

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

Frédérique Barloy-Hubler, Anne-Laure Le Gac, Christophe Boury, Erwan

Guichoux, Dominique Hermine Barloy

▶ To cite this version:

Frédérique Barloy-Hubler, Anne-Laure Le Gac, Christophe Boury, Erwan Guichoux, Dominique Hermine Barloy. Sequencing, de novo assembly of Ludwigia plastomes, and comparative analysis within the Onagraceae family. 2023. hal-04255511

HAL Id: hal-04255511 https://cnrs.hal.science/hal-04255511

Preprint submitted on 24 Oct 2023

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.

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

Barloy-Hubler F.³, Le Gac A.-L.¹, Boury C.², Guichoux E.², and Barloy D.^{1*}

1-DECOD (Ecosystem Dynamics and Sustainability), Institut Agro, IFREMER, INRAE, Rennes, France.

2- Université de Bordeaux, INRAE, BIOGECO, Cestas, France

3- Université de Rennes 1, CNRS, UMR 6553 ECOBIO, Rennes, France.

* Corresponding author: dominique.barloy@agrocampus-ouest.fr

Abstract

Background

The Onagraceae family, which belongs to the order Myrtales, consists of approximately 657 species and 17 genera, including the genus *Ludwigia* L., which is comprised of 82 species. There are few genomic resources for Onagraceae, which limits phylogenetic and population genetics, as well as genomic studies. In this study, new complete plastid genomes of *Ludwigia grandiflora subps. hexapetala* (*Lgh*) and *Ludwigia peploides subsp montevidensis* (*Lpm*) were generated using a combination of different sequencing technologies. These plastomes were then compared to the published *Ludwigia octovalvis* (*Lo*) plastid genome, which was re-annotated by the authors.

Results

We initially sequenced and assembled the chloroplast (cp) genomes of Lpm and Lgh using a hybrid strategy. 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, 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.

Conclusion

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.

Keywords

Water primrose, Ludwigia, Onagraceae, chloroplast genome, long and short reads, hybrid assembly, haplotype

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 [1]. *Ludwigia* L. is composed of 83 species [2, 3]. 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. [4].

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 (1987) subsp. *hexapetala* (Hook. & Arn.) (Nesom & Kartesz 2000) (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 described as aquatic invasive plants in Europe [5]. In France, both species occupied aquatic habitats, such as static or slow-flowing waters, riversides, and have recently been observed in wet meadows [6]. The transition from an aquatic to a terrestrial habitat has led to the emergence of two *Lgh* morphotypes [7]. 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 [8].

Adequate genomic resources are necessary in order to be able to analyze the different (epi)genomic mechanisms of this acclimation. However, even though the number of terrestrial plant genomes has increased considerably over the last 20 years, only a small fraction (~ 0.16%) have been sequenced, with some clades being significantly more represented than others [9]. 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. As of April 2023, there are 10,712 reference plastomes (RefSeq) listed on GenBank, with the vast majority (10,392 genomes) belonging to Viridiplantae (green plants). However, 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 Onagraee, 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 of the genus Ludwigia while *Lpm* and *Lgh* belong to jussieae section [10],[11]. 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 [12, 13]. Both maternal and biparental inheritance are described in the Onagraceae family. In tribe Onagreae, *Oenothera* subsect. *Oenothera* are known to have biparental plastid inheritance [1, 14]. In tribe Epilobieae, biparental plastid inheritance was also reported in *Epilobium* L. with mainly maternal transmission, and very low proportions of paternally transmitted chloroplasts [15].

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) [16]. In general, chloroplast genomes (plastomes) are highly conserved in size, structure, and genetic content. They are rather small (120-170 kb, [17]), 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) [18]. Plastomes are generally circular but linear shapes also exist [19]. 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 [20].

