



HAL
open science

Sequencing, de novo assembly of *Ludwigia* plastomes, and comparative analysis within the Onagraceae family

F Barloy-Hubler, A.-L Le Gac, Christophe C. Boury, Erwan Guichoux, D Barloy

► To cite this version:

F Barloy-Hubler, A.-L Le Gac, Christophe C. Boury, Erwan Guichoux, D Barloy. Sequencing, de novo assembly of *Ludwigia* plastomes, and comparative analysis within the Onagraceae family. *Peer Community In Genomics*, 2025, e43, <10.24072/pci.genomics.100334>. <hal-05068858v2>

HAL Id: hal-05068858

<https://hal.science/hal-05068858v2>

Submitted on 15 May 2025

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



Distributed under a Creative Commons CC BY 4.0 - Attribution - International License

Research article

Published
2025-04-23

Cite as

F. Barloy-Hubler, A.-L. Le Gac, C. Boury, E. Guichoux and D. Barloy (2025) *Sequencing, de novo assembly of Ludwigia plastomes, and comparative analysis within the Onagraceae family*, Peer Community Journal, 5: e43.

Correspondence

dominique.barloy@agrocampus-ouest.fr

Peer-review

Peer reviewed and
recommended by
PCI Genomics,

<https://doi.org/10.24072/pci.genomics.100334>



This article is licensed
under the Creative Commons
Attribution 4.0 License.

Sequencing, de novo assembly of *Ludwigia* plastomes, and comparative analysis within the Onagraceae family

F. Barloy-Hubler¹, A.-L. Le Gac², C. Boury³, E. Guichoux³, and D. Barloy⁴

Volume 5 (2025), article e43

<https://doi.org/10.24072/pcjournal.536>

Abstract

The Onagraceae family, which belongs to the order Myrtales, consists of approximately 657 species and 17 genera. This family includes the genus *Ludwigia* L., which is comprised of 82 species. In this study, we focused on the two aquatic invasive species *Ludwigia grandiflora* subsp. *hexapetala* (*Lgh*) and *Ludwigia peploides* subsp. *montevidensis* (*Lpm*) largely distributed in aquatic environments in North America and in Europe. Both species have been found to degrade major watersheds leading ecological and economical damages. Genomic resources for Onagraceae are limited, with only *Ludwigia octovalvis* (*Lo*) plastid genome available for the genus *Ludwigia* L. at the time of our study. This scarcity constrains phylogenetic, population genetics, and genomic studies. To brush up genomic resources, new complete plastid genomes of *Ludwigia grandiflora* subsp. *hexapetala* (*Lgh*) and *Ludwigia peploides* subsp. *montevidensis* (*Lpm*) were generated using a combination of MiSeq (Illumina) and GridION (Oxford Nanopore) sequencing technologies. These plastomes were then compared to the published *Ludwigia octovalvis* (*Lo*) plastid genome, which was re-annotated by the authors. We initially sequenced and assembled the chloroplast (cp) genomes of *Lpm* and *Lgh* using a hybrid strategy combining short and long reads sequences. We observed the existence of two *Lgh* haplotypes and two potential *Lpm* haplotypes. *Lgh*, *Lpm*, and *Lo* plastomes were similar in terms of genome size (around 159 Kb), gene number, structure, and inverted repeat (IR) boundaries, comparable to other species in the Myrtales order. A total of 45 to 65 SSRs (*simple sequence repeats*), were detected, depending on the species, with the majority consisting solely of A and T, which is common among angiosperms. Four chloroplast genes (*matK*, *accD*, *ycf2* and *ccsA*) were found under positive selection pressure, which is commonly associated with plant development, and especially in aquatic plants such as *Lgh*, and *Lpm*. Our hybrid sequencing approach revealed the presence of two *Lgh* plastome haplotypes which will help to advance phylogenetic and evolutionary studies, not only specifically for *Ludwigia*, but also the Onagraceae family and Myrtales order. To enhance the robustness of our findings, a larger dataset of chloroplast genomes would be beneficial.

¹CNRS, UMR 6553 ECOBIO, Université de Rennes, Rennes 35000, France, ²Institut Curie, 7500, Paris, France, ³Université de Bordeaux, INRAE, BIOGECO, Cestas, France, ⁴DECOD (Ecosystem Dynamics and Sustainability), Institut Agro, IFREMER, INRAE, Rennes, France

Introduction

The Onagraceae family belongs to the order Myrtales which includes approximately 657 species of herbs, shrubs, and trees across 17 genera grouped into two subfamilies: subfam. Ludwigioideae W. L. Wagner and Hoch, which only has one genus (*Ludwigia* L.), and subfam. Onagroideae which contains six tribes and 21 genera (Wagner et al. 2007). *Ludwigia* L. is composed of 83 species (Levin et al. 2003, Levin et al. 2004). The current classification for *Ludwigia* L., which are composed of several hybrid and/or polyploid species, lists 23 sections. A recent molecular analysis is clarified and supported several major relationships in the genus but has challenged the complex sectional classification of *Ludwigia* L. (Liu et al. 2017).

The diploid species *Ludwigia peploides* (Kunth) Raven subsp. *montevidensis* (Spreng.) (Raven 1963) (named here *Lpm*) ($2n=16$), and the decaploid species, *Ludwigia grandiflora* (Michx.) Greuter & Burdet subsp. *hexapetala* (Hook. & Arn) Nesom & Kartesz (named here *Lgh*) ($2n=80$), reproduce essentially by clonal propagation, which suggests that there is a low genetic diversity within the species (Dandelot et al. 2005). *Lgh* and *Lpm* are native to South America and are considered as one of the most aggressive aquatic invasive plants (Reddy et al., 2021). Largely distributed in aquatic environments in North America and in Europe (Hussner et al. 2016), both species have been found to degrade major watersheds as well as aquatic and riparian ecosystems (Grewell et al. 2016) leading ecological and economical damages. In France, both species occupied aquatic habitats, such as static or slow-flowing waters, riversides, and have recently been observed in wet meadows (Lambert et al. 2010). The transition from an aquatic to a terrestrial habitat has led to the emergence of two *Lgh* morphotypes (Haury et al. 2014a). The appearance of metabolic and morphological adaptations could explain the ability to acclimatize to terrestrial conditions, and this phenotypic plasticity involves various genomic and epigenetic modifications (Billet et al. 2018).

Adequate genomic resources are necessary in order to be identify the genes and metabolic pathways involved in the adaptation process leading to plant invasion (Gioria et al. 2023) with genomic information making it possible to predict and control invasiveness (Moravcová et al. 2015). However, even though the number of terrestrial plant genomes has increased considerably over the last 20 years, only a small fraction ($\sim 0.16\%$) have been sequenced, with some clades being significantly more represented than others (Marks et al. 2021). Thus, for the Onagraceae family (which includes *Ludwigia* sp.), only a handful of chloroplast sequences (plastomes) are available, and the complete genome has not yet been sequenced. If *Lpm* is a diploid species ($2n=2x=16$) with a relatively small genome size (262 Mb), *Lgh* is a decaploid species ($2n=10x=80$) with a large size genome of 1419 Mb (Barloy et al. 2024). Obtaining a reference genome for these two non-model species without having a genome close to the *Ludwigia* species is challenging and development of plastome and/or mitogenome will be a first step to generate genomic resource. As of April 2023, there are 10,712 reference plastomes listed on GenBank (Release 255: April 15 2023), with the vast majority (10,392 genomes) belonging to Viridiplantae (green plants). However, in release 255, the number of plastomes available for the Onagraceae family is limited, with only 36 plastomes currently listed. Among these, 15 plastomes are from the tribe Epilobieae, with 11 in the *Epilobium* genus and 4 in the *Chamaenerion* genus. Additionally, there are 23 plastomes from the tribe Onagreae, with 17 in the *Oenothera* genus, 5 in the *Circaea* genus, and only one in the *Ludwigia* genus. The *Ludwigia octovalvis* chloroplast genome was released in 2016 as a unique haplotype of approximately 159 kb (Liu et al. 2016). *L. octovalvis* belongs to sect. *Macrocarpon* (Micheli) H.Hara while *Lpm* and *Lgh* belong to *Jussiaea* section (Zardini and Raven 1992, Hoch et al. 2015). Generally, the inheritance of chloroplast genomes is considered to be maternal in angiosperms. However, biparentally inherited chloroplast genomes could potentially exist in approximately 20% of angiosperm species (Hu et al. 2008, Zhang and Sodmergen 2010). Both maternal and biparental inheritance are described in the Onagraceae family. In tribe Onagreae, *Oenothera* subsect. *Oenothera* are known to have biparental plastid inheritance (Wagner et al. 2007, Jones and Cleland 1974). In tribe Epilobieae, biparental plastid inheritance was also

reported in *Epilobium* L. with mainly maternal transmission, and very low proportions of paternally transmitted chloroplasts (Schmitz and Kowallik 1986).

The chloroplast is the symbolic organelle of plants and plays a fundamental role in photosynthesis. Chloroplasts evolved from cyanobacteria through endosymbiosis and thereby inherited components of photosynthesis reactions (photosystems, electron transfer and ATP synthase) and gene expression systems (transcription and translation, Sato 2021). In general, chloroplast genomes (plastomes) are highly conserved in size, structure, and genetic content. They are rather small (120-170 kb, Gualberto et al. 2014), with a quadripartite structure comprising two long identical inverted repeats (IR, 10–30 kb) separated by large and a small single copy regions (LSC and SSC, respectively). They are also rich in genes, with around 100 unique genes encoding key proteins involved in photosynthesis, and a comprehensive set of ribosomal RNAs (rRNAs) and transfer RNAs (tRNAs, Tonti-Filippini et al. 2017). Plastomes are generally circular but linear shapes also exist (Oldenburg and Bendich 2016). Chloroplast DNA usually represents 5-20% of total DNA extracted from young leaves and therefore low-coverage whole genome sequencing can generate enough data to assemble an entire chloroplast genome (Twyford and Ness 2017).

If we refer to their GenBank records, more than 95% of these plastomes were sequenced by so-called short read techniques (mostly Illumina). However, in most seed plants, the plastid genome exhibits two large inverted repeat regions (60 to 335 kb, Twyford and Ness 2017), which are longer than the short read lengths (< 300 bp). This leads to incomplete or approximate assemblies (Wang et al. 2018). Recent long-read sequencing (> 1000 bp) provides compelling evidence that terrestrial plant plastomes exhibit two structural haplotypes. These haplotypes are present in equal proportions and differ in their inverted repeat (IR) orientation (Wang and Lanfear 2019). This shows the importance of using the so-called third generation sequence (TGS, PacBio or Nanopore) to correctly assemble the IRs of chloroplasts and to identify any different structural haplotypes. The current problem with PacBio or Nanopore long read sequencing is the higher error rate compared to short read technology (Ferrari et al. 2013, Jain et al. 2018, Rang et al. 2018). Thus, a hybrid strategy which combines long reads (to access the genomic structure) and short reads (to correct sequencing errors) could be effective (Wang et al. 2018, Scheunert et al. 2020).

Here, we report the newly sequenced complete plastid genomes of *Ludwigia grandiflora* subsp. *hexapetala* (*Lgh*) and *Ludwigia peploides* subsp. *montevidensis* (*Lpm*), using a combination of different sequencing technologies, as well as a re-annotated comparative genomic analysis of the published *Ludwigia octovalvis* (*Lo*) plastid. The main objectives of this study are (1) to assemble and annotate the plastomes of two new species of *Ludwigia* sp., (2) to reveal the divergent sequence hotspots of the plastomes in this genus and in the Onagraceae (3) to identify the genes under positive selection.

To achieve this, we utilized long read sequencing data from Oxford Nanopore and short read sequencing data from Illumina to assemble the *Lgh* plastomes and compared these assemblies with those obtained solely from long reads of *Lpm*. We also compared both plastomes to the published plastome of *Lo*. Our findings demonstrated the value of de novo assembly in reducing assembly errors and enabling accurate reconstruction of full heteroplasmy. We also evaluated the performance of a variety of software for sequence assembly and correction in order to define a workflow that will be used in the future to assemble *Ludwigia* sp. mitochondrial and nuclear genomes. Finally, the three new *Ludwigia* plastomes generated by our study make it possible to extend the phylogenetic study of the Onagraceae family and to compare it with previously published analyses (Liu et al. 2017, Bedoya and Madriñán 2015, Liu et al. 2020).

Material and Methods

Plant sampling and experimental design

The original plant materials were collected in June of 2018 near to Nantes (France) and formally identified by D. Barloy. *L. grandiflora* subsp. *hexapetala* (*Lgh*) plants were taken from the Mazerolles swamps (N47 23.260, W1 28.206), and *L. peploides* subsp. *montevidensis* (*Lpm*)

plants from La Musse (N 47.240926, W -1.788688). Plants were cultivated in a growth chamber in a mixture of $\frac{1}{3}$ soil, $\frac{1}{3}$ sand, $\frac{1}{3}$ loam with flush water level, at 22°C and a 16 h/8 h (light/dark) cycle. A single stem of 10 cm for each species was used for vegetative propagation in order to avoid potential genetic diversity. *De novo* shoots, taken three centimeters from the apex, were sampled for each species. Samples for gDNA extraction were pooled and immediately snap-frozen in liquid nitrogen, then lyophilized over 48 h using a Cosmos 20K freeze-dryer (Cryotec, Saint-Gély-du-Fesc, France) and stored at room temperature. All the plants were destroyed after being used as required by French authorities for invasive plants (article 3, prefectorial decree n°2018/SEE/2423).

Due to high polysaccharide content and polyphenols in *Lpm* and *Lgh* tissues and as no standard kit provided good DNA quality for sequencing, genomic DNA extraction was carried out using a modified version of the protocol proposed by Panova et al in 2016, with three purification steps (Panova et al. 2016).

40 mg of lyophilized buds were ground at 30 Hz for 60 s (Retsch MM200 mixer mill, Fisher). The ground tissues were lysed with 1 ml CF lysis buffer (Macherey-Nagel) supplemented with 20 µl RNase and incubated for 1 h at 65°C under agitation. 20 µl proteinase K was then added before another incubation for 1 h at 65°C under agitation. To avoid breaking the DNA during pipetting, the extracted DNA was recovered using a Phase-lock gel tube as described in Belser (Belser et al. 2018). The extracts were transferred to 2 ml tubes containing phase-lock gel, and an equal volume of PCIA (Phenol, Chloroform, Isoamyl Alcohol; 25:24:1) was added. After shaking for 5 min, tubes were centrifuged at 11000 g for 20 min. The aqueous phase was transferred into a new tube containing phase-lock gel and extraction with PCIA was repeated. DNA was then precipitated after addition of an equal volume of binding buffer C4 (Macherey-Nagel) and 99% ethanol overnight at 4°C or 1 h in ice then centrifuged at 800 rpm for 10 min. After removal of the supernatant, 1 ml of CQW buffer was added then the pellet of DNA was re-suspended. Next, DNA purification was carried out by adding a 2 ml mixture of wash buffer PW2 (Macherey-Nagel), wash buffer B5 (Macherey-Nagel), and ethanol at 99% in equal volumes, followed by centrifugation at 800 rpm for 10 min. This DNA purification step was carried out twice. Finally, the DNA pellet was dried in the oven at 70°C for 30 min then re-suspended in 100 µl elution buffer BE (Macherey-Nagel) (5 mM Tris solution, pH 8.5) after 10 min incubation at 65°C under agitation.

A second purification step was performed using a PCR product extraction from gel agarose kit from Macherey-Nagel (MN) NucleoSpin® Gel and PCR Clean-up kit and restarting the above protocol from the step with the addition of CQW buffer then PW2 buffer.

The third purification step consisted of DNA purification using a Macherey-Nagel (MN) NucleoMag kit for clean-up and size selection. Finally, the DNA was resuspended after a 5 min incubation at 65°C in 5 mM TRIS at pH 8.5.

The quantity and quality of the gDNA was verified using a NanoDrop spectrometer, electrophoresis on agarose gel and ethidium bromide staining under UV light and Fragment Analyzer (Agilent Technologies) of the University of Rennes1.

Library preparation and sequencing

MiSeq (Illumina) and GridION (Oxford Nanopore Technologies, referred to here as ONT) sequencing were performed at the PGTB (doi:10.15454/1.5572396583599417E12). *Lgh* and *Lpm* genomic DNA were re-purified using homemade SPRI beads (1.8X ratio). *Lgh* has a large genome size of 1419 Mb, 5-fold larger than *Lpm* genome 262 Mb (Barloy et al. 2024). SR (Illumina, one run) and LR (Oxford Nanopore, three runs) sequencing were therefore carried out for *Lgh* and only LR sequencing for *Lpm* (one run). For Illumina sequencing, 200 ng of *Lgh* DNA was used according to the QIAseq FX DNA Library Kit protocol (Qiagen). The final library was checked on TapeStation D5000 screentape (Agilent Technologies) and quantified using a QIAseq Library Quant Assay Kit (Qiagen). The pool was sequenced on an Illumina MiSeq using V3 chemistry and 600 cycles (2x300bp). For ONT sequencing, around 8 µg of *Lgh* and *Lpm* DNA were size selected using a Circulomics SRE kit (according to the manufacturer's instructions) before library preparation using a SQK-LSK109 ligation sequencing kit following ONT recommendations. Basecalling in High

Accuracy - Guppy version: 4.0.11 (MinKNOW GridION release 20.06.9) was performed for the 48 h of sequencing. Long reads (LR) and short reads (SR) were available for *Lgh* and only LR for *Lpm*.

Chloroplast assemblies

Quality controls and preprocessing of sequences were conducted using Guppy v4.0.14 for long reads (via Oxford Nanopore Technology Client access) and fastp v0.20.0 (Chen et al. 2018) for short reads, using Q15, since increasing the Phred quality to 20 or higher has no effect on the number of sequences retained (66%). A preliminary draft assembly was performed using *Lgh* short-reads (SR, 2*23,067,490 reads) with GetOrganelle v1.7.0 (Jin et al. 2020) and NOVOPlasty v4.2.1 (Dierckxsens et al. 2017), and chloroplastic short and long reads were extracted by mapping against this draft genome. Chloroplastic short reads were then *de novo* assemble using Velvet (version 1.2.10) (Zerbino and Birney 2008), ABySS (version 2.1.5, Simpson et al. 2009, Jackman et al. 2017), MEGAHIT (1.1.2, Li et al. 2016), and SPAdes (version 3.15.4, Bankevich et al. 2012), without and with prior error correction. The best k-mer parameters were tested using kmergenie (Chikhi and Medvedev 2014) and k=99 was found to be optimal. For ONT reads, *Lgh* (550,516 reads) and *Lpm* (68,907 reads) reads were self-corrected using CANU 1.8 (Koren et al. 2017) or SR-corrected using Ratatosk (Holley et al. 2021) and *de novo* assembly using CANU (Koren et al. 2017) and FLYE 2.8.2 (Kolmogorov et al. 2019) run with the option—meta and —plasmids. For all these assemblers, unless otherwise specified, we used the default parameters.

Post plastome assembly validation

As we used many assemblers and different strategies, we produced multiple contigs that needed to be analyzed and filtered in order to retain only the most robust plastomes. For that, all assemblies were evaluated using the QUality ASsessment Tool (QUAST) for quality assessment (Gurevich et al. 2013) and visualized using BANDAGE (Wick et al. 2015), both using default parameters. BANDAGE compatible graphs (.gfa format) were created with the megahit_toolkit for MEGAHIT (Li et al. 2016) and with gfatools for ABySS (Jackman et al. 2017). Overlaps between fragments were manually checked and ambiguous “IUPAC or N” nucleotides were also biocured with Illumina reads when available.

Chloroplast genome annotation

Plastomes were annotated via the GeSeq (Tillich et al. 2017) using ARAGORN and tRNAscan_SE to predict tRNAs and rRNAs and tRNAscan_SE to predict tRNAs and rRNAs and via Chloe prediction site (Zhong 2020). The previously reported *Lo* chloroplast genome was also similarly re-annotated to facilitate genomic comparisons. Gene boundaries, alternative splice isoforms, pseudogenes and gene names and functions were manually checked and biocured using Geneious (v.10). Finally, plastomes were represented using OrganellarGenomeDRAW (OGDRAW, Greiner et al. 2019). These genomes were submitted to GenBank at the National Center of Biotechnology Information (NCBI) with specific accession numbers (for *Lgh* haplotype 1, (LGH1) OR166254 and *Lgh* haplotype 2, (LGH2) OR166255; for *Lpm* haplotype, (LPM) OR166256) using annotation tables generated through GB2sequin (Lehwark and Greiner 2019).

