineages.

Given these observations, the comprehensive study by Pedro DLF et al. provides an essential context, emphasizing an association between genome size and TE content across diverse plant lineages. While their exhaustive analysis across 67 plant genomes reveals a general trend of larger genomes having higher occurrences of TEs, our study contributes by examining a broader range of samples. This expanded dataset further underscores the complexity and variability within plant genomes. The insights from Pedro DLF et al. regarding the pervasive nature of TE landscapes, combined with our extensive sampling, collectively emphasize the intricate interplay between genome size and TE content, highlighting its multifaceted role in shaping genomic evolution across plant species.

3. TE Activity and Variability in Rapidly Evolving Lineages

The hypothesis put forth in this research suggests that plant lineages undergoing rapid evolution, particularly those with elevated %Divergence in LTR elements, might manifest increased TE activity or mutation rates. Such a hypothesis is rooted in the premise that heightened evolutionary dynamics, as indicated by increased %Divergence in LTR elements, could be accompanied by escalated TE mobility or activity. This heightened activity might serve as a genomic response mechanism, aiding these rapidly evolving lineages in adapting to a myriad of environmental challenges.

Upon delving into the data, certain findings align with this hypothesis.

Starting with the LTR elements, a clear distinction between Chlorophytes and other plant lineages emerged. Chlorophyte species exhibited a broader range of LTR element %Divergence, predominantly clustering below 15%. This pattern suggests a distinct evolutionary trajectory within this ancestral group. Conversely, advanced plant lineages, including Angiosperms, Hornwort, Moss, Gymnosperms, Lycopods, and Liverworts, consistently demonstrated higher median %Divergence values, hovering around 15-20%. Notably, *Ginkgo biloba*, a Gymnosperm, stood out with a median value reaching 25%, underscoring elevated TE activity or mutation rates within this advanced lineage.

Similar patterns were observed for Copia and Gypsy elements. While most advanced plant lineages displayed median values between 15-20%, Chlorophytes exhibited greater variability, with median values predominantly below 15%.

Analysis of LINE and SINE elements reinforced these patterns. Advanced plant lineages consistently exhibited higher median %Divergence values, particularly within Angiosperms. In contrast, Chlorophyte species demonstrated lower median values, reinforcing the distinct evolutionary dynamics associated with these ancestral lineages. These findings align with our hypothesis, suggesting that rapidly evolving plant lineages may indeed exhibit increased TE diversity and variability.

The Unclassified Elements present an interesting pattern across various plant lineages. Angiosperms, Hornwort, Moss, Lycopods, and Liverworts consistently manifest median %Divergence values predominantly between 15 and 20. Notably, Gymnosperms deviate slightly with a median %Divergence peaking at 23. In stark contrast, the Chlorophyte lineage stands apart, reflecting a more subdued median %Divergence that spans between 5 and 15. These observations further emphasize the heightened TE activity or mutation rates in advanced plant lineages compared to their more primitive counterparts.

The analysis of Simple Repeat Elements presented a contrasting scenario. Advanced plant lineages, including Angiosperms, Hornwort, Moss, Lycopods, and Liverworts, exhibited a stable pattern with median %Divergence values ranging from 5 to 15%. In contrast, Chlorophyte species displayed higher variability, with median values fluctuating between 10 to 20%. This divergence in patterns suggests that while some TE elements may be influenced by evolutionary rates, others might be governed by different evolutionary constraints or mechanisms.

Lastly, the Low Complexity Elements showcased consistent median values across all lineages, hovering around 20%. This uniformity suggests that whether situated at centromeric regions or telomeres, the properties of these components remain nearly identical. Such consistency may arise from the similar functional or structural roles that LCEs play in plant evolution.

In summation, the findings of this study shed light on the nuanced dynamics of TEs across plant lineages. Advanced, rapidly evolving lineages, like Angiosperms and Gymnosperms, seem to harbor increased TE diversity, hinting at potential adaptive evolutionary strategies. Conversely, ancestral lineages maintain a more stable TE landscape, underscoring the delicate balance between TE activity and genomic stability. These insights further underscore the multifaceted roles of TEs in plant genome evolution and the adaptive strategies employed by different lineages.

Conclusions

The exploration of transposable elements (TEs) across a spectrum of plant lineages, from Chlorophytes to Angiosperms, offers profound insights into the evolutionary dynamics and functional implications of these genomic entities. TEs, as dynamic components, play pivotal roles in sculpting the genomic architecture, influencing gene regulation, and potentially driving evolutionary innovations.

Lineage-Specific TE Dynamics: Our study underscores the intricate and unique TE landscapes inherent to distinct plant lineages. Advanced lineages, notably Angiosperms and Gymnosperms, exhibit pronounced TE variability, potentially reflecting adaptive evolutionary strategies. In contrast, ancestral lineages like Chlorophytes maintain a more conserved TE profile, highlighting the interplay between evolutionary stability and TE activity.

Genome Size and TE Interplay: The scatter plot unequivocally confirms a positive correlation between genome size and TE content across diverse plant lineages. This correlation underscores the important role of genomic expansiveness in facilitating TE insertions.

Rapidly Evolving Lineages and TE Activity: Our findings substantiate the hypothesis that rapidly evolving lineages may harbor increased TE activity or mutation rates. Elevated %Divergence values in advanced lineages, coupled with distinct TE profiles, suggest an evolutionary response mechanism, enabling these lineages to navigate diverse environmental challenges.