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, [20]), which are longer than the short read lengths (< 300 bp). This leads to incomplete or approximate assemblies [21]. Recent long-read sequencing (> 1000 bp) provides strong evidence that terrestrial plant plastomes have two structural haplotypes, present in equal proportions and differing in IR orientation of the [22]. 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 [23-25]. Thus, a hybrid strategy which combines long reads (to access the genomic structure) and short reads (to correct sequencing errors) could be effective [21, 26].

Here, we report the newly sequenced complete plastid genomes of *Ludwigia grandiflora* subps. *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. mitochlondrial 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 [4, 27, 28].

Material and Methods

Plant sampling and experimental design

The original plant materials were collected in June of 2018 near to Nantes (France) and formal 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 ¹/₃ soil, ¹/₃ sand, 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, genomic DNA extraction was carried out using a modified version of the protocol proposed by Panova et al in 2016, with three purification steps [29].

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 [30]. 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 μ l 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). 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 μ g 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 performed using Guppy v4.0.14 base call software for ONT and fastp v0.20.0 [31] using the default quality filter (phred quality \geq = Q15) for Illumina reads. After filtering and trimming, Lgh Illumina-reads (SR, 2*23,067,490 reads) were used for reference-guided assembly using GetOrganelle v1.7.0 [31] and NOVOPlasty v4.2.1 [32] using the parameters recommended for green plants. As we sequenced total DNA, chloroplast reads are mixed with nuclear and mitochondrial reads. Chloroplastic reads were then extracted by mapping against draft Lgh plastomes generated by GetOrganelle in order to reduce computational costs. Chloroplastic R1 and R2 reads were used with and without prior error correction using ONT reads with BayesHammer [33] for *de novo* assemblies using ABySS (version 2.1.5 [34, 35]), MEGAHIT (1.1.2, [36]), Velvet (version 1.2.10) using VelvetOptimiser and SPAdes (version 3.15.4, [37]). The best k-mer parameters were tested using kmergenie [38] and Velvetoptimizer. The best results were obtained with k=99, and was therefore the parameter used for all MiSeq assemblers. For ONT reads, Lgh (550,516 reads) and Lpm (68,907 reads) LR reads were self-corrected using CANU 1.8 [39] and de novo assembly using CANU and FLYE 2.8.2 run with the option --meta and -plasmids [40]. Hybrid correction and assemblies of LR reads using SR data were realized for *Lgh* using Ratatosk [41].

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 (<u>http://cab.cc.spbu.ru/quast/</u>) and visualized using BANDAGE [42]. BANDAGE

compatible graphs (.gfa format) were created with the megahit_toolkit for MEGAHIT and with gfatools for ABySS. 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 [43] using ARAGORN and tRNAscan_SE to predict tRNAs and rRNAs and tRNAscan_SE to predict tRNAs and rRNAs and via Chloe prediction site (https://chloe.plantenergy.edu.au). 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 biocurated using Geneious (v.10). Finally, plastomes were represented using OrganellarGenomeDRAW (OGDRAW) [44]. 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 [45].

SSRs and Repeat Sequences Analysis

Simple Sequence Repeats (SSRs) were analyzed through the MISA web server [46], 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 <u>RepEx</u>. 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 [47], 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 [48], with the LAGAN alignment algorithm [49] and a cut-off of 70% identity.

Nucleotide diversity (Pi) was analyzed using the software DnaSP v.6.12.01 [50, 51] with step size set to 200 bp and window length to 300 bp. IRscope [52] 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 (Ka) and synonymous (Ks) substitutions and their ratios ($\omega = Ks/Ks$) using TBtools [53] 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.

Plastome phylogenomic analyses

We reconstructed phylogenetic relationships among plastomes of Onagraceae. The FFT-NS-2 method in MAFFT 7 [54] was used to align all plastomes with one of the IRs removed to avoid data duplication. Phylogenetic tree analysis, based on the maximum likelihood (ML) method, was conducted in RAxML 8.2.9 [55] using the GTR+G model with node support assessed by fast-bootstrap (-f a) using 1,000 non-parametric bootstrap pseudo-replicates.