SSRs and Repeat Sequences Analysis

Simple Sequence Repeats (SSRs) were analyzed through the MISA web (MISA-web) server (Beier et al. 2017), with parameters set to 10, 5, 4, 3, 3, and 3 for mono-, di-, tri-, tetra-, penta-, and hexa-nucleotides, respectively. Direct, reverse and palindromic repeats were identified using RepEx (Gurusaran, Ravella, and Sekar 2013). Parameters used were: for inverted repeats (min size 15 nt, spacer = local, class = exact); for palindromes (min size 20 nt); for direct repeats (minimum size 30 nt, minimum repeat similarity 97%). Tandem repeats were identified using Tandem Repeats Finder (Benson 1999), with parameters set to two for the alignment parameter match and seven for mismatches and indels. The IRa region was removed for all these analyses to avoid over representation of the repeats.

Comparative chloroplast genomic analyses

Lgh and *Lpm* plastomes were compared with the reannotated and biocurated *Lo* plastome using mVISTA program (Frazer et al. 2004), with the LAGAN alignment algorithm (Brudno et al. 2003) and a cut-off of 70% identity. Nucleotide diversity (π) was analyzed using the software DnaSP v.6.12.01 (Rozas and Rozas 1999, Rozas et al. 2017) with step size set to 200 bp and window length to 300 bp. IRscope (Amiryousefi et al. 2018) was used for the analyses of inverted repeat (IR) region contraction and expansion at the junctions of chloroplast genomes. To assess the impact of environmental pressures on the evolution of these three *Ludwigia* species, we calculated the nonsynonymous (K_a) and synonymous (K_s) substitutions and their ratios ($\omega = K_s/K_a$) using TBtools (Chen et al. 2020) to measure the selective pressure. Genes with $\omega < 1$, $\omega = 1$, and $1 < \omega$ were considered to be under purifying selection (negative selection), neutral selection, and positive selection, respectively.

Phylogenetic analysis of *Ludwigia* based on MatK sequences

We performed a phylogenetic analysis on the *Ludwigia* genus using the MatK, only protein coding barcode available for a large number of *Ludwigia* species. All MatK amino acid sequences were aligned with the FFT-NS-2 (Fast Fourier Transform-based Narrow Search) algorithm and BLOSUM62 scoring matrix using MAFFT 7 (Katoh et al. 2002). The phylogenetic tree analysis was conducted using the rapid hill-climbing algorithm (command line: -f d) in RAxML 8.2.11 (Stamatakis 2014), with GAMMA JTT (Jones-Taylor-Thornton) protein model. Node support was assessed through fast bootstrapping (-f a) with 1,000 non-parametric bootstrap pseudo-replicates. *Circaea* MatK were selected as outgroup, and all accession numbers are indicated on the phylogenetic tree labels.

Graphic representation

Statistical analyses were performed using R software in RStudio integrated development environment (R Core Team 2015; RStudio: Integrated Development for R. RStudio, Inc., Boston, MA). Figures were realized using ggplot2, ggpubr, tidyverse, dplyr, gridExtra, reshape2, and viridis packages. SNPs were represented using trackViewer (Ou and Zhu 2019) and genes represented using gggenes packages.

Results

Plastome short read assembly

The chloroplastic fraction of *Lgh* short reads (SR) was extracted by mapping against the two draft haplotypes generated by GetOrganelle, which differ only by a “flip-flop” of the SSC region (Figure 1). Since the assembly by NOVOplasty did not provide any additional information compared to GetOrganelle, it was not retained. This subset (1,360,507 reads) was assembled using ABySS, Velvet, MEGAHIT and SPAdes in order to identify the best assembler for this plant model.

As shown in Figure 2, both the number and size of contigs depend greatly on the algorithms used and the correction step. The effect of prior read correction is notable for MEGAHIT and Velvet, especially concerning the increase in the size of the large alignment (Figure A1-A), loss of misassemblies, and reduction of the number of mismatches (Figure 1A-B). Investigating results via BANDAGE (Figure A2), we observed that ABySS and SPAdes suggest the tripartite structure with the long single-copy (LSC) region as the larger circle in the graph (blue), joined to the small single-copy region (green) by one copy of the inverted repeats (IRs, red), both IRs being collapsed in a segment of approximately twice the coverage. For Velvet and MEGAHIT, graphs confirm the significant fragmentation of the assemblies, which is improved by prior correction of the reads. In conclusion, none of the short-read assemblers tested in our study produced a complete plastome. The best result was achieved by SPAdes using corrected short reads (mean coverage 1900 X) to assemble a plastome consisting of three contigs: 90,272 bp (corresponding to LSC),

19,788 bp (corresponding to SSC), and 24,762 bp (corresponding to one of the two copies of the IR).



Figure 1 - Two structural haplotypes of *L. grandiflora* subsp. *hexapetala* plastomes representing the flip-flop organization of SSC segment

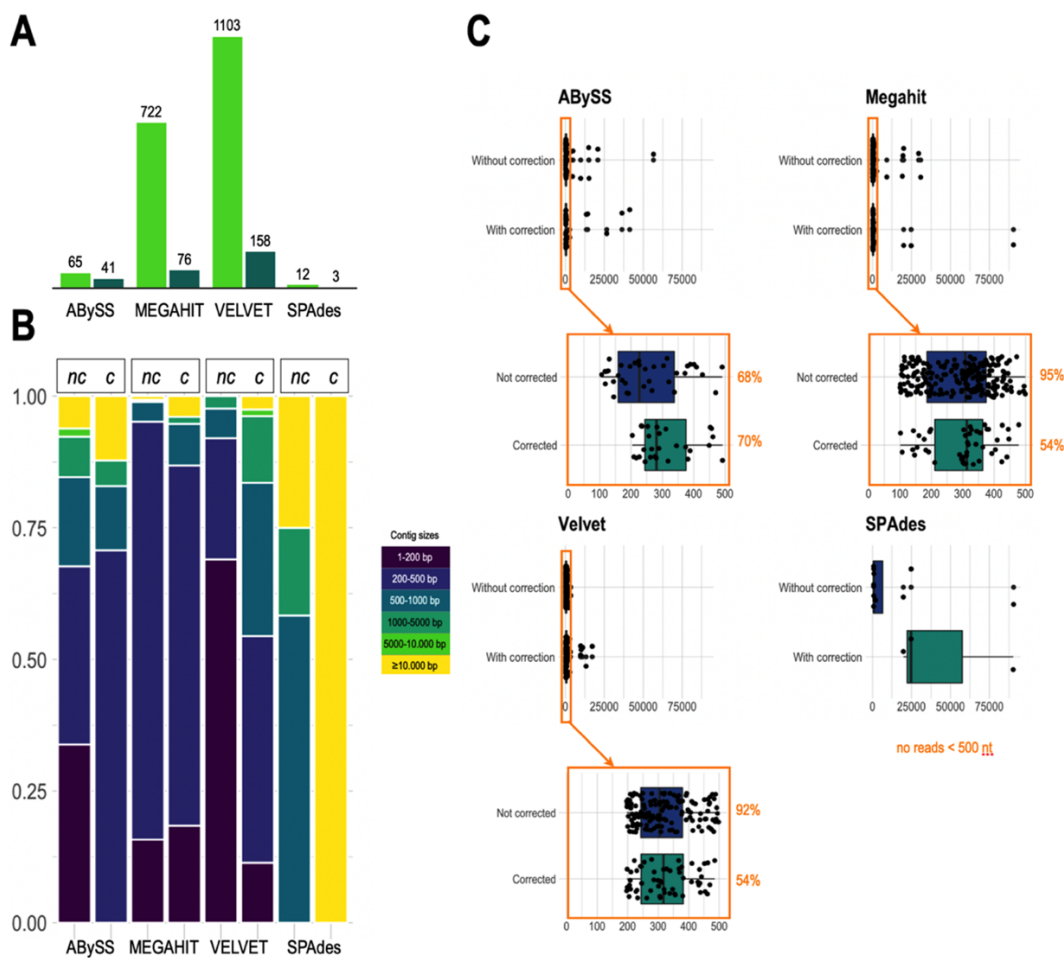


Figure 2 - Comparative results of *L. grandiflora* subsp. *hexapetala* short read (SR) assemblies. **A:** Total number of contigs obtained with the uncorrected (dark green) and corrected (light green) chloroplast SRs for the 4 assemblers (ABYSS, MEGAHIT, Velvet and SPAdes). **B:** Comparison of the size of contigs assembled by the 4 tools using corrected or uncorrected SRs. **C:** Boxplot showing the distribution of these contigs by size and the improvement brought by the prior correction of the SRs with the long reads for each tool.

Plastome long read assembly

Chloroplast fractions of *Lgh* long reads (28,882 reads) were assembled using CANU or FLYE. With raw data, CANU generates a unique contig (NGA50 112648) corresponding to haplotype 2, whereas FLYE makes two contigs (NGA50 133687) that reconstruct haplotype 1. Self-corrected LR leads to fragmentation into two (CANU) or three (FLYE) contigs which both reconstruct haplotype 1, with a large gap corresponding to one of the IR copies for CANU. Finally, SR-correction by RATATOSK allows CANU to assemble two redundant contigs reproducing haplotype 2 while FLYE makes two contigs corresponding to haplotype 1 (Figure A3). In conclusion, the two *Lgh* haplotypes were reconstructed (average coverage 700X) and the most complete and accurate hybrid assemblies (99.94% accuracy, Additional Figure 3B) were submitted to GenBank.

Unfortunately, due to the absence of short read data, we could only perform self-corrected long read assembly for *Lpm* using CANU. We also compared CANU and FLYE assembler efficiency, and found that assembly using CANU produces 13 contigs whereas FLYE produces 12 contigs. In both cases, only three contigs are required to reconstitute a complete cpDNA assembly (no gap, no N), with an SSC region oriented like those of the *Lgh* haplotype 2 and the *Lo* plastome. Although it is more than likely that these two SSC region orientations also exist for *Lpm*, the low number of nanopore sequences generated (68907 reads) and absence of Illumina short reads prevented us from demonstrating the existence of both haplotypes. As a result, only the “haplotype 2” generated sequence was deposited to Genbank.

Annotation and comparison of *Ludwigia* plastomes

General Variations

Plastomes of the three species of *Ludwigia* sp., *Lgh*, *Lpm* and *Lo*, are circular double-stranded DNA molecules (Figure 3) which are all (as shown in Table 1) approximately the same size: *Lo* is 159,396 bp long, making it the smallest, while *Lgh* is the largest with 159,584 bp, and *Lpm* is intermediate at 159,537 bp. The overall GC content is almost the same for the three species (37.4% for *Lo*, 37.3 % for *Lgh* and *Lpm*) and the GC contents of the IR regions are higher than those of the LSC and SSC regions (approximately 43.5 % compared to 35% and ca.32% respectively). Between the three species, the lengths of the total chloroplasts, LSC, SSC, and IR are broadly similar (approximately 90.2 kb for LSC, 19.8 kb for SSC and 24.8 kb for IB, see details Table 1) and the three plastomes are perfectly syntenic if we orient the SSC fragments the same way.

All three *Ludwigia* sp. plastomes contain the same number of functional genes (134 in total) encoding 85 proteins (embracing 7 duplicated in the IR region: *ndhB*, *rpl2*, *rpl23*, *rps7*, *rps12*, *ycf2*, *ycf15*), 37 tRNAs (including *trnK-UUU* which contains *matK*), and 8 rRNAs (16S, 23S, 5S, and 4.5S as duplicated sets in the IR). Among these genes, 18 contain introns, of which six are tRNAs (Table 2). Only the *rps12* gene is a trans-spliced gene. A total of 46 genes are involved in photosynthesis, and 71 genes related to transcription and translation, including a bacterial-like RNA polymerase and 70S ribosome, as well as a full set of transfer RNAs (tRNAs) and ribosomal RNAs (rRNAs). Six other protein-coding genes are involved in essential functions, such as *accD*, which encodes the β -carboxyl transferase subunit of acetyl-CoA carboxylase, an important enzyme for fatty acid synthesis; *matK* encodes for maturase K, which is involved in the splicing of group II introns; *cemA*, a protein located in the membrane envelope of the chloroplast is involved in the extrusion of protons and thereby indirectly allows the absorption of inorganic CO₂ in the plastids; *clpP1* which is involved in proteolysis, and; *ycf1*, *ycf2*, two ATPases members of the TIC translocon. Finally, a highly pseudogenized *ycf15* locus was annotated in the IR even though premature stop codons indicate loss of functionality.

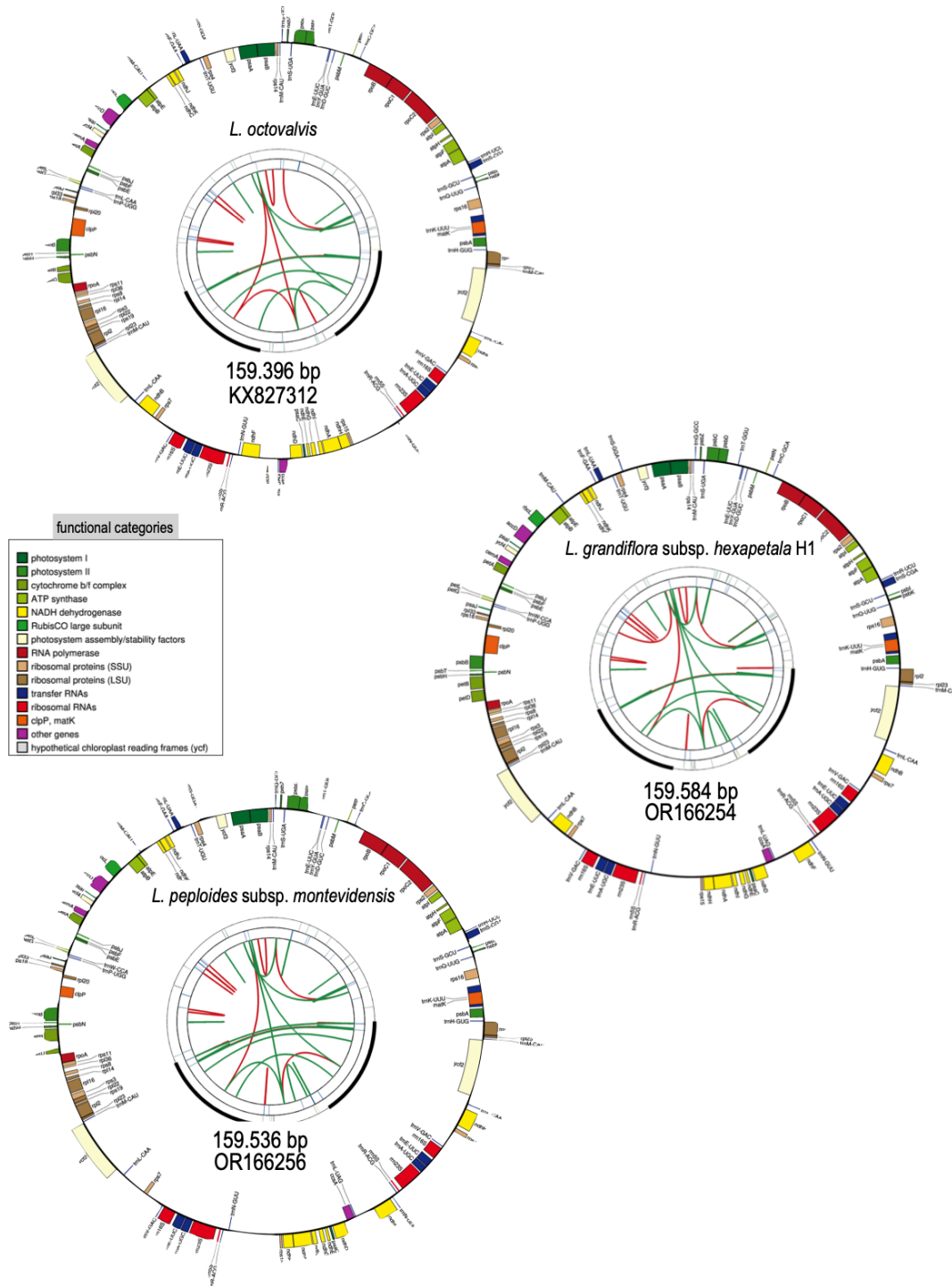


Figure 3 - Circular representation of annotations plastomes in *Ludwigia octovalvis*, *Ludwigia grandiflora* subsp. *hexapetala* and *Ludwigia peploides* subsp. *montevidensis* using ogdraw. Each card contains four circles. From the center outwards, the first circle shows forward and reverse repeats (red and green arcs, respectively). The next circle shows tandem repeats as bars. The third circle shows the microsatellite sequences. Finally, the fourth and fifth circles show the genes colored according to their functional categories (see colored legend). Only the haplotype 1 of *L. grandiflora* subsp. *hexapetala* is represented as haplotype 2 only diverge by the orientation of the SSC segment. Accession numbers are indicated for each plastome.

Table 1 - The general characteristics of the 3 *Ludwigia* plastomes

	<i>L. octovalvis</i> *	<i>L. grandiflora</i> subsp. <i>hexapetala</i>	<i>L. peploides</i> subsp. <i>montevidensis</i>
Size (bp)			
	159,396	159,584	159,537
LSC	90,183	90,272	90,156
SSC	19,703	19,788	19,799
IR	24,755	24,762	24,791
GC%			
	37.4	37.3	37.3
LSC	35.2	35.1	35.1
SSC	32	31.7	31.7
IR	43.5	43.5	43.4

* KX827312 (ref)

Table 2 - Genes present in the plastomes of *Ludwigia* sp.

Function	Name
Photosynthesis	
Rubisco	<i>rbcL</i>
Photosystem I (PSI)	<i>psaA, psaB, psaC, psal, psaJ</i>
PSI assembly factors	<i>ycf3[#] (pafl), ycf4 (pafl1)</i>
Photosystem II	<i>psbA, psbB, psbC, psbD, psbE, psbF, psbH, psbl, psbJ, psbK, psbL, psbM, pbf1 (psbN) psbT, psbZ</i>
ATP synthase	<i>atpA, atpB, atpE, atpF[#], atpH, atpI</i>
Cytochrome <i>b6f</i>	<i>petA, petB[#], petD[#], petG, petL, petN</i>
Cytochrome biogenesis	<i>ccsA</i>
NADPH dehydrogenase	<i>ndhA[#], ndhB^{**}, ndhC, ndhD, ndhE, ndhF, ndhG, ndhH, ndhI, ndhJ</i>
Transcription and translation	
Transcription	<i>rpoA, rpoB, rpoC1[#], rpoC2</i>
Small ribosomal proteins	<i>rps2, rps3, rps4, rps7^{**}, rps8, rps11, rps12^{**}, rps14, rps15, rps16[#], rps18, rps19</i>
Large ribosomal proteins	<i>rpl2^{**}, rpl14, rpl16[#], rpl20, rpl22, rpl23^{**}, rpl32, rpl33, rpl36</i>
Translation initiation	<i>infA</i>
Ribosomal RNA	<i>rrn5^{**}, rrn4,5^{**}, rrn16^{**}, rrn23^{**}</i>
Transfer RNA	<i>trnA-UGC^{**}, trnC-GCA, trnD-GUC, trnE-UUC, trnF-GAA, trnFM-CAU, trnG-GCC, trnG-UCC[#], trnH-GUG, trnI-CAU^{**}, trnI-GAU^{**}, trnK-UUU[#], trnL-CAA^{**}, trnL-UAA[#], trnL-UAG, trnM-CAU, trnN-GUU^{**}, trnP-UGG, trnQ-UUG, trnR-ACG^{**}, trnR-UCU, trnS-GCU, trnS-GGA, trnS-UGA, trnT-GGU, trnT-UGU, trnV-GAC^{**}, trnV-UAC[#], trnW-CCA, trnY-GUA</i>
Other functions	
Group II intron splicing	<i>matK</i>
Inorganic carbon uptake	<i>cemA</i>
Protease	<i>clpP1[#]</i>
Fatty acid synthesis/Heat tolerance	<i>accD</i>
TIC machinery (protein import)	<i>ycf1 (Tic214), ycf2^{**}</i>
Unknown function pseudogene	<i>ycf15^{**}</i>

** duplicated in IR region, # spliced genes

Segments Contractions/Expansion

The junctions between the different chloroplast segments were compared between three *Ludwigia* sp. (*Lpm*, *Lgh* and *Lo*), and we found that the overall resemblance of *Ludwigia* sp. plastomes was confirmed at all junctions (Figure 4A). In all three genomes, *rpl22*, *rps19*, and *rpl2* were located around the LSC/IRb border, and *rpl2*, *trnH*, and *psbA* were located at the IRa/LSC edge. The JSB (junction between IRb and SSC) is either located in the *ndhF* gene or the *ycf1* gene depending on the orientation of the SSC region (Figure 4B). The *ycf1* gene was initially annotated

as a 1139 nt pseudogene that we biocurate as a larger gene (5302 nt) with a frameshift due to a base deletion, compared to *Lgh* and *Lo* which both carry a complete *ycf1* gene.