Functional and Evolutionary Implications: Beyond mere genomic components, TEs hold functional significance, potentially influencing specific genomic features or evolutionary traits. Their pervasive presence across diverse lineages underscores their indispensable role in shaping plant genome architecture and function.

In essence, this research illuminates the multifaceted roles of TEs in plant evolution, genome dynamics, and adaptation. As genomic studies continue to unravel the complexities of plant genomes, TEs remain central players, necessitating continued exploration to decipher their full spectrum of functions and evolutionary impacts.

References

- Bourque, G., Burns, K.H., Gehring, M. *et al.* Ten things you should know about transposable elements. *Genome Biol* **19**, 199 (2018). <https://doi.org/10.1186/s13059-018-1577-z>
- Ramakrishnan, M., Satish, L., Sharma, A. *et al.* Transposable elements in plants: Recent advancements, tools and prospects. *Plant Mol Biol Rep* **40**, 628–645 (2022). <https://doi.org/10.1007/s11105-022-01342-w>
- Muñoz-López M, García-Pérez JL (2010) DNA transposons: nature and applications in genomics. *Curr Genom* 11:115–128
- Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, Flavell A, Leroy P, Morgante M, Panaud O, Paux E, SanMiguel P, Schulman AH (2007) A unified classification system for eukaryotic transposable elements. *Nat Rev Genet* 8(12):973–982. <https://doi.org/10.1038/nrg2165>
- Brown PO, Bowerman B, Varmus HE, Bishop JM. Correct integration of retroviral DNA in vitro. *Cell*. 1987;49:347–56.
- Luan DD, Korman MH, Jakubczak JL, Eickbush TH. Reverse transcription of R2Bm RNA is primed by a nick at the chromosomal target site: a mechanism for non-LTR retrotransposition. *Cell*. 1993;72:595–605.
- Malik HS, Eickbush TH. Phylogenetic analysis of ribonuclease H domains suggests a late, chimeric origin of LTR retrotransposable elements and retroviruses. *Genome Res*. 2001;11:1187–97.
- Mhiri, C., Borges, F. and Grandbastien, M.A., 2022. Specificities and dynamics of transposable elements in land plants. *Biology*, 11(4), p.488. <https://doi.org/10.3390/biology11040488>
- Schubert I, Vu GTH. Genome stability and evolution: attempting a holistic view. *Trends Plant Sci*. 2016;21:749–57.
- Haley, A.L. and Mueller, R.L., 2022. Transposable element diversity remains high in gigantic genomes. *Journal of Molecular Evolution*, 90 (5), pp.332-341. <https://doi.org/10.1007/s00239-022-10063-3>
- Almeida, M.V., Vernaz, G., Putman, A.L. and Miska, E.A., 2022. Taming transposable elements in vertebrates: from epigenetic silencing to domestication. *Trends in Genetics*, 38(6), pp.529-553 <https://doi.org/10.1016/j.tig.2022.02.009>

Alla Mikheenko, Andrey Prjibelski, Vladislav Saveliev, Dmitry Antipov, Alexey Gurevich,

Versatile genome assembly evaluation with QUASt-LG,

Bioinformatics (2018) 34 (13): i142-i150. doi: 10.1093/bioinformatics/bty266

Jacques Dainat, Darío Hereñú, Dr. K. D. Murray, Ed Davis, Kathryn Crouch, LucileSol, Nuno Agostinho, pascal-git, Zachary Zollman, & tayyrov. (2023). NBISweden/AGAT: AGAT-v1.2.0 (v1.2.0). Zenodo. <https://doi.org/10.5281/zenodo.8178877>

Flynn, J. M., Hubley, R., Goubert, C., Rosen, J., Clark, A. G., Feschotte, C., & Smit, A. F. (2020). RepeatModeler2 for automated genomic discovery of transposable element families. *Proceedings of the National Academy of Sciences*, 117(17), 9451-9457. <https://doi.org/10.1073/pnas.1921046117>

A.F.A. Smit, R. Hubley & P. Green. RepeatMasker at <http://repeatmasker.org>.

Zhang, R.-G., Li, G.-Y., Wang, X.-L., Dainat, J., Wang, Z.-X., Ou, S., & Ma, Y. (2022). TEsorter: An accurate and fast method to classify LTR-retrotransposons in plant genomes. *Horticulture Research*, 9, uhac017. <https://doi.org/10.1093/hr/uhac017>

Mokhtar, M. M., Alsamman, M., & El Allali, A. (2023). PlantLTRdb: An interactive database for 195 plant species LTR-retrotransposons. *Frontiers in Plant Science*, 14, Sec. Plant Systematics and Evolution. <https://doi.org/10.3389/fpls.2023.1134627>

Pedro, D. L. F., Amorim, T. S., Varani, A., Guyot, R., Domingues, D. S., & Paschoal, A. R. (2021). An Atlas of Plant Transposable Elements. *F1000Research*, 10, 1194. <https://doi.org/10.12688/f1000research.74524.1>

Hadley Wickham, Winston Chang, Lionel Henry, Thomas Lin Pedersen, Kohske Takahashi, Claus Wilke, and Kara Woo. (2022). ggplot2: Create Elegant Data Visualisations Using the Grammar of Graphics (Version 3.4.4) [Software]. <https://CRAN.R-project.org/package=ggplot2>