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, <u>http://www.rstudio.com/</u>). Figures were realized using ggplot2, ggpubr, tidyverse, dplyr, gridExtra, reshape2, and viridis packages. SNPs were represented using trackViewer [52] and genes represented using gggenes packages.

Results

Plastome short read assembly

We first used GetOrganelle and Novoplasty with the short reads of *Lgh*. GetOrganelle reconstructed two identical haplotypes that differ only by a "flip-flop" of the SSC region, as shown in Figure 1, while NOVOplasty generated three contigs that can be reassembled to form either one of the two haplotypes, as the SSC region can be oriented in both directions.

In order to evaluate the best strategies for assembly of *Lgh* plastomes, we used these haplotypes to map and extract short (Illumina) and long (Nanopore) chloroplast reads. Chloroplastic short reads, with or without prior correction using long reads, were assembled using ABySS, Velvet, MEGAHIT and SPAdes. As shown in Figure 2A, the number of contigs obtained is highly variable depending on the tools and strategies. Without correction, Velvet and MEGAHIT generate highly fragmented assemblies with 1103 and 722 contigs respectively, while ABySS and SPAdes produce a more reasonable number of contigs (65 and 12

respectively). In addition, all contigs produced by all assemblers except SPAdes, are predominantly less than 500 nt long (Figure 2B from 68 to 95% of the contigs). By correcting the reads before assembly, we saw a noticeable improvement in the number of contigs, which decreased for all assemblers while increasing their sizes. However, except for SPAdes, the number of contigs below 500 nt remained preponderant (from 54 to 70%, Figure 2C).

An assessment of the assembly quality using QUAST, and based on an alignment on *Lgh* haplotype 1 generated by GetOrganelle as a reference, shows that none of the assemblers produces a complete plastome. By comparing the contribution of small contigs (less than 200, 500 and 1000 nt), we noticed an invariant coverage rate for SPAdes (85%), while the fraction of covered genome decreases drastically for Velvet (from 92% to 22%), correlated with a duplication rate going from 2 to 1 (Add. Figure 1A). We noted that compared to an assembly with all reads, the use of reads superior to 1 kb slightly decreases the coverage rate to reach an optimal duplication ratio equal to 1. The effect of prior read correction is notable for MEGAHIT and Velvet, especially concerning the increase in the size of the large alignment (Add. Figure 1A), loss of misassemblies, and reduction of the number of mismatches (Add. Figure 1B).

Investigating all possible scaffolds generated by BANDAGE visualization (Add. Figure 2) shows that ABySS and SPAdes assembly graphs suggest a 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). In all graphs, the two IRs are a collapsed repeat 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, SPAdes is the most efficient tool, and gives the best result using corrected short reads to assemble a resolved plastome, with three contigs of 90,272 bp (corresponding to LSC), 19,788 bp (corresponding to SSC) and 24,762 bp (corresponding to a copy of the IR).

Plastome long read assembly

Chloroplast fractions of *Lgh* long reads were selected through mapping. Raw, selfcorrected reads were assembled using CANU, and short-read-corrected comparative assembly was carried out using CANU or FLYE. Using raw data, CANU generates a unique contig corresponding to haplotype 2, whereas FLYE makes two contigs that reconstruct haplotype 1. Assemblies using long reads corrected by CANU leads to fragmentation into two (CANU assembler) or three (FLYE assembler) contigs which both reconstruct haplotype 1, with an large gap corresponding to one of the IR copies for the CANU assembler. Finally, correction by RATATOSK allows CANU to assemble two redundant contigs, with contig 1 reproducing the entire haplotype 2, while FLYE makes two contigs corresponding to haplotype 1 (Add. Figure 3A). In conclusion, the two haplotypes proposed by GetOrganelle using the short reads were likewise reconstructed randomly using the long reads, with or without correction. Consequently, both haplotypes were retained as *Lgh* plastomes and submitted to Genbank. Hybrid assembly is the most complete and the most faithful, with an average of 1.6 contigs (all assemblers combined), and optimal correction using RATATOSK, which permits plastomes with 99.94% accuracy compared to biocurated molecules (Add. Figure 3B).