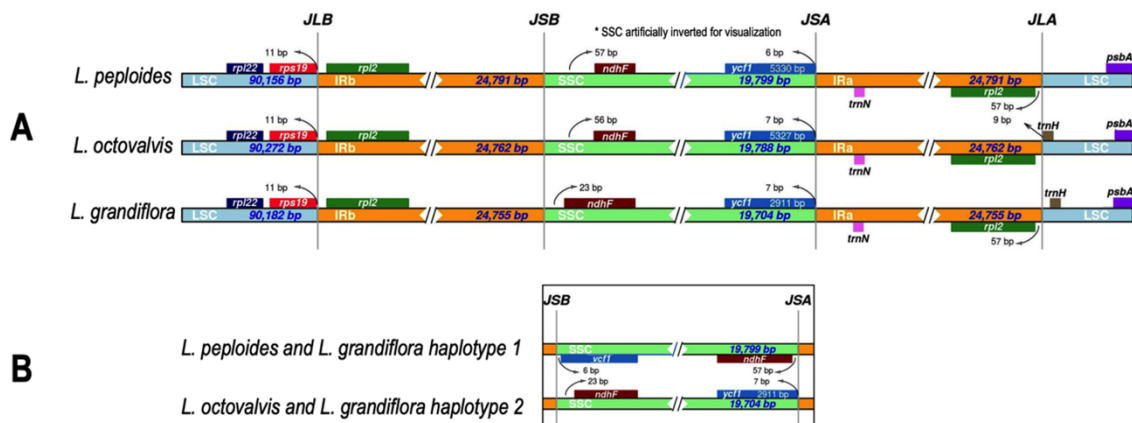


Figure 4 - Comparison of the borders of LSC, SSC, and IR regions in Onagraceae plastomes. **A**: Comparison of the junction between large single-copy (LSC, light blue), inverted repeat (IR, orange) and short single-copy (SSC, light green) regions among the chloroplast genomes of *L. octovalvis*, *L. peploides* subsp. *montevidensis* and *L. grandiflora* subsp. *hexapetala* (both haplotypes). Genes are denoted by colored boxes and the gaps between genes and boundaries are indicated by base lengths (bp). JLB: junction line between LSC and IRb; JSB: junction line between IRb and SSC; JSA: junction line between SSC and IRa; JLA: junction line between IRa and LSC. **B**: Comparison of SSC boundaries in haplotype 1 (*L. peploides* subsp. *montevidensis* and *L. grandiflora* subsp. *hexapetala* haplotype 1) and haplotype 2 (*L. octovalvis* and *L. grandiflora* subsp. *hexapetala* haplotype 2) plastomes.

If we compare *Ludwigia* sp. chloroplastic LSC/SSC/IR junctions (via IRscope) with representative Onagraceae plastomes of *Chamaenerion conspersum* (MZ353638) and *chamaenerion angustifolium* (NC_052848), *Circaea cordata* (NC_060876) and *Circaea alpina* (NC_061010), *Epilobium amurense* (NC_061015) and *Oenothera villosa* subsp. *strigosa* (NC_061365) and *Oenothera lindheimeri* (MW538951) (Figure 5), we can observe that the gene positions at the JLB (junction of LSC/IRb) and JLA (junction of IRa/LSC) boundary regions are well-preserved throughout the entire family, whereas those at the JSB and JSA regions differ. Concerning JSB (junction of IRb/SSC), in the five *Onagraceae* genera studied, *ndhF* is duplicated, with the exception of *Circaea* sp. and *Ludwigia* sp. For *Oenothera villosa*, the first copy of *ndhF*, which is located in the IRb, overlaps the JSB border, whereas for *Oenothera lindheimeri*, *Epilobium amurense* and *Chamaenerion* sp., *ndhF* is only located in inverted repeats. Only *Circaea* sp. and *Ludwigia* sp. have a unique copy of this locus, and it is found in the SSC segment (Figure 5). At the JSA border (junction of SSC/IRa), in *Circaea* sp., the *ycf1* gene crosses the IRa/SSC boundary and extends into the IRa region.

When comparing the respective sizes of chloroplast fragments (IR/SSC/LSC) in Onagraceae, it can be observed that *Ludwigia* species exhibit expansions in the SSC and LSC regions which are not compensated by significant contractions in the IR regions. This is likely due to the relocation of the *ndhF* in the SSC region and *rps19* in the LSC region. Additionally, there may be significant size variations in the intergenic region between *trnI* and *ycf2*, as well as the intergenic segment containing the *ycf15* pseudogene (Figure A4).

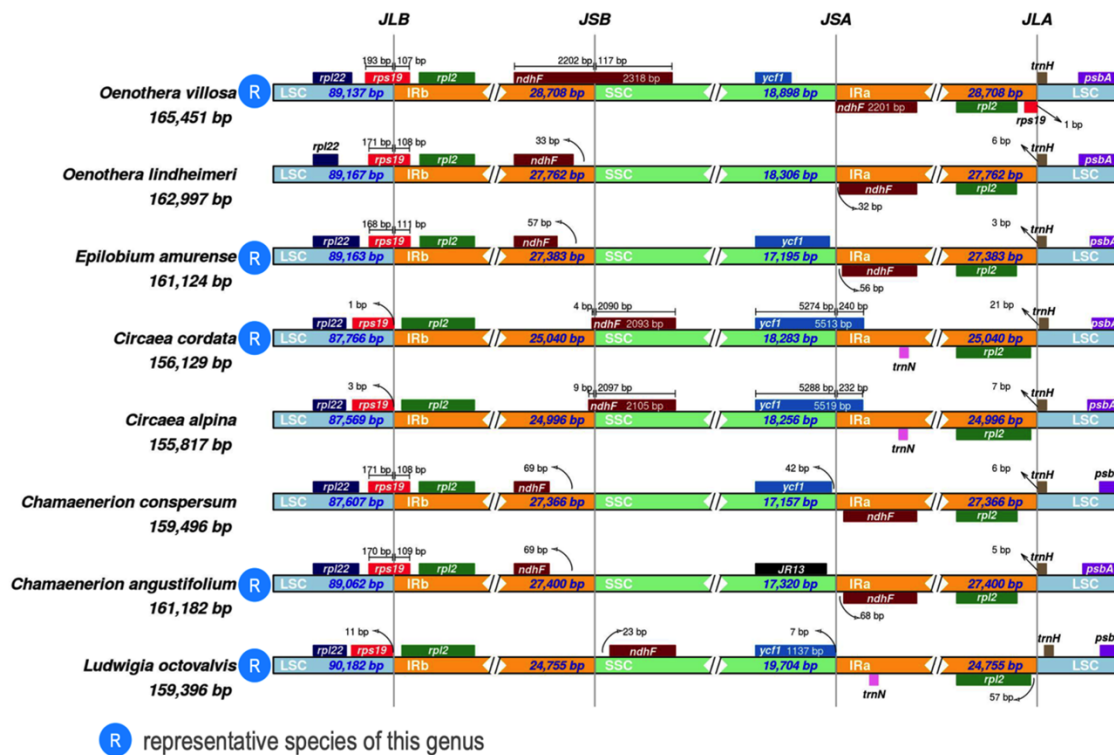


Figure 5 - Comparison of LSC, SSC and IR regions boundaries in Onagraceae chloroplast genomes. Representative sequences from each genus have been chosen (noted R on the diagram) except for *Oenothera lindheimeri* (only 89.35 % identity with others *Oenothera*), *Circaea alpina* (99.5 % identity but all others *Circaea* are 99.9% identical) and *Chamaenerion conspersum* (99% but all others *Chamaenerion* are ca. 99.7 identical). As shown in Figure 7, the 3 *Ludwigia* plastomas had the same structure, *L. octovalvis* was chosen as a representative of this genus. JLB: junction of LSC/IRb; JSB: junction of IRb/SSC; JSA: junction of SSC/IRa; JLA: junction of IRa/LSC. Accession numbers: *Chamaenerion conspersum* (MZ353638), *Chamaenerion angustifolium* (NC_052848), *Circaea cordata* (NC_060876), *Circaea alpina* (NC_061010), *Epilobium amurense* (NC_061015), *Oenothera villosa* subsp. *strigosa* (NC_061365) and *Oenothera lindheimeri* (MW538951).

Repeats and SSRs analysis

In this study, we analyzed the nature and distribution of single sequence repeats (SSR), as their polymorphism is an interesting indicator in phylogenetic analyses. A total of 65 (*Lgh*), 48 (*Lpm*) and 45 (*Lo*) SSRs were detected, the majority being single nucleotide repeats (38–21), followed by tetranucleotides (12–10) and then di-, tri- and penta-nucleotides (Figure A5-A). Mononucleotide SSRs are exclusively composed of A and T, indicating a bias towards the use of the A/T bases, which is confirmed for all SSRs (Figure A5-B). In addition, the SSRs are mainly distributed in the LSC region for the three species, which is probably biased by the fact that LSC is the longest segment of the plastome (Figure A5-C). The analysis of SRR locations revealed that most were distributed in non-coding regions (intergenic regions and introns, Figure A5-D).

The chloroplast genomes of the three *Ludwigia* species were also screened for long repeat sequences. They were counted in a non-redundant way (if smaller repetitions were included in large repeats, only the large ones were considered). Four types of repeats (tandem, palindromic inverted and direct) were surveyed in the three *Ludwigia* sp. plastomes. No inverted repeats were detected with the criteria used.

For the three other types of repeats, here are their distributions:

Tandem repeats (Table 3: Perfect tandem repeats (TRs) with more than 15 bp were examined. Twenty-two loci were identified in the three *Ludwigia* sp. plastomes (*Lgh*, *Lpm*, *Lo*), heterogeneously distributed as shown in Table 3: 13 loci (plus one imperfect) in *Lo*, nine loci (plus one imperfect) in *Lgh* and seven loci (plus two imperfect) in *Lpm*. It can therefore be seen that the TR distributions (occurrence and location) are specific to each plastome, since only four pairs are common to the three species. Thus, nine TRs are unique to *Lo*, three to *Lpm* and three to *Lgh*. Two pairs are common to *Lgh* and *Lpm* and one is common to *Lo* and *Lgh*. TRs are mainly intergenic or intronic but are detected in two genes (*accD* and *ycf1*). These genes have accelerated substitution rates, although this does not generate a large difference in their lengths. This point will be developed later in this article.

Direct repeats (Table 4): There are few direct (non-tandem) repeats (DRs) in the chloroplast genomes of *Ludwigia* sp. A single direct repeat of 41 nt is common to the three species, at 2 kb intervals, in *psaB* and *psaA* genes. This DR corresponds to an amino acid repeat [WLTDIAHHHLAIA] which corresponds to a region predicted as transmembrane. We then observe three direct repeats conserved in *Lpm* and *Lgh* in *ycf1*, *accD* and *clpP1* respectively, two unique DRs in *Lo* (in the *accD* gene and *rps12-clpP1* intergene) and one in *Lgh* (in the *clpP1* intron 1 and *clpP1* intron 2).

Palindromes (Table 5): Palindromic repeats make up the majority of long repetitions, with the numbers of perfect repeats varying from 19, 24 and 26 in *Lo*, *Lgh* and *Lpm*, respectively, and the number of quasi-palindromes (1 mutation) varying between 8, 3 and 6. They are mainly found in the intronic and intergenic regions, with the exception of six genic locations in *psbD*, *ndhK*, *ccsA* and *rpl22*, and two palindromic sequences in *ycf2*. These gene palindromic repeats do not seem to cause genetic polymorphism in *Ludwigia* and can be considered as silent.

Thirteen palindromes are common to the three species (including 2 with co-variations in *Lo*). 13 others present in *Lpm* and *Lgh* correspond to quasi-palindromes (QPs) in *Lo* due to mutated bases, and conversely, three *Lo* perfect palindromes are mutated in *Lpm* and *Lgh*. Finally, only five palindromes are species specific. Two in particular are located in the hypervariable intergenic spacer *ndhF-rpl32*, and are absent in *Lo* due to a large deletion of 160 nt.

Repeat distribution in LSC, SSC and IR segments

In the IRa/IRb regions, repeats are only identified in the first 9 kb region between *rpl2* and *ycf2*: a tandem repeat in the *Lpm rpl2* intron, and a tetranucleotide repeat, [TATC]*3, located in the *ycf2* gene in the three species. In *ycf2* we also found 1 common palindrome (16 nt), a single palindrome in *Lo* (20 nt, absent following an A:G mutation in the 2 other species), as well as a shared tandem repeat (24 nt), and an additional 15 nt tandem repeat in *Lo* which adds 4 amino acids to protein sequence.

In the SSC region, the repeats are almost all located in the intergenic and/or intronic regions, with a hotspot between *ndhF* and *ccsA*. There is also a shared microsatellite in *ndhF*, and a palidrome (16 nt) in *ccsA* which is absent in *Lo* (due to an A:C mutation), resulting in a synonymous mutation (from isoleucine to leucine). We also observed multiple and various repeats in the *ycf1* gene: 3 common poly-A repeats (from 10 to 13 nt), 3 species-specific microsatellites (ATAG)*3 and (ACCA)*4 in *Lgh* and (CAAC)*3 in *Lo*, as well as two direct repeats of 32 nt (37 nt spacing), which were absent from *Lo* due to a G:T SNP. Two tandem repeats were also observed in *Lo* and *Lgh*. Neither of these repeats are at the origin of the frameshift causing the pseudogenization of *ycf1* in *Lo*, this latter being due to a single deletion of an A at position 3444 of the gene.

Table 3 - Tandem repeats

Sequence	<i>L. octovalvis</i> (L.o)	<i>L. grandiflora</i> (L.g)	<i>L. peplioide</i> s (L.p)	Length	Region	Locus	Comments
TTGTAGTCAGGGGTGTAGTACTAT				24	IRs	<i>ycf2</i>	
TAGAAGAGAGTGCAG		X	X	15	IRs	<i>ycf2</i>	15 nt deletion in L.g and L.p
ATGAAATATCGTATAATGAAGTACCACACGAGTGGATAT	X	X		39	IRs	<i>rp2</i> intron	39 nt deletion in L.g and L.o
AAAAATAGGATAGGAT		X	X	16	LSC	<i>ycf1-trmH-GUG</i>	56 nt deletion in L.g and L.p
TAAATTAATATCTATATA		X	X	18	LSC	<i>psbZ-trnG-GCC</i>	18 nt deletion in L.g and L.p
TTTTCTATCTATCTTATATCAA		X	X	22	LSC	<i>trnK-UUU-rps16</i>	22 nt deletion in L.g and L.p
AGATCCATAACATCATCAAA		X	X	20	LSC	<i>rps16</i> intron	22 nt deletion in L.g and L.p
TATTAGTTATTAAATATTATAGA		X	X	23	LSC	<i>trnP-UGG-psaJ</i>	23 nt deletion in L.g and L.p
AATAATATATAAATACTAAATA		X	X	23	LSC	<i>rp133-rps18</i>	33 et 44 nt deletion in L.g et L.p, respectively
TTTTTATTTAACATGCTATCAAAATCAACAATGCCATACCGTAGGGCATCTGTT		X	X	53	LSC	<i>rp20-clpP1</i>	107 nt deletion in L.g and L.p 3 copies in a 57 nt deletion in L.o and L.p
ATATATTTTCGATTCAATTC	X		X	19	LSC	<i>trmH-GUG-psbA</i>	
ATAGAAAATATCAGTATTTGAGTG	X		X	23	LSC	<i>atpH-atpI</i>	23 nt deletion in L.o and L.p 17 and 24 nt deletion in L.o
TTAATTTTAAITGAAGAA	X		X	18	LSC	<i>psbJ-psbL</i>	and L.p, respectively A -> C mutation in second copy in L.o
TTAAAGAAATATTAATATTC	imperfect TR			19	LSC	<i>trnR-UCU-atpA</i>	
TATTATTATTATTAAT	X	X		16	LSC	<i>atpH-atpI</i>	16 nt deletion in L.g and L.o
TCTAAGGCTGAAATAAGG	X	X		18	LSC	<i>pafl</i> intron	18 nt deletion in L.g and L.o
TGTGAATCTATCTAT			X	15	LSC	<i>trnS-UGA-psbZ</i>	8 nt deletion in L.p
TTTTTTCTAGTA	imperfect TR		imperfect TR	12	LSC	<i>pafl</i> intron	G -> A mutation in second in L.p et L.g
CTAGTTATTGACATGG			TR	16	LSC	<i>psaJ-rp133</i>	T->A mutation in first copy in L.p, other sequence in first copy in L.o
ATTTTTTAACTCT	X		imperfect TR	15	SSC	<i>ycf1</i>	other sequence in first copy of L.p and L.g
AATCAAAATAGTTGAT		X	X	15	SSC	<i>ycf1</i>	
ATAATAATATATTTTATTATTAATAATA	X			28	SSC	<i>ndhF-rp132</i>	160 nt deletion in L.o

Lo = *Ludwigia octovalvis*; Lgh = *L. grandiflora* subsp. *hexapetala*; Lpm = *L. peplioides* subsp. *montevidensis*.

Table 4 - Direct repeats

Sequence	<i>L. octovalis</i> (Lo)	<i>L. grandiflora</i> subsp. <i>hexapetala</i> (Lgh)	<i>L. peploides</i> subsp. <i>montevicensis</i> (Lpm)	Size (nt)	Spacers (nt)	Region	Locus	Comments
TTCAATTGGAACGGACGATTTCGTCAATCATCT				32	37	SSC	<i>ycf1</i>	2 copies. In <i>Lo</i> , one mutation (G->A) in the second copie
CATCGATGATGAAAGTGAAAACAGTAATGAAGAGG	X			35	28 - 22 - 11	LSC	<i>accD</i>	3 perfects copies and 1 mutated (G->A) copie in <i>Lgh</i> and <i>Lpm</i> Region of 174 nt deleted in <i>Lo</i>
TTAAGAGCCGTACAGGCACCTTTTGATGCATACGG	X				408 in <i>Lpm</i> , 406 in <i>Lgh</i>	LSC	<i>clpP1</i>	2 copies. In <i>Lgh</i> , one mutation (C->T) in the second copie
AGATGGTGAAGAACCCTTATGAAGATGGTGAAGAACCCTTATG		X	X	41	22	LSC	<i>accD</i>	Region of 147 nt deleted in <i>Lgh</i> and <i>Lpm</i>
TATCAAATCAACAATGCCATACCGTAGGGCAT		X	X	32	22 - 21	LSC	<i>rps12-clpP1</i>	3 copies
TTAAGAGCCGTACAGGCACCTTTTGATGCATACGG	X		X	35	811	LSC	<i>clpP1</i> intron 1- intron 2	
TGCAATAGCCAAATGATGATGAGCAAATATCAGTCAGCCATA				41	2178	LSC	<i>psaB & psaA</i>	

Lo = *Ludwigia octovalis*; Lgh = *L. grandiflora* subsp. *hexapetala*; Lpm = *L. peploides* subsp. *Montevicensis*

Table 5 - Palindromic repeats

Common perfect palindromic repeats	
AGACTCTCATGAGAGTCT	<i>trnC-GCA - petN</i>
ATTAATAGAAATATTCTATTTAAT	<i>trnE-UUC-trnT-GGU</i>
TTGGTAAATTTACCAA	<i>psbD</i>
TTCAATTTCAATTTCAATTTGAAATTTGAAATGAA	<i>trnI-CAU-ycf2</i> 2 copies in IR
GAAAAGGCCCTTTTTC	<i>ycf2</i> 2 copies in IR
TCTCAAATGATTAATCATTGAGA	<i>trnL-UAA</i> intron
GGATTACTAGTAATCC	<i>trnD-GUC-trnY-GUA</i>
TTTGAATGCATTCAAA	<i>trnG-UCC</i> intron
ATATATTCGAATATAT	<i>trnG-UCC -trnR-UUCU</i>
TAGTAATTAATTACTA	<i>trnG-GCC-trnifM-CAU</i>
CCAGTATGCATACTGG	<i>ndhK</i>
Common palindromic repeats with covariation	
<i>in L. octovalvis</i>	
ATAGAATCTATATTCTATTAGAATATAGATTCAT	<i>ndhC-trnV-UAC</i>
ATGTATATATCGAT	<i>trnE-UUC-trnT-GGU</i>
<i>in L. grandiflora et L. peploides</i>	
ATCGAATCTATATTCTATTAGAATATAGATTCGAT	
ATCTATATATATAGAT	
Common palindromic and quasi-palindromic repeats	
<i>in L. octovalvis</i>	
TTTAACGAAATTAATATT † GTTAAA	<i>trnR-UUCU-atpA</i>
TTAA c GAATTAATTAATTTCTTTAA	<i>trnR-UUCU-atpA</i>
AATTGTA c TTACAATT	<i>ccsA</i>
AGGAAGATTGATCAATCTT † CT	<i>trnL-UAG-rpl32</i>
TTA c TAATATTACTAA	<i>trnK-UUU</i> intron
ATATAGAATAT c CTATAT	<i>psbZ-trnG-GCC</i>
ACATATCATGATA g GT	<i>rpl22</i>
<i>in L. grandiflora and L. peploides</i>	
TTTAACGAAATTAATTAATTTCTGTTAAA	
TTAAAGAAATTAATTAATTTCTTTAA	
AATTGTAATTACAA TT	
AGGAAGATTGATCAATCTTCTCT	
TTAGTAATATTACTAA	
ATATAGAATATTCTATAT	
ACATATCATGATATGT	

Table 5 - Palindromic repeats - Continued

AATTACTAATTTCTATTACTAGTAAATGAAACATAGTAAATAGTAATAA	AATTACTAATTTCTATTACT † TGTTCAATTTGAACATAGTAATAGAAAATTAGTAATT	<i>atpH-atpI</i>
TAGTTAGAATTTCTAACTA	TAGTT † c GAATTTCTAACTA	<i>trnT-UGU-trnL-UAA</i>
TATTTTTTCTAGAAAAATA	TATTTTTTCTAGAA † g AAATA	<i>ycf2</i> 2 copies in IR
<i>in L. octovalvis and L. peplioides</i>	<i>in L. grandiflora</i>	
CCCATCAATCATGATT † TGGG	CCCATCAATCATGATTGATGGG	<i>psbN-trnD-GUC</i>
<i>in L. octovalvis and L. grandiflora</i>	<i>in L. peplioides</i>	
ATGAAAAAATCGATTTTTTTCAT	ATGATAAAAAATAGATTTTT † a TCAT	<i>trnK-UUU-rps16</i>
ATGAAAAAATCGATTTTTTTCAT- ATGATAAAAAATCGATTTTTTTCAT	ATGATAAAAAATA † g ATTTTTTATCAT	<i>trnK-UUU-rps16</i>
Unique palindromic repeats		
<i>L. peplioides</i>		
TTATATATATATATATAA		Full deletion in <i>L.o</i> 6 bases deletion in <i>Lgh</i>
<i>L. octovalvis</i>		
ATTGAAATTCGAAATTTCAAT		Full deletion in <i>Lgh</i> and <i>Lpm</i>
<i>L. peplioides and L. grandiflora</i>		
AAAAATGGATCCATTTTTT		3 bases deleted and 3 bases mutated in <i>Lo</i>
AATATATTATTATAATAATATAT		Full deletion in <i>Lo</i>
TATATTTATTATTAATAATAATAATA		Full deletion in <i>L.o</i>

Lo = *Ludwigia octovalvis*; *Lgh* = *L. grandiflora* subsp. *hexapetala*; *Lpm* = *L. peplioides* subsp. *montevicensis*.

Finally, in the LSC region, the longest segment, which consequently contains the maximum number of repeats, we still observed a preferential localization in the intergenic and intronic regions since only genes *atpA*, *rpoC2*, *rpoB*, *psbD*, *psbA*, *psbB*, *ndhK* and *clpP1* contain either mononucleotic repeats (poly A and T), palindromes, or microsatellites (most often common to the three species and without affecting the sequences of the proteins produced). As mentioned earlier, the only exception is the *accD* gene, which contains several direct and tandem repeats in *Lgh* and *Lpm*, corresponding to a region of 174 nt (58 amino acids) missing in *Lo* and, conversely, a direct repeat of 40 nucleotides, in a region of 147 nt (49 aa), which is present in *Lo* and missing in the other two species. These tandem repeats lead to the presence of four copies of nine amino acids [DESENSNEE] in *Lgh* and *Lpm*, two of which form a larger duplication of 17 aa [FLSDSDIDDESENSNEE]. Similarly, the TRs present only in *Lo* generate two perfect nine amino acid repeats [EELSEDGEE], included in two longer degenerate repeats of 27 nt (Figure A6). It should be noted that though these TRs do not disturb the open reading phases, it is still possible for them to form an intron which is not translated. Different functional studies will be necessary to clarify this point. The presence of polymorphisms of the *accD* gene between *Lo* and the two species (*Lpm*, *Lgh*) is interesting because *accD*, that encodes a subunit of acetyl-CoA carboxylase (EC 6.4.1.2). This enzyme is essential in fatty acid synthesis and also catalyzes the synthesis of malonyl-CoA, which is necessary for the growth of dicots, plant fitness and leaf longevity, and is involved in the adaptation to specific ecological niches (Konishi and Sasaki 1994). Large *accD* expansions due to TRs have also been described in other plants such as *Medicago* (Wu et al. 2021) and *Cupressophytes* (Li et al. 2018). Some authors have suggested that these inserted repeats are not important for acetyl-CoA carboxylase activity as the reading frame is always preserved, and they assume that these repeats must have a regulatory role (Gurdon and Maliga 2014).

Sequence Divergence Analysis and Polymorphic Loci Identification

Determination of divergent regions by MVista, using *Lo* as a reference, confirmed that the three *Ludwigia* sp. plastomes are well preserved if the SSC segment is oriented in the same way (Figure A7). Sliding window analysis (Figure 6) indicated variations in definite coding regions, notably *clpP*, *accD*, *ndh5*, *ycf1* with high Pi values, and to a lesser extent, *rps16*, *matK*, *ndhK*, *petA*, *ccsA* and four tRNAs (*trnH*, *trnD*, *trnT* and *trnN*). These polymorphic *loci* could be suitable for inferring genetic diversities in *Ludwigia* sp.

A comparative analysis of the sizes of protein coding genes sizes also shows that the *rps11* gene initially annotated in *Lo* is shorter than those which have been newly annotated in *Lgh* and *Lpm* (345 bp instead of 417 bp). Comparative analysis by BLAST shows that it is the long form which is annotated in other Myrtales, and the observation of the locus in *Lo* shows a frameshift mutation (deletion of a nucleotide in position 311). Functional analysis would be necessary to check whether the *rps11* frameshift mutation produces shorter proteins that have lost their function. And only obtaining the complete genome will verify whether copies of some of these genes have been transferred to mitochondrial or nuclear genomes. Such *rps11* horizontal transfers have been reported for this gene in the mitochondrial genomes of various plant families (Richardson and Palmer 2007). This also applies to *ycf1*, found as a pseudogene in *Lo* (as specified previously), although it is not known if this reflects a gene transfer or a complete loss of function (de Vries et al. 2015, Filip and Skuza 2021). Moreover, there is a deletion of nine nucleotides in the 3' region of the *rp132* gene in *Lgh* and *Lpm*, leading to a premature end of the translation and the deletion of the last four amino acids [QRLD], which are replaced by a K. However, if we look carefully at the preserved region as defined by the RPL32 domain (CHL00152, member of the superfamily CL09115), we see that the later amino acids are not important for *rp132* function since they are not found in the orthologs.

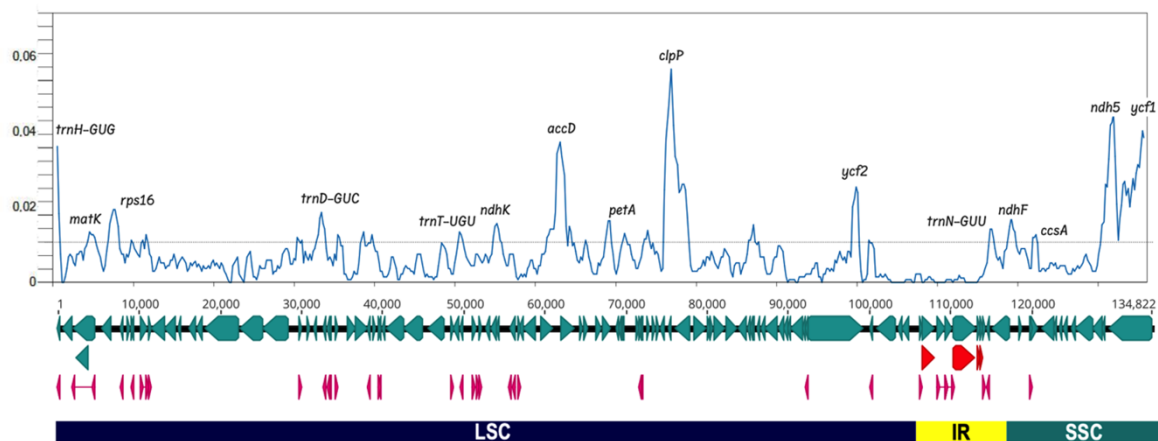


Figure 6 - Illustration of nucleotide diversity of the three *Ludwigia* chloroplast genome sequences. The graph was generated using DnaSP software version 6.0 (windows length: 800 bp, step size: 200 bp) (Rozas et al. 2017, Rozas and Rozas 1999). The x-axis corresponds to the base sequence of the alignment, and the y-axis represents the nucleotide diversity (π value). LSC, SSC and IR segments were indicated under the line representing the genes coding the proteins (in light blue) the tRNAs (in pink) and the rRNAs (in red). The genes marking diversity hotspots are noted at the top of the peaks.

Our results show that the K_a/K_s ratio is less than 1 for most genes (Figure 7). This indicates adaptive pressures to maintain the protein sequence except for *matK* (1.17 between *Lgh* and *Lpm*), *accD* (2.48 between *Lgh* and *Lo* and 2.16 between *Lpm* and *Lo*), *ycf2* (4.3 between both *Lgh-Lp* and *Lo*) and *ccsA* (1.4 between both *Lgh-Lpm* and *Lo*), showing a positive selection for these genes, and a possible key role in the processes of the species' ecological adaptations. As we have already described the variability in the *accD* sequence, we will focus on *ycf2*, *matK*, and *ccsA* variations.

Concerning *ccsA*, the variations observed, although significant, concern only five amino acids, and modifications do not seem to affect the C-type cytochrome synthase gene function.

Concerning *ycf2*, our analysis shows that this gene is highly polymorphic with 256 SNPs that provoke 10 deletions, 7 insertions, 21 conservative and 49 non-conservative substitutions in *Lo* (Figure A8), compared to *Lgh* and *Lpm* (100 % identical). This gene has been shown as "variant" in other plant species such as *Helianthus tuberosus* (Zhong et al. 2019).

The *matK* gene has been used as a universal barcoding locus to enable species discrimination of terrestrial plants (Antil et al. 2023), and is often, together with the *rbcl* gene, the only known genetic resource for many plants. Thus, we propose a phylogenetic tree from *Ludwigia matK* sequences (Figure 8). It should however be noted that this tree contains only 149 amino acids common to all the sequences (out of the 499 in the complete protein). As only three complete *Ludwigia* plastomes are available at the time of our study, we cannot specify whether these barcodes are faithful to the phylogenomic history of *Ludwigia* in the same way as the complete plastome. In any case, for this tree, we can see that *Lo* stands apart from the other *Ludwigia* sp., *Lpm* and *Lgh*, and that the *L. grandiflora* subsp. *hexapetala* belongs to the same branch as the species *L. ovalis* (aquatic taxon used in aquariums (Li J et al. 2022), *L. stolonifera* (native to the Nile, found in a variety of habitats, from freshwater wetlands to brackish and marine waters) (Soliman et al. 2018) and *L. adscendens* (common weed of rice fields in Asia) (Kamoshita et al. 2016). *Lpm* is in a sister branch, close to the *L. grandiflora* subsp. *hexapetala*, forming a phylogenetic group corresponding to subsect *Jussiaea* (in green, Figure 8).

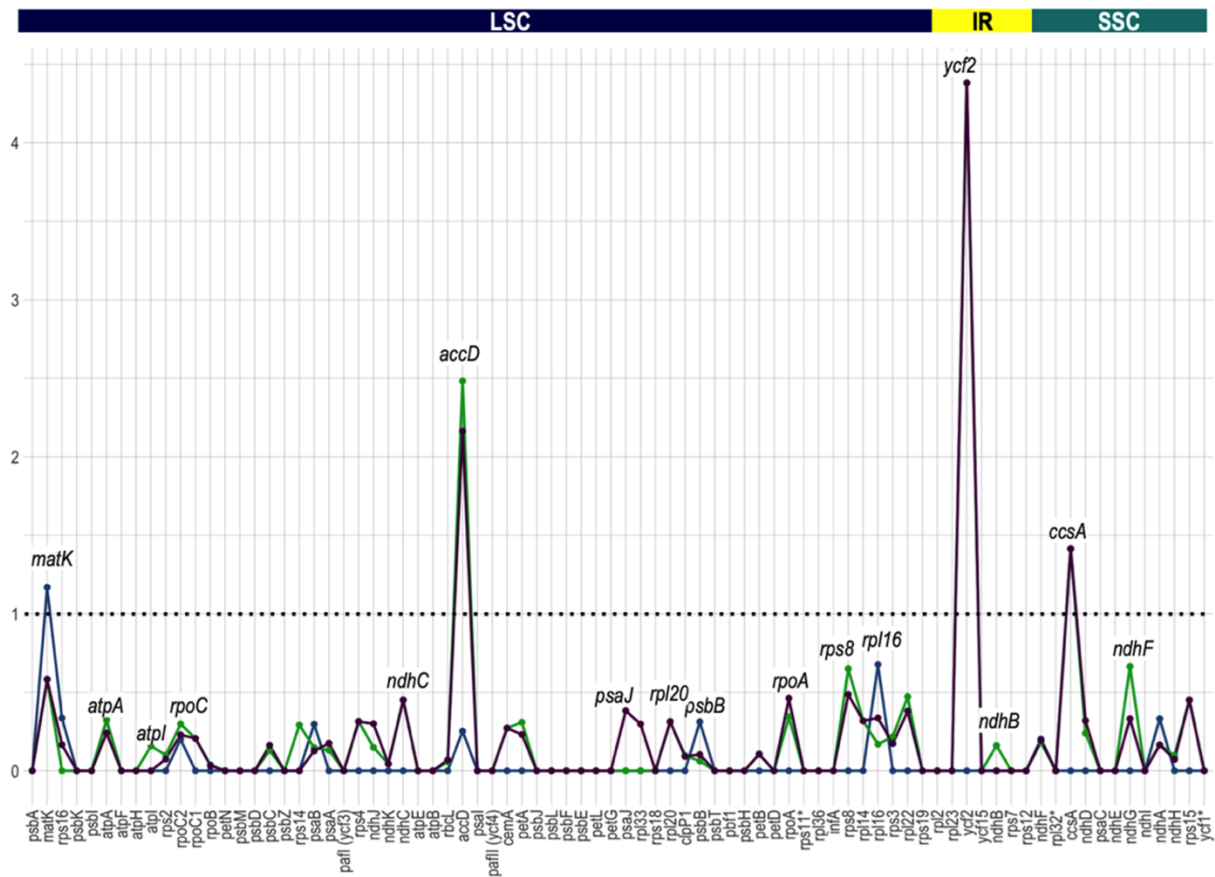


Figure 7 - The Ka/Ks ratios of the 80 protein-coding genes of *Ludwiglia* plastomes. The blue curve represents *L. grandiflora* subsp. *hexapetala* versus *L. peplodes* subsp. *montevidensis*, purple curve denotes *L. grandiflora* subsp. *hexapetala* versus *L. octovalvis* and green curve *L. peplodes* subsp. *montevidensis* versus *L. octovalvis*. Four genes (*matK*, *accD*, *ycf2* and *ccsA*) have Ka/Ks ratios greater than 1.0, whereas the Ka/Ks ratios of the other genes were less than 1.0.

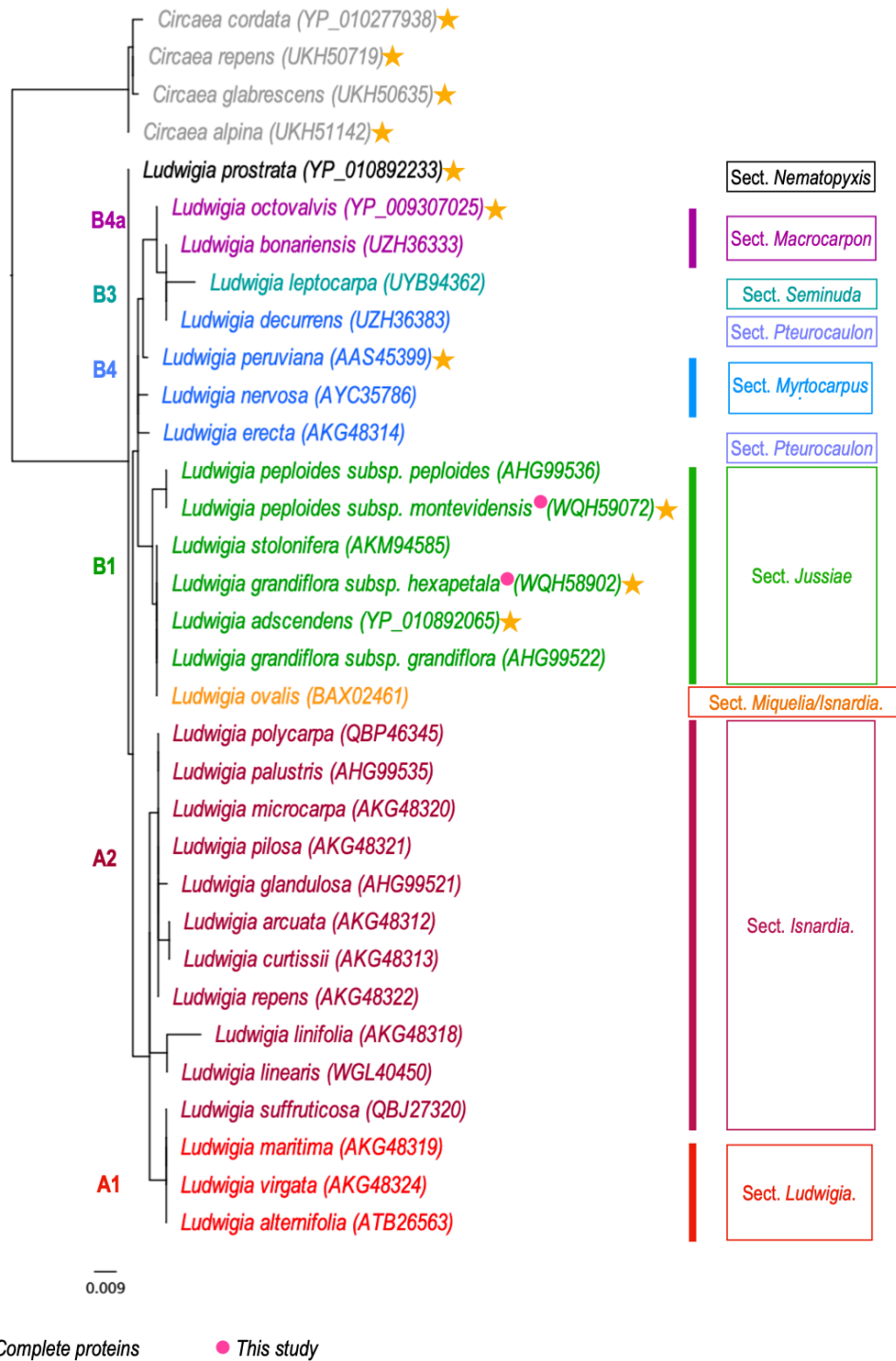


Figure 8 - Phylogenetic tree based on Ludwigia MatK protein sequences. Only six Ludwigia sequences are complete (yellow star), the others correspond to amino acids ranging from 128 to 289 aa, with an average of 244 aa. Clades are named and colored regarding the Ludwigia phylogeny proposed by Liu et al. (2017). The sections are based on the works of Raven (1963), Wagner et al. (2007) and Liu et al. (2023). The scale bar indicates the branch length.

Discussion

In the present study, we first sequenced and de novo assembled the chloroplast (cp) genomes of *Ludwigia peploides* subsp. *montevidensis* (*Lpm*) and *Ludwigia grandiflora* subsp. *hexapetala* (*Lgh*), two species belonging to the Onagraceae family. We employed a hybrid strategy and demonstrated the presence of two cp haplotypes in *Lgh* and one haplotype in *Lpm*, although the presence of both haplotypes in *Lpm* is likely. Furthermore, we compared these genomes with those of other species in the Onagraceae family to expand our knowledge of genome organization and molecular evolution in these species.