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

1. General Variations

Plastomes of the three species of Ludwigia sp., *Lgh*, *Lpm* and *Lo*, are circular doublestranded DNA molecules (Figure 3) which are all (as shown in Table 1) approximately the same size: *Lo* being slightly smallest and *Lgh* the biggest. The overall GC content is almost the same for the three species and the GC contents of the IR regions are higher than those of the LSC and SSC regions. Between the three species, the lengths of the total chloroplasts, LSC, SSC, and IR are broadly similar (Table 2) 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 therby 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.

2. 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 *Lg* and *Lo* which both carry a complete *ycf1* gene.

If we compare Ludwigia sp. junctions with those of other Onagraceae plastomes (Figure 5), we can observe that the JLB and JLA connections are well-preserved in the whole family, whereas JSB and JLA differ. Concerning JSB, in the five Onagraceae family genera with available plastomes, *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*, *Epibolium 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, 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 (Add. Figure 4).

3. 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 (Add. Figure 5A). 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 (Add. Figure 5B). 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 (Add. Figure 5C). The analysis of SRR locations revealed that most were distributed in non-coding regions (intergenic regions and introns, Add. Figure 5D).

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 3A): 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 3A: 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 3B): 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 3C): 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 palidromes 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.

4. 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 3 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.

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 9 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 9 amino acid repeats [EELSEDGEE], included in two longer degenerate repeats of 27 nt (Add. Figure 6). 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 [56]. Large accD expansions due to TRs have also been described in other plants such as Medicago [57] and *Cupressophytes* [58]. 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 [59].

5. 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 (Add. Figure 7). 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 [60]. 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 [61, 62]. Moreover, there is a deletion of nine nucleotides in the 3' region of the *rpl32* gene in *Lgh* and *Lpm*, leading to a premature end of the translation and the deletion of the last 4 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 *rpl32* function since they are not found in the orthologs.

Our results show that the Ka/Ks 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* (Add. Figure 8), compared to *Lgh* and *Lpm* (100 % identical). This gene has been shown as "variant" in other plant species such as *Helianthus tuberosus* [63].

The *matK* gene has been used as a universal barcoding locus to enable species discrimination of terrestrial plants [64], 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, 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 [65]), *L. stolonifera* (native to the Nile, found in a variety of habitats, from freshwater wetlands to brackish and marine waters) [66] and *L.*

adscendens (common weed of rice fields in Asia) [67]. *Lpm* is in a branch unique to its species but close to the *L. grandiflora* subsp. *hexapetala* branch.

Discussion

In the present study, we first sequenced and *de novo* assembled the chloroplast (cp) genomes of *Ludwigia peploides* (*Lpm*) and *Ludwigia grandiflora* (*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 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 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. Our findings corroborate with similar results obtained in studies on other chloroplasts, such as those in *Eucalyptus* [21], *Falcataria* [68], *Carex* [69] or *Cypripedium* [70]. As a result of this strategy, we were able to successfully identify the presence of two haplotypes in Lgh, which is a first for Ludwigia, as the plastome of L. octovalvis was only delivered in the form of one haplotype [71]. Due to the unavailability of sequence data for *Ludwigia octovalvis* and our exclusive use of long reads for *Ludwigia peploides*, 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. Notably, the existence of two plastome haplotypes has been widely identified in the closely related order of Myrtales (Eucalyptus sp.), and more commonly in Angiosperms, with 58 species across 16 other orders showing similar patterns [22]. Recent studies have also revealed the presence of two cp haplotypes in other orders, such as Asparagales (Ophrys apifera orchid, [72]), Brassicales (Carica papaya,

Vasconcellea pubescens, [73]), Solanales (*Solanum tuberosum*, [74]), and even in a new order, Laurales (Avocado *Persea americana*, [75]). Furthermore, Wanichthanarak *et al.* (2023) revisited 24 Ramanaceae published chloroplast genomes and identified a likely form of heteroplasmy in *Rhamnus crenata* [76]. These emerging findings highlight the potential underestimation of the occurrence of two chloroplast haplotypes in angiosperms and emphasize the importance of identifying and considering such heteroplasmy during *de novo* or revisited chloroplast genome assemblies.

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 [77]. 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, [78]). 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%, [77]) and more generally for the order Myrtales (36.9–38.9%, with the average GC content being 37%, [78]). 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, [79]) or Lamiaceae (order Lamiales, [80]), 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 [78]. Chloroplast genes have been selected during evolution due to their functional importance [81]. 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. in 2007 [1]. 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 [7]. Additionally, *Lo* is a

wetland plant that typically grows in gullies and at the edges of ponds, as documented by Wagner et al. in 2007 [1]. 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*, Barthet et al [82] 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., [83]), and more generally in terrestrial plants (Pinus sp, [84] or *Chrysosplenium sp.* [85]). The *accD* gene has been described as an essential gene required for leaf development [86] and longevity [87] in tobacco (Nicotiana tabacum). Under drought stress, plant resistance can be increased by inhibiting accD [88], and conversely, enhanced in response to flooding stress by upregulating accD accumulation [89]. 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 vcf2 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 [90, 91]), as well as plant cell survival [92, 93]. The *ccsA* gene positive selection is found in some aquatic plants such as Anubia sp. [83], marine flowering plants as Zostera species [94], and some species of Lythraceae [88]. The ccsA gene is required for cytochrome c biogenesis [95] and this hemoprotein plays a key role in aerobic and anaerobic respiration, as well as photosynthesis $[\underline{96}]$. 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 [8] 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.* [4]. 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 [97, 98]. 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 [99, 100]. 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 [101-103]. 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.

Conclusions

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 both *Lpm* and *Lgh*, while the absence of the *Lo* genome precluded further investigation for this species. 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 [104].

Acknowledgements

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 (doi:10.15454/1.5572396583599417E12).

Declarations

An official authorization from "prefecture de Loire Atlantique" was obtained to collect and transport *L. grandiflora* subsp. *hexapetala* (*Lgh*) from the Mazerolles swamps (commune Petit Mars, France), and *L. peploides* subsp. montevidensis (*Lpm*) from La Musse (commune Saint Etienne de Montluc, France) (prefectorial decree n°2018/SEE/2423). In accordance with French legislation on invasive species and as indicated in prefectorial decree n°2018/SEE/2423 (article 3), the samples were duly destroyed and under no circumstances could they be placed in a herbarium, which is strictly forbidden regarding local and national guidelines.

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).

The authors declare that they have no competing interests.

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

Authors' contributions

F.B-H and D.B conceived and designed the study. A-L.L-G carried out DNA extractions. C.B and E.G made sequencing run. A-L.L-G and F.B-H assembled and annotated the genomes. F.B-H and D.B analyzed the data and drafted the original manuscript. F.B-H prepared all illustrations. All authors reviewed and approved the manuscript.