Our findings demonstrate that the utilization of solely short reads has failed to produce complete *Ludwigia* plastomes, likely due to challenges posed by long repeats and rearrangements. On the other hand, relying solely on long reads resulted in a lower quality sequence due to insufficient coverage and sequencing errors. After conducting our research, we discovered that, for *Lgh* plastomes, hybrid assembly, which incorporates both long and short read sequences, resulted in the most superior complete assemblies. This innovative approach capitalizes on the advantages of both sequencing technologies, harnessing the accuracy of short read sequences and the length of long read sequences. In the case of our study on *Lgh* plastome reconstruction, hybrid assembly was the most complete and effective, similarly to studies on other chloroplasts, such as those in *Eucalyptus* (Wang et al. 2018), *Falcataria* (Anita et al. 2023), *Carex* (Xu et al. 2023) or *Cypripedium* (Guo et al. 2021).

In our study, we were able to identify the presence of two haplotypes in *Lgh*, which is a first for *Ludwigia* (and more broadly within Onagraceae), as the plastome of *L. octovalvis* was only delivered in one haplotype (Liu et al. 2016).

Due to the unavailability of sequence data for *Ludwigia octovalvis* and the fact that we only have long reads for *Ludwigia peploides* subsp. *montevidensis*, none of which large enough to cover the SSC/IR junctions, we are unable to conclusively identify the presence of these two forms in the *Ludwigia* genus. However, we believe that they are likely to be present. Unfortunately, the current representation of plastomes in GenBank primarily consists of short-read data, which may result in an underrepresentation of this polymorphism. It is unfortunate that structural heteroplasmy, which is expected to be widespread in angiosperms, has been overlooked. Existence of two plastome haplotypes has been identified in the related order of Myrtales (*Eucalyptus* sp.), in 58 species of Angiosperms (Wang and Lanfear 2019), Asparagales (*Ophrys apifera* orchid, Bateman et al. 2021), Brassicales (*Carica papaya*, *Vasconcellea pubescens*, Lin et al. 2020), Solanales (*Solanum tuberosum*, Lihodeevskiy and Shanina 2022), Laurales (*Avocado Persea americana*, Nath et al. 2022) and Rhamnaceae (*Rhamnus crenata*, Wanichthanarak et al. 2023). However, the majority of reference plastomes in the current GenBank database (Release 260: April 15, 2024) are described as a single haplotype, indicating an underrepresentation of structural heteroplasmy in angiosperm chloroplasts. This underscores the importance of sequencing techniques, as the database is predominantly composed of short-read data (98%), which are less effective than long reads or hybrid assemblies at detecting flip-flop phenomena in the LSC region.

The chloroplast genome sizes for the three genera of Onagraceae subfam. Onagroideae varied as follows: *Circaea* sp. ranged from 155,817 bp to 156,024 bp, *Chamaenerion* sp. ranged from 159,496 bp to 160,416 bp, and *Epilobium* sp. ranged from 160,748 bp to 161,144 bp (Luo et al. 2021). Our study revealed that the size of the complete chloroplast of *Ludwigia* (Onagraceae subfamily Ludwigioideae) ranged from 159,369 bp to 159,584 bp, which is remarkably similar to other Onagraceae plants (average length of 162,030 bp). Furthermore, *Ludwigia* plastome sizes are consistent with the range observed in Myrtales (between 152,214 to 171,315 bp, Zhang et al. 2021). In the same way, similar overall GC content was found in *Ludwigia* sp. (from 37.3 to 37.4%), *Circaea* sp. (37.7 to 37.8%), *Chamaenerion* sp. and *Epilobium* sp. (38.1 to 38.2%, Luo et al. 2021) and more generally for the order Myrtales (36.9–38.9%, with the average GC content being 37%, (Zhang et al. 2021)). Higher GC content of the IR regions (43.5%) found in *Ludwigia* sp. has already been shown in the Myrtales order (39.7–43.5%) and in other families/orders such as

Amaranthaceae (order Caryophyllales, Xu et al. 2020) or Lamiaceae (order Lamiales, Lian et al. 2022), and is mainly due to the presence of the four GC rich rRNA genes.

The complete chloroplast genomes of the three *Ludwigia* species encoded an identical set of 134 genes including 85 protein-coding genes, 37 tRNA genes and eight ribosomal RNAs, consistent with gene content found in the Myrtales order, with a gene number varying from 123 to 133 genes with 77–81 protein-coding genes, 29–31 tRNA gene and four rRNA genes (Zhang et al. 2021). Chloroplast genes have been selected during evolution due to their functional importance (Mohanta et al. 2020). In our current study, we made the noteworthy discovery that *matK*, *accD*, *ycf2*, and *ccsA* genes were subjected to positive selection pressure. These genes have frequently been reported in literature as being associated with positive selection, and are known to play crucial roles in plant development conditions. *Lgh* and *Lpm* are known to thrive in aquatic environments, where they grow alongside rooted emergent aquatic plants, with their leaves and stems partially submerged during growth, as reported by Wagner et al. (2007). Both species possess the unique ability of vegetative reproduction, enabling them to establish themselves rapidly in diverse habitats, including terrestrial habitats, as noted by Haury et al. (2014b). Additionally, *Lo* is a wetland plant that typically grows in gullies and at the edges of ponds, as documented by Wagner et al. (2007). Given their ability to adapt to different habitats, these species may have evolved specialized mechanisms to cope with various abiotic stresses, such as reduced carbon and oxygen availability or limited access to light in submerged or emergent conditions. Concerning *matK*, Barhet and Hilu (2007) demonstrated the relationship between light and developmental stages, and *MatK* maturase activity, suggesting important functions in plant physiology. This gene has recently been largely reported to be under positive selection in an aquatic plant (*Anubias* sp., Li L et al. 2022), and more generally in terrestrial plants (*Pinus* sp., Zeb et al. 2022) or *Chrysosplenium* sp., Wu et al. 2020)). The *accD* gene has been described as an essential gene required for leaf development (Kode et al. 2005) and longevity in tobacco (*Nicotiana tabacum*, Madoka et al. 2002). Under drought stress, plant resistance can be increased by inhibiting *accD* (Gu et al. 2020), and conversely, enhanced in response to flooding stress by upregulating *accD* accumulation (Bharadwaj et al. 2023). Hence, we can hypothesize that the positive selection observed on the *accD* gene can be explained by the submerged and emerged constraints undergone by *Ludwigia* species. The *ycf2* gene seems to be subject to adaptive evolution in *Ludwigia* species. Its function, although still vague, would be to contribute to a protein complex generating ATP for the TIC machinery (proteins importing into the chloroplasts (Kikuchi et al. 2018, Schreier et al. 2018), as well as plant cell survival (Drescher et al. 2000, Xing et al. 2022). The *ccsA* gene positive selection is found in some aquatic plants such as *Anubia* sp. (Li L et al. 2022), marine flowering plants as *Zostera* species (Chen et al. 2023), and some species of Lythraceae (Gu et al. 2020). The *ccsA* gene is required for cytochrome c biogenesis (Xie and Merchant 1996) and this hemoprotein plays a key role in aerobic and anaerobic respiration, as well as photosynthesis (Kranz et al. 1998). Furthermore, we showed that *Lgh* colonization is supported by metabolic adjustments mobilizing glycolysis and fermentation pathways in terrestrial habitats, and the aminoacyl-tRNA biosynthesis pathway, which are key components of protein synthesis in aquatic habitats (Billet et al. 2018). It can be assumed that the ability of *Ludwigia* to invade aquatic and wet environments, where the amount of oxygen and light can be variable, leads to a high selective pressure on genes involved in respiration and photosynthesis.

Molecular markers are often used to establish population genetic relationships through phylogenetic studies. Five chloroplasts (*rps16*, *rpl16*, *trnL-trnF*, *trnL-CD*, *trnG*) and two nuclear markers (ITS, *waxy*) were used in previous phylogeny studies of *Ludwigia* sp. (Liu et al. 2017). However, no SSR markers had previously been made available for the *Ludwigia* genus, or more broadly, the Onagraceae. In this study, we identified 45 to 65 SSR markers depending on the *Ludwigia* species. Most of them were AT mononucleotides, as already recorded for other angiosperms (Maheswari et al. 2021, Zhang et al. 2016). In addition, we identified various genes with highly mutated regions that can also be used as SNP markers. Chloroplast SSRs (cpSSRs) represent potentially useful markers showing high levels of intraspecific variability due to the non-recombinant and uniparental inheritance of the plastomes (Huang et al. 2018, Leontaritou et al.

2021). Chloroplast SSR characteristics for *Ludwigia* sp. (location, type of SSR) were similar to those described in most plants. While the usual molecular markers used for phylogenetic analysis are nuclear DNA markers, cpSSRs have also been used to explore cytoplasmic diversity in many studies (Snoussi et al. 2022, Song et al. 2014, Wheeler et al. 2014). To conclude, the 13 highly variable loci and cpSSRs identified in this study are potential markers for population genetics or phylogenetic studies of *Ludwigia* species, and more generally, Onagraceae.

Concerning the *MatK*-based phylogenetic tree, its topology is generally congruent with the first molecular classification of Liu et al. (2017) as all *Ludwigia* from sect. *Jussiaea* (clade B1) and sect. *Ludwigia* (clade A1) and sect. *Isnardia* (clade A2) branched together. In this *MatK*-based tree, *Ludwigia prostrata*, a species absent from previously published phylogenetic studies, positions itself alone at the root of the *Ludwigia* tree. This species, sole member of section Nematopyxis, is related as having no close relatives (Raven and Tai 1979), finding supported by our work. We also observed that *Ludwigia ovalis* branches within sect. *Jussiaea*, as its 258 amino acids partial *MatK* sequence (ca. half of the complete sequence) is identical to the *MatK* proteins of *L. grandiflora* subsp. *hexapetala*, *L. stolonifera* and *L. adscendens*. Its phylogenetic placement remains unresolved: classified alone by Raven (1963) and Wagner et al. (2007) in sect. *Miquelia*, later positioned by Liu et al. (2017) within the *Isnardia-Microcarpium* section (using nuclear DNA) or as sister to it (using plastid DNA). For this reason, conducting a whole plastome analysis would be valuable to provide insights into *L. ovalis* phylogenetic positioning. Another species positioned on the margins of sect. *Isnardia* (clade A2) is *Ludwigia suffruticosa* (previously classified in sect. *Microcarpium*), which branches within sect. *Ludwigia* (clade A1). This positioning raises questions about the current grouping of sections *Isnardia*, *Michelia*, and *Microcarpium* into a single section *Isnardia* as proposed by Liu et al. (2023) and highlights that plastid protein coding markers can provide differing phylogenetic insights. Finally, the last species positioned differently of this clade (clade B4) is *Ludwigia decurrens* (sect. *Pterocaulon*) which clusters with *L. leptocarpa* (clade B3) and *L. bonariensis* (clade B4a). However, it is important to note that in their study, Liu et al. (2017) indicate that clade B4 is moderately supported and that the two members of sect. *Pterocaulon*, *L. decurrens* and *L. nervosa*, diverge in all trees (Liu et al. 2017). In summary, acquiring complete plastomes for *Ludwigia* sp. could significantly enhance our understanding of the phylogeny of this complex genus. Furthermore, comparing nuclear and plastid phylogenies would help determine if they reflect the same evolutionary history and whether plastid phylogeny alone can accurately reconstruct the phylogeny of *Ludwigia* genus.

Conclusion

In this study, we conducted the first-time sequencing and assembly of the complete plastomes of *Lpm* and *Lgh*, which are the only available genomic resources for functional analysis in both species. We were able to identify the existence of two haplotypes in *Lgh*, but further investigations will be necessary to confirm their presence in *Lo* and *Lpm*, and more broadly, within the *Ludwigia* genus. Comparison of all 10 Onagraceae plastomes revealed a high degree of conservation in genome size, gene number, structure, and IR boundaries. However, to further elucidate the phylogenetic analysis and evolution in *Ludwigia* and Onagraceae, additional chloroplast genomes will be necessary, as highlighted in recent studies of *Iris* and *Aristidoideae* species (Feng et al. 2022).

Availability of data and materials

The datasets generated and/or analysed during the current study were available in GenBank (for *Lgh* haplotype 1, (LGH1) OR166254 and *Lgh* haplotype 2, (LGH2) OR166255; for *Lpm* haplotype, (LPM) OR166256). Chloroplastic short and long reads are available at EBI-ENA database (<https://www.ebi.ac.uk/ena/browser/home>) under these accession numbers for LGH plastomes (Long reads: Experiment: ERX13439011 ; Run: ERR14035997 and short reads:

Experiment: ERX13439002 ; Run: ERR14035988) and for LPM plastomes (Long reads: Experiment: ERX13439014 ; Run: ERR14036000).

Conflict of interest disclosure

The authors declare that they comply with the PCI rule of having no financial conflicts of interest in relation to the content of the article.

Funding

The post-doctoral research grant of Anne-Laure Le Gac was supported by the Conseil regional Bretagne (SAD18001).

Acknowledgements

Preprint version v4 of this article has been peer-reviewed and recommended by Peer CommunityIn Genomics (<https://doi.org/10.24072/pci.genomics.100334>; Sabot, 2025). We are grateful to Luis Portillo-Lemus for developing the high molecular weight genomic DNA extraction protocol. All sequencing experiments were performed at the PGTB (<https://doi.org/10.15454/1.5572396583599417E12>).

References

- Amiryousefi, Ali, Jaakko Hyvönen, and Peter Poczai. 2018. "IRscope: An Online Program to Visualize the Junction Sites of Chloroplast Genomes." *Bioinformatics* 34 (17). <https://doi.org/10.1093/bioinformatics/bty220>.
- Anita, V P D, D D Matra, and U J Siregar. 2023. "Chloroplast Genome Draft Assembly of *Falcataria Moluccana* Using Hybrid Sequencing Technology." *BMC Res Notes* 16 (1): 31. <https://doi.org/10.1186/s13104-023-06290-6>.
- Antil, S, J S Abraham, S Sripoorna, S Maurya, J Dagar, S Makhija, P Bhagat, et al. 2023. "DNA Barcoding, an Effective Tool for Species Identification: A Review." *Mol Biol Rep* 50 (1): 761–75. <https://doi.org/10.1007/s11033-022-08015-7>.
- Bankevich, Anton, Sergey Nurk, Dmitry Antipov, Alexey A. Gurevich, Mikhail Dvorkin, Alexander S. Kulikov, Valery M. Lesin, et al. 2012. "SPAdes: A New Genome Assembly Algorithm and Its Applications to Single-Cell Sequencing." *Journal of Computational Biology* 19 (5). <https://doi.org/10.1089/cmb.2012.0021>.
- Barloy, Dominique, Luis Portillo-Lemus, Stacy Krueger-Hadfield, Virginie Huteau, and Olivier Coriton. 2024. "Genomic Relationships among Diploid and Polyploid Species of the Genus *Ludwigia* L. Section *Jussiaea* Using a Combination of Molecular Cytogenetic, Morphological, and Crossing Investigations." *Peer Community Journal* 4. <https://doi.org/10.24072/pcjournal.364>.
- Barthel, Michelle M, and Khidir W Hilu. 2007. "Expression of MatK: Functional and Evolutionary Implications." *American Journal of Botany* 94 (8): 1402–12.
- Bateman, R M, P J Rudall, A R M Murphy, R S Cowan, D S Devey, and O A Perez-Escobar. 2021. "Whole Plastomes Are Not Enough: Phylogenomic and Morphometric Exploration at Multiple Demographic Levels of the Bee Orchid Clade *Ophrys* Sect. *Sphegodes*." *J Exp Bot* 72 (2): 654–81. <https://doi.org/10.1093/jxb/eraa467>.
- Bedoya, Ana M., and Santiago Madriñán. 2015. "Evolution of the Aquatic Habit in *Ludwigia* (Onagraceae): Morpho-Anatomical Adaptive Strategies in the Neotropics." *Aquatic Botany* 120 (PB). <https://doi.org/10.1016/j.aquabot.2014.10.005>.
- Beier, Sebastian, Thomas Thiel, Thomas Münch, Uwe Scholz, and Martin Mascher. 2017. "MISA-Web: A Web Server for Microsatellite Prediction." *Bioinformatics* 33 (16). <https://doi.org/10.1093/bioinformatics/btx198>.

- Belser, Caroline, Benjamin Istace, Erwan Denis, Marion Dubarry, Franc Christophe Baurens, Cyril Falentin, Mathieu Genete, et al. 2018. "Chromosome-Scale Assemblies of Plant Genomes Using Nanopore Long Reads and Optical Maps." *Nature Plants*. <https://doi.org/10.1038/s41477-018-0289-4>.
- Benson, Gary. 1999. "Tandem Repeats Finder: A Program to Analyze DNA Sequences." *Nucleic Acids Research* 27 (2). <https://doi.org/10.1093/nar/27.2.573>.
- Bharadwaj, Bhiolina, Avetis Mishegyan, Sanjeevi Nagalingam, Alex Guenther, Nirmal Joshee, Samantha H Sherman, and Chhandak Basu. 2023. "Physiological and Genetic Responses of Lentil (*Lens Culinaris*) under Flood Stress." *Plant Stress*, 100130.
- Billet, K, J Genitoni, M Bozec, D Renault, and D Barloy. 2018. "Aquatic and Terrestrial Morphotypes of the Aquatic Invasive Plant, *Ludwigia Grandiflora*, Show Distinct Morphological and Metabolomic Responses." *Ecol Evol* 8 (5): 2568–79. <https://doi.org/10.1002/ece3.3848>.
- Billet, Kevin, Julien Genitoni, Michel Bozec, David Renault, and Dominique Barloy. 2018. "Aquatic and Terrestrial Morphotypes of the Aquatic Invasive Plant, *Ludwigia Grandiflora*, Show Distinct Morphological and Metabolomic Responses." *Ecology and Evolution* 8 (5). <https://doi.org/10.1002/ece3.3848>.
- Brudno, Michael, Chuong B. Do, Gregory M. Cooper, Michael F. Kim, Eugene Davydov, Eric D. Green, Arend Sidow, and Serafim Batzoglou. 2003. "LAGAN and Multi-LAGAN: Efficient Tools for Large-Scale Multiple Alignment of Genomic DNA." *Genome Research*. <http://www.genome.org/cgi/doi/10.1101/gr.926603>.
- Chen, Chengjie, Hao Chen, Yi Zhang, Hannah R. Thomas, Margaret H. Frank, Yehua He, and Rui Xia. 2020. "TBtools: An Integrative Toolkit Developed for Interactive Analyses of Big Biological Data." *Molecular Plant* 13 (8). <https://doi.org/10.1016/j.molp.2020.06.009>.
- Chen, J, Y Zang, S Shang, Z Yang, S Liang, S Xue, Y Wang, and X Tang. 2023. "Chloroplast Genomic Comparison Provides Insights into the Evolution of Seagrasses." *BMC Plant Biol* 23 (1): 104. <https://doi.org/10.1186/s12870-023-04119-9>.
- Chen, Shifu, Yanqing Zhou, Yaru Chen, and Jia Gu. 2018. "Fastp: An Ultra-Fast All-in-One FASTQ Preprocessor." In *Bioinformatics*. Vol. 34. <https://doi.org/10.1093/bioinformatics/bty560>.
- Chikhi, Rayan, and Paul Medvedev. 2014. "Informed and Automated K-Mer Size Selection for Genome Assembly." *Bioinformatics* 30 (1). <https://doi.org/10.1093/bioinformatics/btt310>.
- Dandelot, Sophie, Régine Verlaque, Alain Dutartre, and Arlette Cazaubon. 2005. "Ecological, Dynamic and Taxonomic Problems Due to *Ludwigia* (Onagraceae) in France." In *Hydrobiologia*. Vol. 551. <https://doi.org/10.1007/s10750-005-4455-0>.
- Dierckxsens, Nicolas, Patrick Mardulyn, and Guillaume Smits. 2017. "NOVOPlasty: De Novo Assembly of Organelle Genomes from Whole Genome Data." *Nucleic Acids Research* 45 (4). <https://doi.org/10.1093/nar/gkw955>.
- Drescher, A, S Ruf, T Calsa Jr., H Carrer, and R Bock. 2000. "The Two Largest Chloroplast Genome-Encoded Open Reading Frames of Higher Plants Are Essential Genes." *Plant J* 22 (2): 97–104. <https://doi.org/10.1046/j.1365-313x.2000.00722.x>.
- Feng, Jing Lu, Li Wei Wu, Qing Wang, Yun Jia Pan, Bao Li Li, Yu Lin Lin, and Hui Yao. 2022. "Comparison Analysis Based on Complete Chloroplast Genomes and Insights into Plastid Phylogenomic of Four Iris Species." *BioMed Research International* 2022. <https://doi.org/10.1155/2022/2194021>.
- Ferrarini, Marco, Marco Moretto, Judson A. Ward, Nada Šurbanovski, Vladimir Stevanović, Lara Giongo, Roberto Viola, et al. 2013. "An Evaluation of the PacBio RS Platform for Sequencing and de Novo Assembly of a Chloroplast Genome." *BMC Genomics* 14 (1). <https://doi.org/10.1186/1471-2164-14-670>.
- Filip, E, and L Skuza. 2021. "Horizontal Gene Transfer Involving Chloroplasts." *Int J Mol Sci* 22 (9). <https://doi.org/10.3390/ijms22094484>.
- Frazer, Kelly A., Lior Pachter, Alexander Poliakov, Edward M. Rubin, and Inna Dubchak. 2004. "VISTA: Computational Tools for Comparative Genomics." *Nucleic Acids Research* 32 (WEB SERVER ISS.). <https://doi.org/10.1093/nar/gkh458>.