References

- 1. Wagner WL, Hoch PC, Raven PH: **Revised classification of the Onagraceae**. *Systematic Botany Monographs* 2007.
- Levin RA, Wagner WL, Hoch PC, Nepokroeff M, Pires JC, Zimmer EA, Sytsma KJ: Family-level relationships of Onagraceae based on chloroplast rbcL and ndhF data. American Journal of Botany 2003, 90(1):107-115.
- 3. Levin RA, Wagner WL, Hoch PC, Hahn WJ, Rodriguez A, Baum DA, Katinas L, Zimmer EA, Sytsma KJ: Paraphyly in tribe Onagreae: insights into phylogenetic relationships of Onagraceae based on nuclear and chloroplast sequence data. Systematic Botany 2004, 29(1):147-164.
- 4. Liu S-H, Hoch PC, Diazgranados M, Raven PH, Barber JC: Multi-locus phylogeny of Ludwigia (Onagraceae): insights on infra-generic relationships and the current classification of the genus. *Taxon* 2017, 66(5):1112-1127.
- 5. Hussner A, Windhaus M, Starfinger U: From weed biology to successful control: an example of successful management of Ludwigia grandiflora in Germany. Weed Research 2016, 56(6):434-441.
- 6. Lambert E, Dutartre A, Coudreuse J, Haury J: Relationships between the biomass production of invasive Ludwigia species and physical properties of habitats in France. *Hydrobiologia* 2010, 656:173-186.
- Haury J, Druel A, Cabral T, Paulet Y, Bozec M, Coudreuse J: Which adaptations of some invasive Ludwigia spp.(Rosidae, Onagraceae) populations occur in contrasting hydrological conditions in Western France? *Hydrobiologia* 2014, 737:45-56.
- 8. Billet K, Genitoni J, Bozec M, Renault D, Barloy D: Aquatic and terrestrial morphotypes of the aquatic invasive plant, Ludwigia grandiflora, show distinct morphological and metabolomic responses. *Ecol Evol* 2018, **8**(5):2568-2579.
- 9. Marks RA, Hotaling S, Frandsen PB, VanBuren R: **Representation and participation** across 20 years of plant genome sequencing. *Nat Plants* 2021, 7(12):1571-1578.

- 10. Zardini E, Raven PH: A new section of Ludwigia (Onagraceae) with a key to the sections of the genus. *Systematic botany* 1992:481-485.
- 11. Hoch PC, Wagner WL, Raven PH: The correct name for a section of Ludwigia L.(Onagraceae). *PhytoKeys* 2015(50):31.
- 12. Hu Y, Zhang Q, Rao G: Occurrence of plastids in the sperm cells of Caprifoliaceae: biparental plastid inheritance in angiosperms is unilaterally derived from maternal inheritance. *Plant and cell physiology* 2008, **49**(6):958-968.
- 13. Zhang Q: Why does biparental plastid inheritance revive in angiosperms? *Journal* of plant research 2010, **123**(2):201-206.
- 14. Cleland RE: **Oenothera. Cytogenetics and evolution**. 1972.
- Schmitz UK, Kowallik K-V: Plastid inheritance in Epilobium. *Current genetics* 1986, 11(1):1-5.
- 16. Sato N: Are Cyanobacteria an Ancestor of Chloroplasts or Just One of the Gene Donors for Plants and Algae? *Genes (Basel)* 2021, **12**(6).
- 17. Gualberto JM, Mileshina D, Wallet C, Niazi AK, Weber-Lotfi F, Dietrich A: **The plant mitochondrial genome: dynamics and maintenance**. *Biochimie* 2014, **100**:107-120.
- 18. Tonti-Filippini J, Nevill PG, Dixon K, Small I: What can we do with 1000 plastid genomes? In., vol. 90: Wiley Online Library; 2017: 808-818.
- Oldenburg DJ, Bendich AJ: The linear plastid chromosomes of maize: terminal sequences, structures, and implications for DNA replication. *Current genetics* 2016, 62:431-442.
- 20. Twyford AD, Ness RW: Strategies for complete plastid genome sequencing. *Molecular ecology resources* 2017, **17**(5):858-868.
- 21. Wang W, Schalamun M, Morales-Suarez A, Kainer D, Schwessinger B, Lanfear R: Assembly of chloroplast genomes with long-and short-read data: a comparison of approaches using Eucalyptus pauciflora as a test case. *BMC genomics* 2018, **19**:1-15.
- 22. Wang W, Lanfear R: Long-reads reveal that the chloroplast genome exists in two distinct versions in most plants. *Genome Biology and Evolution* 2019, **11**(12):3372-3381.
- Ferrarini M, Moretto M, Ward JA, Šurbanovski N, Stevanović V, Giongo L, Viola R, Cavalieri D, Velasco R, Cestaro A: An evaluation of the PacBio RS platform for sequencing and de novo assembly of a chloroplast genome. *BMC genomics* 2013, 14(1):1-12.
- 24. Jain M, Koren S, Miga KH, Quick J, Rand AC, Sasani TA, Tyson JR, Beggs AD, Dilthey AT, Fiddes IT: Nanopore sequencing and assembly of a human genome with ultra-long reads. *Nature biotechnology* 2018, **36**(4):338-345.
- 25. Rang FJ, Kloosterman WP, de Ridder J: From squiggle to basepair: computational approaches for improving nanopore sequencing read accuracy. *Genome biology* 2018, **19**(1):90.
- 26. Scheunert A, Dorfner M, Lingl T, Oberprieler C: Can we use it? On the utility of de novo and reference-based assembly of Nanopore data for plant plastome sequencing. *PLoS One* 2020, **15**(3):e0226234.
- 27. Bedoya AM, Madriñán S: Evolution of the aquatic habit in Ludwigia (Onagraceae): Morpho-anatomical adaptive strategies in the Neotropics. Aquatic Botany 2015, 120:352-362.
- Liu S-H, Yang H-A, Kono Y, Hoch PC, Barber JC, Peng C-I, Chung K-F: Disentangling Reticulate Evolution of North Temperate Haplostemonous Ludwigia (Onagraceae) 1, 2. Annals of the Missouri Botanical Garden 2020, 105(2):163-182.