- Gioria, Margherita, Philip E Hulme, David M Richardson, and Petr Pyšek. 2023. "Annual Review of Plant Biology Why Are Invasive Plants Successful?" *Annu. Rev. Plant Biol.* 2023 74:2023. <https://doi.org/10.1146/annurev-arplant-070522>.
- Greiner, Stephan, Pascal Lehwark, and Ralph Bock. 2019. "OrganellarGenomeDRAW (OGDRAW) Version 1.3.1: Expanded Toolkit for the Graphical Visualization of Organellar Genomes." *Nucleic Acids Research* 47 (W1). <https://doi.org/10.1093/nar/gkz238>.
- Grewell, Brenda J., M. D. Netherland, and M. J. Skaer Thomason. 2016. "Establishing Research and Management Priorities for Invasive Water Primroses (*Ludwigia* Spp.)." Aquatic Plant Control Research Program, US Army Corps of Engineers, Engineer Research and Development Center, Environmental Laboratory Technical Report ERDC/ELTR-15-X, no. February.
- Gu, H, Y Wang, H Xie, C Qiu, S Zhang, J Xiao, H Li, L Chen, X Li, and Z Ding. 2020. "Drought Stress Triggers Proteomic Changes Involving Lignin, Flavonoids and Fatty Acids in Tea Plants." *Sci Rep* 10 (1): 15504. <https://doi.org/10.1038/s41598-020-72596-1>.
- Gualberto, José M., Daria Mileshina, Clémentine Wallet, Adnan Khan Niazi, Frédérique Weber-Loffi, and André Dietrich. 2014. "The Plant Mitochondrial Genome: Dynamics and Maintenance." *Biochimie*. <https://doi.org/10.1016/j.biochi.2013.09.016>.
- Guo, Y Y, J X Yang, H K Li, and H S Zhao. 2021. "Chloroplast Genomes of Two Species of *Cypripedium*: Expanded Genome Size and Proliferation of AT-Biased Repeat Sequences." *Front Plant Sci* 12:609729. <https://doi.org/10.3389/fpls.2021.609729>.
- Gurdon, C, and P Maliga. 2014. "Two Distinct Plastid Genome Configurations and Unprecedented Intraspecies Length Variation in the *accD* Coding Region in *Medicago truncatula*." *DNA Res* 21 (4): 417–27. <https://doi.org/10.1093/dnares/dsu007>.
- Gurevich, Alexey, Vladislav Saveliev, Nikolay Vyahhi, and Glenn Tesler. 2013. "QUAST: Quality Assessment Tool for Genome Assemblies." *Bioinformatics* 29 (8): 1072–75. <https://doi.org/10.1093/bioinformatics/btt086>.
- Gurusaran, M., Dheeraj Ravella, and K. Sekar. 2013. "RepEx: Repeat Extractor for Biological Sequences." *Genomics* 102 (4): 403–8. <https://doi.org/10.1016/j.ygeno.2013.07.005>.
- Haury, Jacques, Arsène Druel, Teipotemarama Cabral, Yann Paulet, Michel Bozec, and Julie Coudreuse. 2014a. "Which Adaptations of Some Invasive *Ludwigia* Spp. (Rosidae, Onagraceae) Populations Occur in Contrasting Hydrological Conditions in Western France?" *Hydrobiologia* 737 (1). <https://doi.org/10.1007/s10750-014-1815-7>.
- Haury, Jacques, Arsène Druel, Teipotemarama Cabral, Yann Paulet, Michel Bozec, and Julie Coudreuse. 2014b. "Which Adaptations of Some Invasive *Ludwigia* Spp. (Rosidae, Onagraceae) Populations Occur in Contrasting Hydrological Conditions in Western France?" *Hydrobiologia* 737:45–56.
- Hoch, Peter C., Warren L. Wagner, and Peter H. Raven. 2015. "The Correct Name for a Section of *Ludwigia* L. (Onagraceae)." *PhytoKeys* 50 (1). <https://doi.org/10.3897/phytokeys.50.4887>.
- Holley, Guillaume, Doruk Beyter, Helga Ingimundardottir, Peter L. Møller, Snædis Kristmundsdottir, Hannes P. Eggertsson, and Bjarni V. Halldorsson. 2021. "Ratatosk: Hybrid Error Correction of Long Reads Enables Accurate Variant Calling and Assembly." *Genome Biology* 22 (1). <https://doi.org/10.1186/s13059-020-02244-4>.
- Hu, Yingchun, Quan Zhang, Guangyuan Rao, and Sodmergen. 2008. "Occurrence of Plastids in the Sperm Cells of Caprifoliaceae: Biparental Plastid Inheritance in Angiosperms Is Unilaterally Derived from Maternal Inheritance." *Plant and Cell Physiology* 49 (6). <https://doi.org/10.1093/pcp/pcn069>.
- Huang, L S, Y Q Sun, Y Jin, Q Gao, X G Hu, F L Gao, X L Yang, J J Zhu, Y A El-Kassaby, and J F Mao. 2018. "Development of High Transferability CpSSR Markers for Individual Identification and Genetic Investigation in Cupressaceae Species." *Ecol Evol* 8 (10): 4967–77. <https://doi.org/10.1002/ece3.4053>.
- Hussner, A., M. Windhaus, and U. Starfinger. 2016. "From Weed Biology to Successful Control: An Example of Successful Management of *Ludwigia grandiflora* in Germany." *Weed Research* 56 (6). <https://doi.org/10.1111/wre.12224>.

- Jackman, Shaun D., Benjamin P. Vandervalk, Hamid Mohamadi, Justin Chu, Sarah Yeo, S. Austin Hammond, Golnaz Jahesh, et al. 2017. "ABYSS 2.0: Resource-Efficient Assembly of Large Genomes Using a Bloom Filter." *Genome Research* 27 (5). <https://doi.org/10.1101/gr.214346.116>.
- Jain, Miten, Sergey Koren, Karen H. Miga, Josh Quick, Arthur C. Rand, Thomas A. Sasani, John R. Tyson, et al. 2018. "Nanopore Sequencing and Assembly of a Human Genome with Ultra-Long Reads." *Nature Biotechnology* 36 (4). <https://doi.org/10.1038/nbt.4060>.
- Jin, Jian Jun, Wen Bin Yu, Jun Bo Yang, Yu Song, Claude W. Depamphilis, Ting Shuang Yi, and De Zhu Li. 2020. "GetOrganelle: A Fast and Versatile Toolkit for Accurate de Novo Assembly of Organelle Genomes." *Genome Biology* 21 (1). <https://doi.org/10.1186/s13059-020-02154-5>.
- Jones, K., and R. E. Cleland. 1974. "Oenothera, Cytogenetics and Evolution." *Kew Bulletin* 29 (1). <https://doi.org/10.2307/4108389>.
- Kamoshita, Akihiko, Hiroyuki Ikeda, Junko Yamagishi, Bunna Lor, and Makara Ouk. 2016. "Residual Effects of Cultivation Methods on Weed Seed Banks and Weeds in Cambodia." *Weed Biology and Management* 16 (3): 93–107.
- Katoh, Kazutaka, Kazuharu Misawa, Kei Ichi Kuma, and Takashi Miyata. 2002. "MAFFT: A Novel Method for Rapid Multiple Sequence Alignment Based on Fast Fourier Transform." *Nucleic Acids Research* 30 (14). <https://doi.org/10.1093/nar/gkf436>.
- Kikuchi, S, Y Asakura, M Imai, Y Nakahira, Y Kotani, Y Hashiguchi, Y Nakai, et al. 2018. "A Ycf2-FtsHi Heteromeric AAA-ATPase Complex Is Required for Chloroplast Protein Import." *Plant Cell* 30 (11): 2677–2703. <https://doi.org/10.1105/tpc.18.00357>.
- Kode, V, E A Mudd, S lamtham, and A Day. 2005. "The Tobacco Plastid AccD Gene Is Essential and Is Required for Leaf Development." *Plant J* 44 (2): 237–44. <https://doi.org/10.1111/j.1365-313X.2005.02533.x>.
- Kolmogorov, Mikhail, Jeffrey Yuan, Yu Lin, and Pavel A. Pevzner. 2019. "Assembly of Long, Error-Prone Reads Using Repeat Graphs." *Nature Biotechnology* 37 (5). <https://doi.org/10.1038/s41587-019-0072-8>.
- Konishi, T, and Y Sasaki. 1994. "Compartmentalization of Two Forms of Acetyl-CoA Carboxylase in Plants and the Origin of Their Tolerance toward Herbicides." *Proc Natl Acad Sci U S A* 91 (9): 3598–3601. <https://doi.org/10.1073/pnas.91.9.3598>.
- Koren, Sergey, Brian P. Walenz, Konstantin Berlin, Jason R. Miller, Nicholas H. Bergman, and Adam M. Phillippy. 2017. "Canu: Scalable and Accurate Long-Read Assembly via Adaptive k-Mer Weighting and Repeat Separation." *Genome Research* 27 (5). <https://doi.org/10.1101/gr.215087.116>.
- Kranz, R, R Lill, B Goldman, G Bonnard, and S Merchant. 1998. "Molecular Mechanisms of Cytochrome c Biogenesis: Three Distinct Systems." *Mol Microbiol* 29 (2): 383–96. <https://doi.org/10.1046/j.1365-2958.1998.00869.x>.
- Lambert, E., A. Dutartre, J. Coudreuse, and J. Haury. 2010. "Relationships between the Biomass Production of Invasive Ludwigia Species and Physical Properties of Habitats in France." *Hydrobiologia* 656 (1). <https://doi.org/10.1007/s10750-010-0440-3>.
- Lehwark, Pascal, and Stephan Greiner. 2019. "GB2sequin - A File Converter Preparing Custom GenBank Files for Database Submission." *Genomics* 111 (4). <https://doi.org/10.1016/j.ygeno.2018.05.003>.
- Leontaritou, P, F N Lamari, V Papatotiropoulos, and G Iatrou. 2021. "Exploration of Genetic, Morphological and Essential Oil Variation Reveals Tools for the Authentication and Breeding of *Salvia Pomifera* Subsp. *Calycina* (Sm.) Hayek." *Phytochemistry* 191:112900. <https://doi.org/10.1016/j.phytochem.2021.112900>.
- Levin, Rachel A., Warren L. Wagner, Peter C. Hoch, William J. Hahn, Aaron Rodriguez, David A. Baum, Liliana Katinas, Elizabeth A. Zimmer, and Kenneth J. Sytsma. 2004. "Paraphyly in Tribe Onagreae: Insights into Phylogenetic Relationships of Onagraceae Based on Nuclear and Chloroplast Sequence Data." *Systematic Botany* 29 (1). <https://doi.org/10.1600/036364404772974293>.

- Levin, Rachel A., Warren L. Wagner, Peter C. Hoch, Molly Nepokroeff, J. Chris Pires, Elizabeth A. Zimmer, and Kenneth J. Sytsma. 2003. "Family-Level Relationships of Onagraceae Based on Chloroplast RbcL and NdhF Data." *American Journal of Botany* 90 (1). <https://doi.org/10.3732/ajb.90.1.107>.
- Li, Dinghua, Ruibang Luo, Chi Man Liu, Chi Ming Leung, Hing Fung Ting, Kunihiko Sadakane, Hiroshi Yamashita, and Tak Wah Lam. 2016. "MEGAHIT v1.0: A Fast and Scalable Metagenome Assembler Driven by Advanced Methodologies and Community Practices." *Methods*. <https://doi.org/10.1016/j.ymeth.2016.02.020>.
- Li, J, Y Su, and T Wang. 2018. "The Repeat Sequences and Elevated Substitution Rates of the Chloroplast AccD Gene in Cupressophytes." *Front Plant Sci* 9:533. <https://doi.org/10.3389/fpls.2018.00533>.
- Li, J, Y Wang, J Cui, W Wang, X Liu, Y Chang, D Yao, and J Cui. 2022. "Removal Effects of Aquatic Plants on High-Concentration Phosphorus in Wastewater during Summer." *J Environ Manage* 324:116434. <https://doi.org/10.1016/j.jenvman.2022.116434>.
- Li, L, C Liu, K Hou, and W Liu. 2022. "Comparative Analyses of Plastomes of Four Anubias (Araceae) Taxa, Tropical Aquatic Plants Endemic to Africa." *Genes (Basel)* 13 (11). <https://doi.org/10.3390/genes13112043>.
- Lian, C, H Yang, J Lan, X Zhang, F Zhang, J Yang, and S Chen. 2022. "Comparative Analysis of Chloroplast Genomes Reveals Phylogenetic Relationships and Intraspecific Variation in the Medicinal Plant *Isodon Rubescens*." *PLoS One* 17 (4): e0266546. <https://doi.org/10.1371/journal.pone.0266546>.
- Lihodeevskiy, Georgiy A, and Elena P Shanina. 2022. "The Use of Long-Read Sequencing to Study the Phylogenetic Diversity of the Potato Varieties Plastome of the Ural Selection." *Agronomy* 12 (4): 846.
- Lin, Z, P Zhou, X Ma, Y Deng, Z Liao, R Li, and R Ming. 2020. "Comparative Analysis of Chloroplast Genomes in *Vasconcellea Pubescens* A.DC. and *Carica Papaya* L." *Sci Rep* 10 (1): 15799. <https://doi.org/10.1038/s41598-020-72769-y>.
- Liu, Shih Hui, Christine Edwards, Peter C. Hoch, Peter H. Raven, and Janet C. Barber. 2016. "Complete Plastome Sequence of *Ludwigia Octovalvis* (Onagraceae), a Globally Distributed Wetland Plant." *Genome Announcements* 4 (6). <https://doi.org/10.1128/genomeA.01274-16>.
- Liu, Shih Hui, Peter C. Hoch, Mauricio Diazgranados, Peter H. Raven, and Janet C. Barber. 2017. "Multi-Locus Phylogeny of *Ludwigia* (Onagraceae): Insights on Infra- Generic Relationships and the Current Classification of the Genus." *Taxon* 66 (5). <https://doi.org/10.12705/665.7>.
- Liu, Shih Hui, Kuo Hsiang Hung, Tsai Wen Hsu, Peter C. Hoch, Ching I. Peng, and Tzen Yuh Chiang. 2023. "New Insights into Polyploid Evolution and Dynamic Nature of *Ludwigia* Section *Isnardia* (Onagraceae)." *Botanical Studies* 64 (1). <https://doi.org/10.1186/s40529-023-00387-8>.
- Liu, Shih Hui, Hsun An Yang, Yoshiko Kono, Peter C. Hoch, Janet C. Barber, Ching I. Peng, and Kuo Fang Chung. 2020. "Disentangling Reticulate Evolution of North Temperate Haplostemonous *Ludwigia* (Onagraceae)." *Annals of the Missouri Botanical Garden* 105 (2). <https://doi.org/10.3417/2020479>.
- Liu, Shih-Hui, Christine Edwards, Peter C Hoch, Peter H Raven, and Janet C Barber. 2016. "Complete Plastome Sequence of *Ludwigia octovalvis* (Onagraceae), a Globally Distributed Wetland Plant." *Genome Announcements* 4 (6): e01274-16.
- Liu, Shih-Hui, Peter C Hoch, Mauricio Diazgranados, Peter H Raven, and Janet C Barber. 2017. "Multi-locus Phylogeny of *Ludwigia* (Onagraceae): Insights on Infra-generic Relationships and the Current Classification of the Genus." *Taxon* 66 (5): 1112–27.
- Luo, Y, J He, R Lyu, J Xiao, W Li, M Yao, L Pei, J Cheng, J Li, and L Xie. 2021. "Comparative Analysis of Complete Chloroplast Genomes of 13 Species in *Epilobium*, *Circaea*, and *Chamaenerion* and Insights Into Phylogenetic Relationships of Onagraceae." *Front Genet* 12:730495. <https://doi.org/10.3389/fgene.2021.730495>.
- Madoka, Y, K Tomizawa, J Mizoi, I Nishida, Y Nagano, and Y Sasaki. 2002. "Chloroplast Transformation with Modified AccD Operon Increases Acetyl-CoA Carboxylase and Causes

- Extension of Leaf Longevity and Increase in Seed Yield in Tobacco." *Plant Cell Physiol* 43 (12): 1518–25. <https://doi.org/10.1093/pcp/pcf172>.
- Maheswari, P, C Kunhikannan, and R Yasodha. 2021. "Chloroplast Genome Analysis of Angiosperms and Phylogenetic Relationships among Lamiaceae Members with Particular Reference to Teak (*Tectona Grandis* L.f)." *J Biosci* 46. <https://www.ncbi.nlm.nih.gov/pubmed/34047286>.
- Marks, Rose A., Scott Hotaling, Paul B. Frandsen, and Robert VanBuren. 2021. "Representation and Participation across 20 Years of Plant Genome Sequencing." *Nature Plants* 7 (12). <https://doi.org/10.1038/s41477-021-01031-8>.
- Mohanta, T K, A K Mishra, A Khan, A Hashem, E F Abd Allah, and A Al-Harrasi. 2020. "Gene Loss and Evolution of the Plastome." *Genes (Basel)* 11 (10). <https://doi.org/10.3390/genes11101133>.
- Moravcová, Lenka, Petr Pyšek, Vojtěch Jarošík, and Jan Pergl. 2015. "Getting the Right Traits: Reproductive and Dispersal Characteristics Predict the Invasiveness of Herbaceous Plant Species." *PLoS ONE* 10 (4). <https://doi.org/10.1371/journal.pone.0123634>.
- Nath, O, S J Fletcher, A Hayward, L M Shaw, A K Masouleh, A Furtado, R J Henry, and N Mitter. 2022. "A Haplotype Resolved Chromosomal Level Avocado Genome Allows Analysis of Novel Avocado Genes." *Hortic Res* 9:uhac157. <https://doi.org/10.1093/hr/uhac157>.
- Oldenburg, Delene J., and Arnold J. Bendich. 2016. "The Linear Plastid Chromosomes of Maize: Terminal Sequences, Structures, and Implications for DNA Replication." *Current Genetics* 62 (2). <https://doi.org/10.1007/s00294-015-0548-0>.
- Ou, Jianhong, and Lihua Julie Zhu. 2019. "TrackViewer: A Bioconductor Package for Interactive and Integrative Visualization of Multi-Omics Data." *Nature Methods*. Nature Publishing Group. <https://doi.org/10.1038/s41592-019-0430-y>.
- Panova, Marina, Henrik Aronsson, R. Andrew Cameron, Peter Dahl, Anna Godhe, Ulrika Lind, Olga Ortega-Martinez, et al. 2016. "DNA Extraction Protocols for Whole-Genome Sequencing in Marine Organisms." In *Methods in Molecular Biology*. Vol. 1452. https://doi.org/10.1007/978-1-4939-3774-5_2.
- RStudio Team. 2015. RStudio: Integrated Development Environment for R, RStudio, PBC. <http://www.rstudio.com/>
- Rang, Franka J., Wigard P. Kloosterman, and Jeroen de Ridder. 2018. "From Squiggle to Basepair: Computational Approaches for Improving Nanopore Sequencing Read Accuracy." *Genome Biology*. <https://doi.org/10.1186/s13059-018-1462-9>.
- Raven, Peter H, and William Tai. 1979. "Observations of Chromosomes in *Ludwigia* (Onagraceae)." Source: *Annals of the Missouri Botanical Garden*. Vol. 66. <https://about.jstor.org/terms>.
- Raven P.H. 1963. "The Old World species of *Ludwigia* including *Jussia*." *Reinwardtia* 6:327–427.
- Reddy, Angelica M, Paul D Pratt, Brenda J Grewell, Nathan E Harms, Guillermo Cabrera Walsh, M Cristina Hern, Ana Faltlhauser, and Ximena Cibils-stewart. 2021. "Biological Control of Invasive Water Primroses, *Ludwigia* Spp., in the United States: A Feasibility Assessment." *J. Aquat. Plant Manage.* Vol. 59.
- Richardson, A O, and J D Palmer. 2007. "Horizontal Gene Transfer in Plants." *J Exp Bot* 58 (1): 1–9. <https://doi.org/10.1093/jxb/erl148>.
- Rozas, J., and R. Rozas. 1999. "DnaSP Version 3: An Integrated Program for Molecular Population Genetics and Molecular Evolution Analysis." *Bioinformatics*. <https://doi.org/10.1093/bioinformatics/15.2.174>.
- Rozas, Julio, Albert Ferrer-Mata, J. C. Sanchez-DelBarrio, Sara Guirao-Rico, Pablo Librado, Sebastian E. Ramos-Onsins, and Alejandro Sanchez-Gracia. 2017. "DnaSP 6: DNA Sequence Polymorphism Analysis of Large Data Sets." *Molecular Biology and Evolution* 34 (12). <https://doi.org/10.1093/molbev/msx248>.
- Sabot, F. 2025. Onagre, monster, invasion and genetics. *Peer Community in Genomics*, 100334. 10.24072/pci.genomics.100334
- Sato, Naoki. 2021. "Are Cyanobacteria an Ancestor of Chloroplasts or Just One of the Gene Donors for Plants and Algae?" *Genes* 12 (6). <https://doi.org/10.3390/genes12060823>.

- Scheunert, Agnes, Marco Dorfner, Thomas Lingl, and Christoph Oberprieler. 2020. "Can We Use It? On the Utility of de Novo and Reference-Based Assembly of Nanopore Data for Plant Plastome Sequencing." *PLoS ONE* 15 (3). <https://doi.org/10.1371/journal.pone.0226234>.
- Schmitz, Udo K., and Klaus V. Kowallik. 1986. "Plastid Inheritance in *Epilobium*." *Current Genetics* 11 (1). <https://doi.org/10.1007/BF00389419>.
- Schreier, T B, A Clery, M Schlafli, F Galbier, M Stadler, E Demarsy, D Albertini, et al. 2018. "Plastidial NAD-Dependent Malate Dehydrogenase: A Moonlighting Protein Involved in Early Chloroplast Development through Its Interaction with an FtsH12-FtsHi Protease Complex." *Plant Cell* 30 (8): 1745–69. <https://doi.org/10.1105/tpc.18.00121>.
- Simpson, Jared T., Kim Wong, Shaun D. Jackman, Jacqueline E. Schein, Steven J.M. Jones, and Inanç Birol. 2009. "ABYSS: A Parallel Assembler for Short Read Sequence Data." *Genome Research* 19 (6). <https://doi.org/10.1101/gr.089532.108>.
- Snoussi, M, L Riahi, M Ben Romdhane, A Mliki, and N Zoghalmi. 2022. "Chloroplast DNA Diversity of Tunisian Barley Landraces as Revealed by CpSSRs Molecular Markers and Implication for Conservation Strategies." *Genet Res (Camb)* 2022:3905957. <https://doi.org/10.1155/2022/3905957>.
- Soliman, Ashraf T, Rim S Hamdy, and Azza B Hamed. 2018. "Ludwigia stolonifera (Guill. & Perr.) PH Raven, Insight into Its Phenotypic Plasticity, Habitat Diversity and Associated Species." *Egyptian Journal of Botany* 58 (3): 605–26.
- Song, S L, P E Lim, S M Phang, W W Lee, D D Hong, and A Prathep. 2014. "Development of Chloroplast Simple Sequence Repeats (CpSSRs) for the Intraspecific Study of *Gracilaria tenuistipitata* (Gracilariales, Rhodophyta) from Different Populations." *BMC Res Notes* 7:77. <https://doi.org/10.1186/1756-0500-7-77>.
- Stamatakis, Alexandros. 2014. "RAxML Version 8: A Tool for Phylogenetic Analysis and Post-Analysis of Large Phylogenies." *Bioinformatics* 30 (9): 1312–13. <https://doi.org/10.1093/bioinformatics/btu033>.
- Tillich, Michael, Pascal Lehwark, Tommaso Pellizzer, Elena S. Ulbricht-Jones, Axel Fischer, Ralph Bock, and Stephan Greiner. 2017. "GeSeq - Versatile and Accurate Annotation of Organelle Genomes." *Nucleic Acids Research* 45 (W1). <https://doi.org/10.1093/nar/gkx391>.
- Tonti-Filippini, Julian, Paul G. Nevill, Kingsley Dixon, and Ian Small. 2017. "What Can We Do with 1000 Plastid Genomes?" *Plant Journal* 90 (4). <https://doi.org/10.1111/tpj.13491>.
- Twyford, Alex D., and Rob W. Ness. 2017. "Strategies for Complete Plastid Genome Sequencing." *Molecular Ecology Resources* 17 (5). <https://doi.org/10.1111/1755-0998.12626>.
- Vries, J de, F L Sousa, B Bolter, J Soll, and S B Gould. 2015. "YCF1: A Green TIC?" *Plant Cell* 27 (7): 1827–33. <https://doi.org/10.1105/tpc.114.135541>.
- Wagner, Warren L, Peter C Hoch, and Peter H Raven. 2007. "Revised Classification of the Onagraceae." *Systematic Botany Monographs*. 83: 1–240.
- Wang, Weiwen, and Robert Lanfear. 2019. "Long-Reads Reveal That the Chloroplast Genome Exists in Two Distinct Versions in Most Plants." *Genome Biology and Evolution* 11 (12): 3372–81.
- Wang, Weiwen, Robert Lanfear, and Brandon Gaut. 2019. "Long-Reads Reveal That the Chloroplast Genome Exists in Two Distinct Versions in Most Plants." *Genome Biology and Evolution* 11 (12). <https://doi.org/10.1093/gbe/evz256>.
- Wang, Weiwen, Miriam Schalamun, Alejandro Morales-Suarez, David Kainer, Benjamin Schwessinger, and Robert Lanfear. 2018. "Assembly of Chloroplast Genomes with Long- and Short-Read Data: A Comparison of Approaches Using *Eucalyptus pauciflora* as a Test Case." *BMC Genomics* 19 (1). <https://doi.org/10.1186/s12864-018-5348-8>.
- Wanichthanarak, Kwanjeera, Intawat Nookaew, Phongthana Pasookhush, Thidathip Wongsurawat, Piroon Jenjaroenpun, Namkhang Leeratsuwan, Songsak Wattanachaisaereekul, Wonnop Visessanguan, Yongyut Sirivatanauksorn, and Narong Nuntasaen. 2023. "Revisiting Chloroplast Genomic Landscape and Annotation towards Comparative Chloroplast Genomes of Rhamnaceae." *BMC Plant Biology* 23 (1): 59.

- Wheeler, G L, H E Dorman, A Buchanan, L Challagundla, and L E Wallace. 2014. "A Review of the Prevalence, Utility, and Caveats of Using Chloroplast Simple Sequence Repeats for Studies of Plant Biology." *Appl Plant Sci* 2 (12). <https://doi.org/10.3732/apps.1400059>.
- Wick, Ryan R., Mark B. Schultz, Justin Zobel, and Kathryn E. Holt. 2015. "Bandage: Interactive Visualization of de Novo Genome Assemblies." *Bioinformatics* 31 (20). <https://doi.org/10.1093/bioinformatics/btv383>.
- Wu, S, J Chen, Y Li, A Liu, A Li, M Yin, N Shrestha, J Liu, and G Ren. 2021. "Extensive Genomic Rearrangements Mediated by Repetitive Sequences in Plastomes of *Medicago* and Its Relatives." *BMC Plant Biol* 21 (1): 421. <https://doi.org/10.1186/s12870-021-03202-3>.
- Wu, Z, R Liao, T Yang, X Dong, D Lan, R Qin, and H Liu. 2020. "Analysis of Six Chloroplast Genomes Provides Insight into the Evolution of *Chrysosplenium* (Saxifragaceae)." *BMC Genomics* 21 (1): 621. <https://doi.org/10.1186/s12864-020-07045-4>.
- Xie, Z, and S Merchant. 1996. "The Plastid-Encoded *CcsA* Gene Is Required for Heme Attachment to Chloroplast c-Type Cytochromes." *J Biol Chem* 271 (9): 4632–39. <https://doi.org/10.1074/jbc.271.9.4632>.
- Xing, J, J Pan, H Yi, K Lv, Q Gan, M Wang, H Ge, et al. 2022. "The Plastid-Encoded Protein Orf2971 Is Required for Protein Translocation and Chloroplast Quality Control." *Plant Cell* 34 (9): 3383–99. <https://doi.org/10.1093/plcell/koac180>.
- Xu, J, X Shen, B Liao, J Xu, and D Hou. 2020. "Comparing and Phylogenetic Analysis Chloroplast Genome of Three *Achyranthes* Species." *Sci Rep* 10 (1): 10818. <https://doi.org/10.1038/s41598-020-67679-y>.
- Xu, S, K Teng, H Zhang, K Gao, J Wu, L Duan, Y Yue, and X Fan. 2023. "Chloroplast Genomes of Four *Carex* Species: Long Repetitive Sequences Trigger Dramatic Changes in Chloroplast Genome Structure." *Front Plant Sci* 14:1100876. <https://doi.org/10.3389/fpls.2023.1100876>.
- Zardini, Elsa, and Peter H. Raven. 1992. "A New Section of *Ludwigia* (Onagraceae) with a Key to the Sections of the Genus." *Systematic Botany* 17 (3). <https://doi.org/10.2307/2419486>.
- Zeb, U, X Wang, A AzizUllah, S Fiaz, H Khan, S Ullah, H Ali, and K Shahzad. 2022. "Comparative Genome Sequence and Phylogenetic Analysis of Chloroplast for Evolutionary Relationship among *Pinus* Species." *Saudi J Biol Sci* 29 (3): 1618–27. <https://doi.org/10.1016/j.sjbs.2021.10.070>.
- Zerbino, Daniel R., and Ewan Birney. 2008. "Velvet: Algorithms for de Novo Short Read Assembly Using de Bruijn Graphs." *Genome Research* 18 (5): 821–29. <https://doi.org/10.1101/gr.074492.107>.
- Zhang, Quan, and Sodmergen. 2010. "Why Does Biparental Plastid Inheritance Revive in Angiosperms?" *Journal of Plant Research* 123 (2). <https://doi.org/10.1007/s10265-009-0291-z>.
- Zhang, X F, J B Landis, H X Wang, Z X Zhu, and H F Wang. 2021. "Comparative Analysis of Chloroplast Genome Structure and Molecular Dating in Myrtales." *BMC Plant Biol* 21 (1): 219. <https://doi.org/10.1186/s12870-021-02985-9>.
- Zhang, Y, L Du, A Liu, J Chen, L Wu, W Hu, W Zhang, et al. 2016. "The Complete Chloroplast Genome Sequences of Five *Epimedium* Species: Lights into Phylogenetic and Taxonomic Analyses." *Front Plant Sci* 7:306. <https://doi.org/10.3389/fpls.2016.00306>.
- Zhong, Q, S Yang, X Sun, L Wang, and Y Li. 2019. "The Complete Chloroplast Genome of the Jerusalem Artichoke (*Helianthus Tuberosus* L.) and an Adaptive Evolutionary Analysis of the *Ycf2* Gene." *PeerJ* 7:e7596. <https://doi.org/10.7717/peerj.7596>.
- Zhong, Xiao. 2020. "Assembly, Annotation and Analysis of Chloroplast Genomes", Doctor of Philosophy, The University of Western Australia. <https://doi.org/10.26182/5f333d9ac2bee>

Appendix

	ABySS		MEGAHIT		VELVET		SPAdes	
	not corrected	corrected	not corrected	corrected	not corrected	corrected	not corrected	corrected
A								
Using all contigs								
Genome fraction (%)	86.868	85.279	86.428	85.158	91.927	86.796	84.682	84.483
Duplication ratio	1.047	1.042	1.796	1.041	2.002	1.128	1.042	1
Largest alignment	56 588	41 262	30 904	90 352	3531	17 235	90 399	90 272
Using contigs > 200 nt								
Genome fraction (%)	86.419	85.279	86.377	85.057	76.589	86.181	84.682	84.483
Duplication ratio	1.028	1.042	1.681	1.029	1.177	1.11	1.042	1
Largest alignment	56 588	41 262	30 904	90 352	3531	17 235	90 399	90 272
Using contigs > 500 nt								
Genome fraction (%)	85.564	84.517	83.503	84.774	45.468	79.279	84.682	84.483
Duplication ratio	1.009	1.012	1.041	1.004	1.015	1.054	1.042	1
Largest alignment	56 588	41 262	30 904	90 352	3531	17 235	90 399	90 272
Using contigs > 1000 nt								
Genome fraction (%)	83.701	84.199	81.256	84.545	22.194	66.438	84.563	84.483
Duplication ratio	1	1.002	1.007	1.001	1	1.011	1.026	1
Largest alignment	56 588	41 262	30 904	90 352	3531	17 235	90 399	90 272
B								
Using all contigs								
NGA50	15 215	26 577	19 986	90 352	469	2796	90 399	90 272
LGA50	3	3	3	1	93	9	1	1
Misassemblies								
# misassemblies	0	0	4	0	0	0	0	0
Misassembled contigs length	0	0	1595	0	0	0	0	0
Mismatches								
# mismatches per 100 kbp	109.53	107.19	1036.93	45.24	499.16	229.11	96.57	0
# indels per 100 kbp	12.4	10.58	62.99	16.26	27.92	74.88	19.17	0
# N's per 100 kbp	0	0	0	0	0	6.1	0	0
Using contigs > 500 nt								
NGA50	15 215	26 577	19 986	90 352	-	2796	90 399	90 272
LGA50	3	3	3	1	-	9	1	1
Misassemblies								
# misassemblies	0	0	1	0	0	0	0	0
Misassembled contigs length	0	0	665	0	0	0	0	0
Mismatches								
# mismatches per 100 kbp	62.39	46.17	123.32	8.1	221.33	148.48	96.57	0
# indels per 100 kbp	2.9	2.2	4.33	1.47	28.51	63.74	19.17	0
# N's per 100 kbp	0	0	0	0	0	7.22	0	0
Using contigs > 1000 nt								
NGA50	15 215	26 577	19 986	90 352	-	2796	90 399	90 272
LGA50	3	3	3	1	-	9	1	1
Misassemblies								
# misassemblies	0	0	0	0	0	0	0	0
Misassembled contigs length	0	0	0	0	0	0	0	0
Mismatches								
# mismatches per 100 kbp	0	0	37.51	0.74	64.94	61.6	67.87	0
# indels per 100 kbp	1.5	0.74	0	0.74	25.41	56.93	19.49	0
# N's per 100 kbp	0	0	0	0	0	9.25	0	0

Figure A1 - QAST evaluation of performance of the four assembly tools (using corrected or uncorrected SRs). A: Comparison of plastome fraction, duplication rate and size of the largest alignment obtained. B: Comparison of classic metrics (NGA50 and LGA50), number of errors (misassemblies and mismatches) produced.

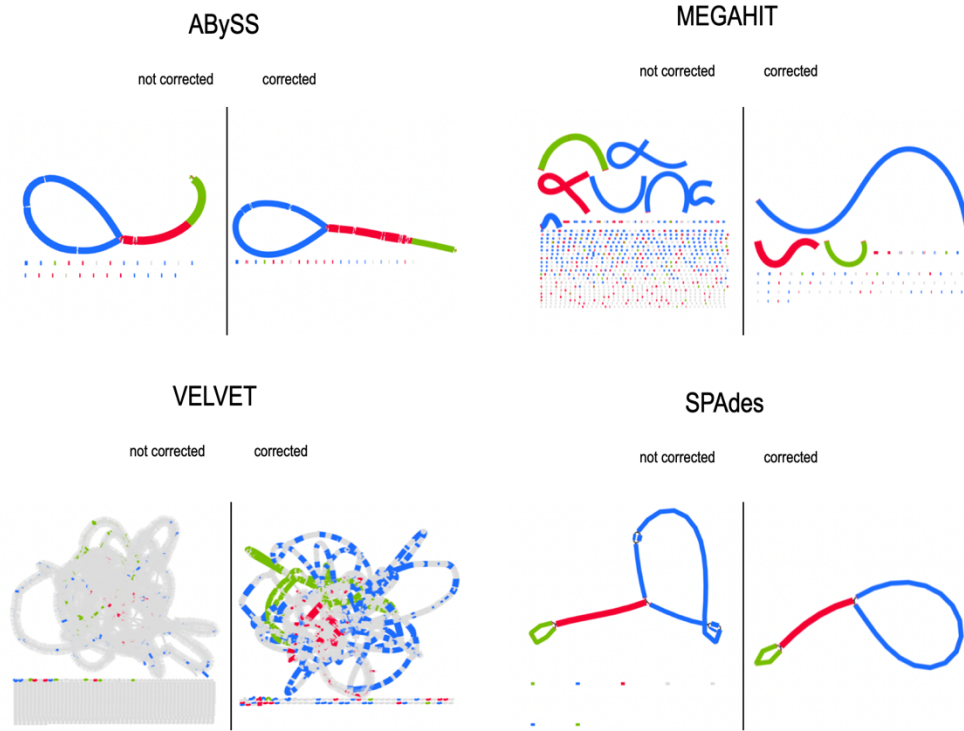


Figure A2 - BANDAGE visualization of the *L. grandiflora* plastome assembly graphs on corrected or uncorrected SRs. Contigs are colored according to their BLAST match to the LSC (blue), SSC (green), and IR (red) segments.

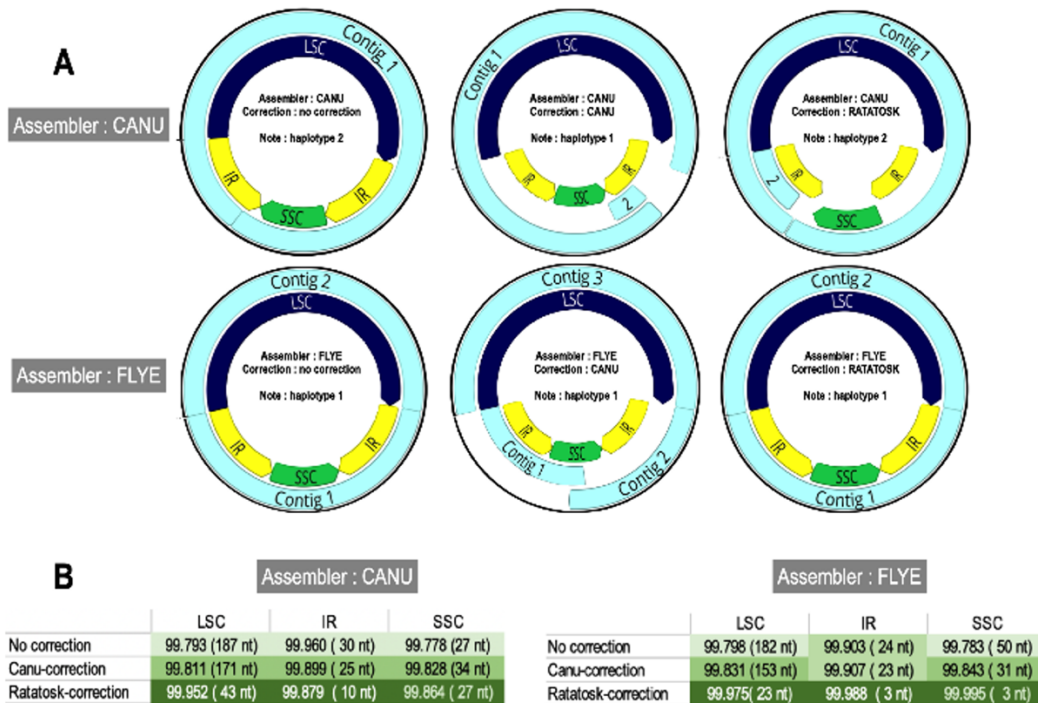


Figure A3 - Graphs representing the assemblies of *L. grandiflora* long reads. **A:** Contigs are represented in light blue and the three segments (LSC, SSC and IR) in dark blue, green and yellow, respectively. **B:** Comparative effectiveness of CANU and RATATOSK correctors.