- 29. Panova M, Aronsson H, Cameron RA, Dahl P, Godhe A, Lind U, Ortega-Martinez O, Pereyra R, Tesson SV, Wrange A-L: **DNA extraction protocols for whole-genome sequencing in marine organisms**. *Marine genomics: Methods and protocols* 2016:13-44.
- 30. Belser C, Istace B, Denis E, Dubarry M, Baurens F-C, Falentin C, Genete M, Berrabah W, Chèvre A-M, Delourme R: Chromosome-scale assemblies of plant genomes using nanopore long reads and optical maps. *Nature plants* 2018, **4**(11):879-887.
- 31. Jin J-J, Yu W-B, Yang J-B, Song Y, DePamphilis CW, Yi T-S, Li D-Z: GetOrganelle: a fast and versatile toolkit for accurate de novo assembly of organelle genomes. *Genome biology* 2020, **21**:1-31.
- 32. Dierckxsens N, Mardulyn P, Smits G: **NOVOPlasty: de novo assembly of organelle** genomes from whole genome data. *Nucleic acids research* 2017, **45**(4):e18-e18.
- 33. Nikolenko SI, Korobeynikov AI, Alekseyev MA: **BayesHammer: Bayesian clustering for error correction in single-cell sequencing**. In: *BMC genomics: 2013*. Springer: 1-11.
- 34. Jackman SD, Vandervalk BP, Mohamadi H, Chu J, Yeo S, Hammond SA, Jahesh G, Khan H, Coombe L, Warren RL: **ABySS 2.0: resource-efficient assembly of large genomes using a Bloom filter**. *Genome research* 2017, **27**(5):768-777.
- 35. Simpson JT, Wong K, Jackman SD, Schein JE, Jones SJ, Birol I: **ABySS: a parallel** assembler for short read sequence data. *Genome research* 2009, **19**(6):1117-1123.
- 36. Li D, Luo R, Liu C-M, Leung C-M, Ting H-F, Sadakane K, Yamashita H, Lam T-W: **MEGAHIT v1. 0: a fast and scalable metagenome assembler driven by advanced methodologies and community practices.** *Methods* 2016, **102**:3-11.
- 37. Bankevich A, Nurk S, Antipov D, Gurevich AA, Dvorkin M, Kulikov AS, Lesin VM, Nikolenko SI, Pham S, Prjibelski AD: **SPAdes: a new genome assembly algorithm and its applications to single-cell sequencing**. *Journal of computational biology* 2012, **19**(5):455-477.
- 38. Chikhi R, Medvedev P: Informed and automated k-mer size selection for genome assembly. *Bioinformatics* 2014, **30**(1):31-37.
- 39. Koren S, Walenz BP, Berlin K, Miller JR, Bergman NH, Phillippy AM: Canu: scalable and accurate long-read assembly via adaptive k-mer weighting and repeat separation. *Genome research* 2017, **27**(5):722-736.
- 40. Kolmogorov M, Yuan J, Lin Y, Pevzner PA: Assembly of long, error-prone reads using repeat graphs. *Nature biotechnology* 2019, **37**(5):540-546.
- 41. Holley G, Beyter D, Ingimundardottir H, Møller PL, Kristmundsdottir S, Eggertsson HP, Halldorsson BV: **Ratatosk: hybrid error correction of long reads enables** accurate variant calling and assembly. *Genome Biology* 2021, **22**(1):1-22.
- 42. Wick RR, Schultz MB, Zobel J, Holt KE: **Bandage: interactive visualization of de novo genome assemblies**. *Bioinformatics* 2015, **31**(20):3350-3352.
- 43. Tillich M, Lehwark P, Pellizzer T, Ulbricht-Jones ES, Fischer A, Bock R, Greiner S: GeSeq-versatile and accurate annotation of organelle genomes. *Nucleic acids research* 2017, 45(W1):W6-W11.
- 44. Greiner S, Lehwark P, Bock R: OrganellarGenomeDRAW (OGDRAW) version 1.3.
 1: expanded toolkit for the graphical visualization of organellar genomes. Nucleic acids research 2019, 47(W1):W59-W64.
- 45. Lehwark P, Greiner S: **GB2sequin-A file converter preparing custom GenBank files** for database submission. *Genomics* 2019, **111**(4):759-761.
- 46. Beier S, Thiel T, Münch T, Scholz U, Mascher M: **MISA-web: a web server for microsatellite prediction**. *Bioinformatics* 2017, **33**(16):2583-2585.