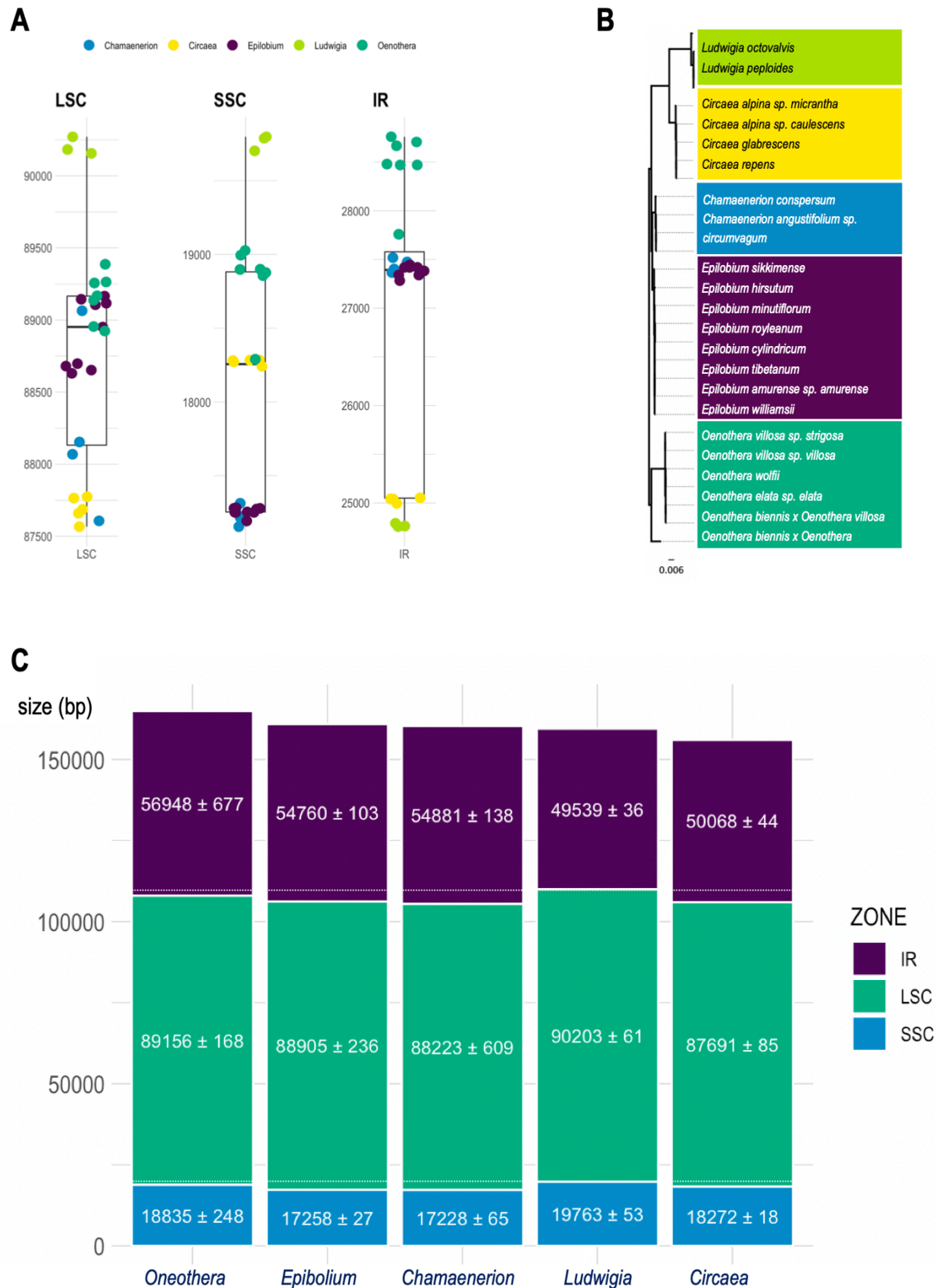


Figure A4 - Comparison of LSC, SSC and IR sizes in the Onagraceae. **A:** Comparison of the sizes of LSC, SSC and IR segments in the Onograceae family (*Chamaenerion* in blue, *Circaea* in yellow, *Epilobium* in dark purple, Ludwigia in light green and *Oenothera* in dark green). **B:** Maximum likelihood tree made using RAxML (model GTR-GAMMA, algorithm Rapid Hill-climbing) on multiple sequences alignment of Onograceae plastomes made using MAFFT. **C:** Average size of the different chloroplast segments (LSC, SSC and IR) for the 5 genres of Onograceae. IR size corresponds to the sum of the two copies.

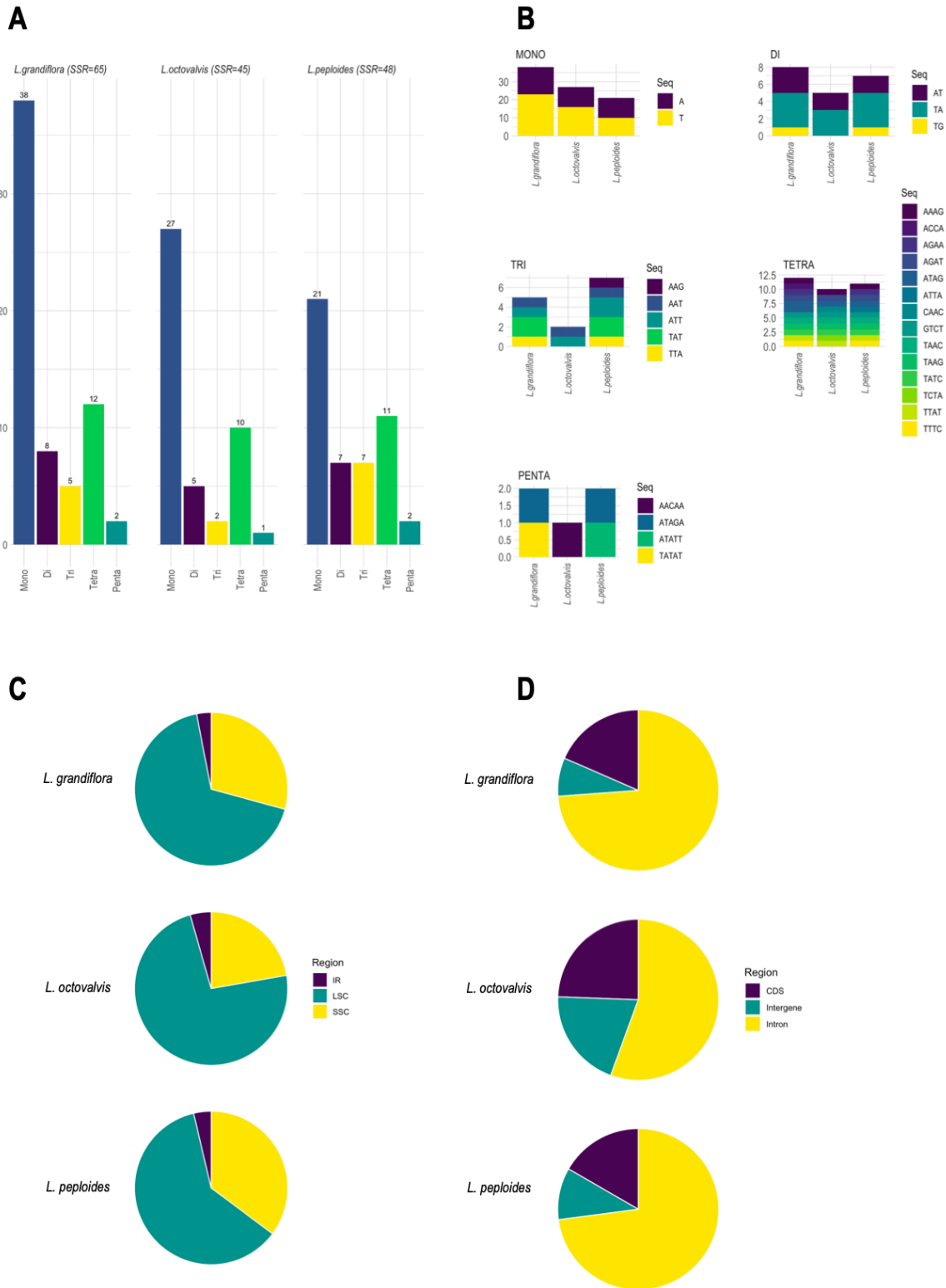


Figure A5 - Comparative analysis of Simple-Sequence Repeats (SSRs) in Ludwigia chloroplast genomes. **A:** SSR numbers detected in the three species, by repeat class types (mono, di-, tri-, tetra and pentanucleotides). **B:** Frequency of SSR motifs by repeat class types. **C:** Frequency of SSRs in LSC, SSC and IR regions. **D:** Repartition of SSRs in intergenic, protein-coding and intronic regions.

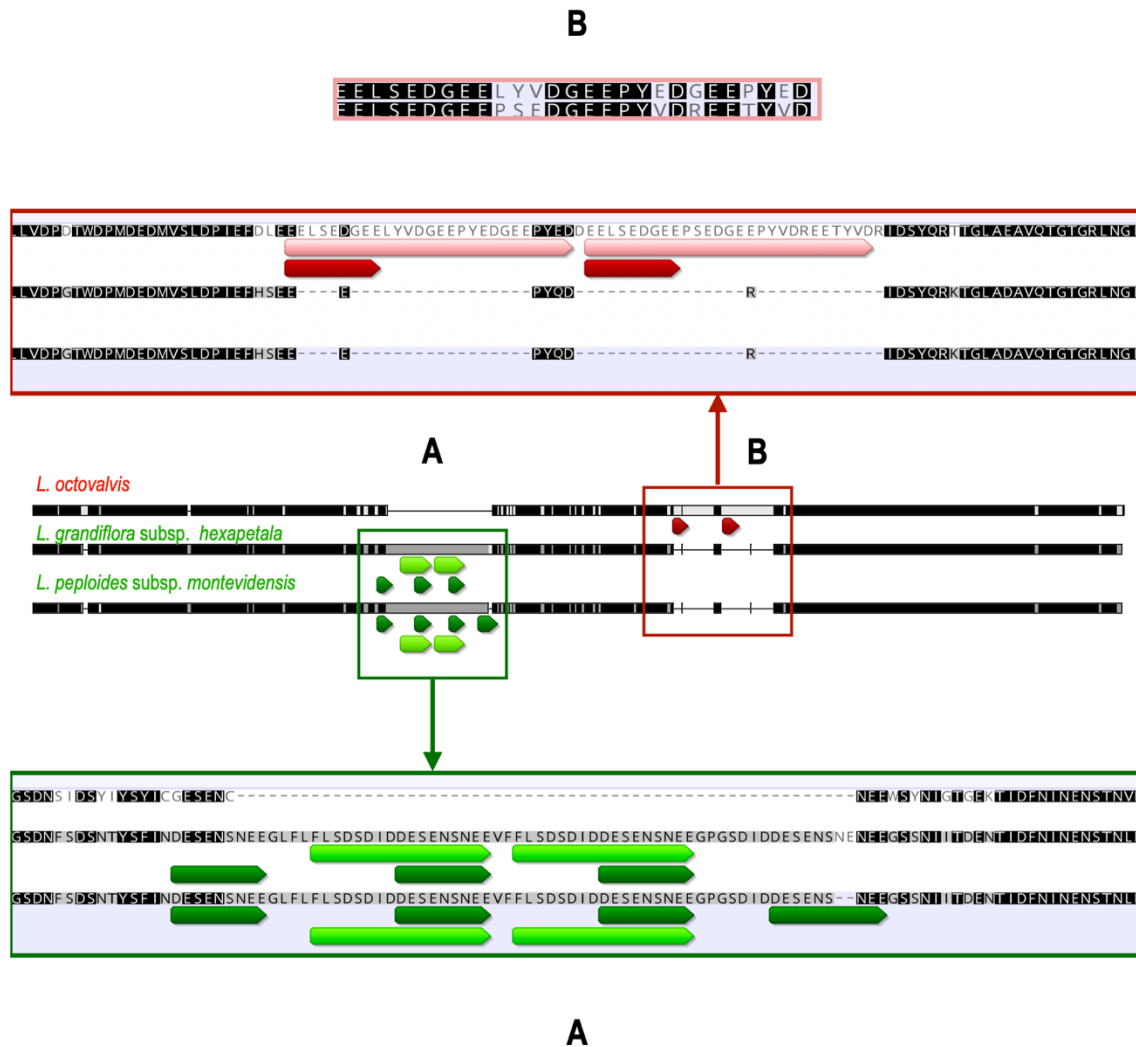


Figure A6 - Diagram showing the position of tandem repeats in the *accD* gene. *L. octovalvis* (in red) and *L. peplodes* and *L. grandiflora* (in green). We also observe the consequences of these repetitions on the insertion of amino acids, also repeated.

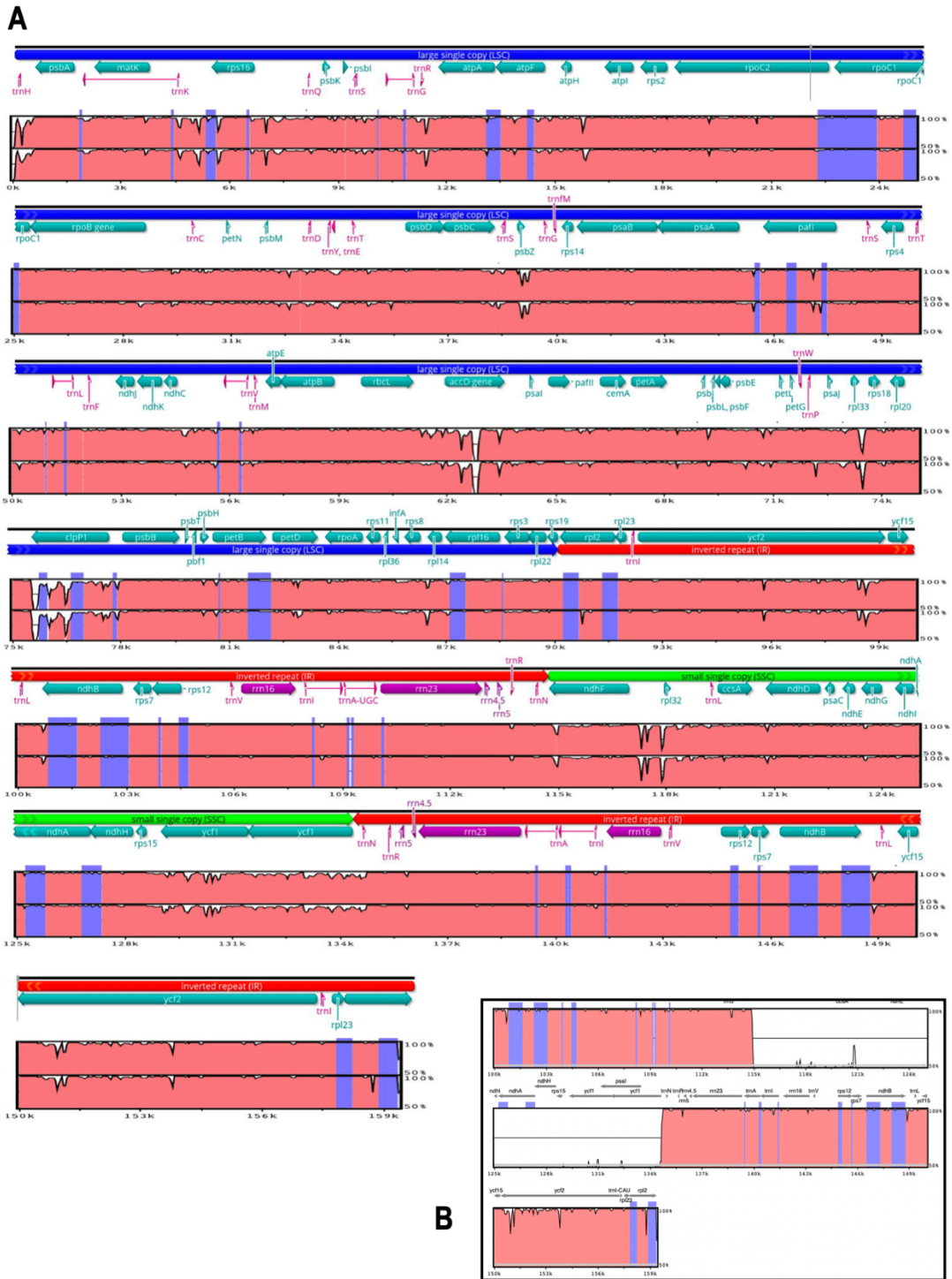


Figure A7 - Comparison of the three Ludwigia plastomes using mVISTA, with the *L. octovalvis* as a reference. **A:** The y-axis represents the identity percentage (between 50 and 100%). The arrows show the genes (in green: proteins genes, in purple: rRNAs and in fuchsia: tRNAs). Blue blocks indicate exonic regions. LCS, IR and SSC regions are also distinguished (in dark blue, red and green, respectively). The second line corresponds to *L. grandiflora* haplotype 2 (For this haplotype, SSC segment is oriented like *L. octovalvis*) and the third line corresponds to *L. peploides* for which the SSC region has been artificially oriented in the same way as the two other plastomes to allow comparison. **B:** Small box showing a part of the alignment and presenting the consequences if we do not artificially orient the SSC segments in the same direction for the analysis.

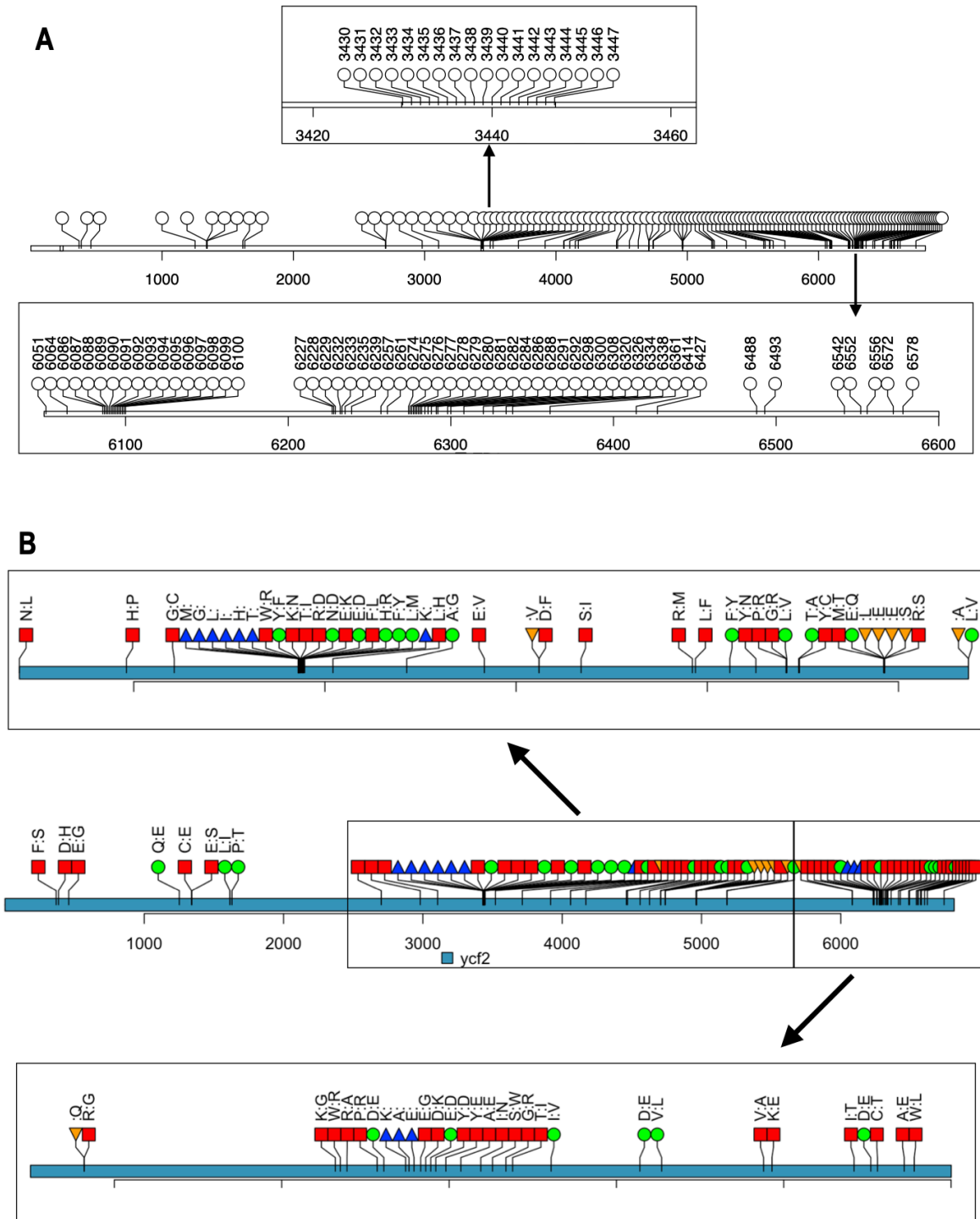


Figure A8 - Lollipop diagram allowing the visualization of SNPs and their translational effects on the *ycf2*. **A**: localization of the 256 single nucleotide polymorphisms (SNP) observed by comparing *L. grandiflora*-*L. peploides* with *L. octovalvis*. Two regions particularly dense in SNPs (between 3420 and 3460 and between 6100 and 6600) have been zoomed into to allow better reading. **B**: Effect of these SNPs on the translated sequence of *L. octovalvis*, compared to Ycf2 of the other two species: non conservative mutation: red square; conservative mutation: circle green; deletion: triangle_point_up blue and insertion: triangle_point_down, orange. As for A, two regions were zoomed into in order to distinguish each mutation.