- 47. Benson G: Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids research* 1999, **27**(2):573-580.
- 48. Frazer KA, Pachter L, Poliakov A, Rubin EM, Dubchak I: VISTA: computational tools for comparative genomics. *Nucleic acids research* 2004, **32**(suppl_2):W273-W279.
- 49. Brudno M, Do CB, Cooper GM, Kim MF, Davydov E, Green ED, Sidow A, Batzoglou S, Program NCS: LAGAN and Multi-LAGAN: efficient tools for large-scale multiple alignment of genomic DNA. *Genome research* 2003, **13**(4):721-731.
- 50. Rozas J, Ferrer-Mata A, Sánchez-DelBarrio JC, Guirao-Rico S, Librado P, Ramos-Onsins SE, Sánchez-Gracia A: **DnaSP 6: DNA sequence polymorphism analysis of large data sets**. *Molecular biology and evolution* 2017, **34**(12):3299-3302.
- 51. Rozas J, Rozas R: DnaSP version 3: an integrated program for molecular population genetics and molecular evolution analysis. *Bioinformatics (Oxford, England)* 1999, **15**(2):174-175.
- 52. Amiryousefi A, Hyvönen J, Poczai P: **IRscope: an online program to visualize the** junction sites of chloroplast genomes. *Bioinformatics* 2018, **34**(17):3030-3031.
- 53. Chen C, Chen H, Zhang Y, Thomas HR, Frank MH, He Y, Xia R: **TBtools: an** integrative toolkit developed for interactive analyses of big biological data. *Molecular plant* 2020, **13**(8):1194-1202.
- 54. Katoh K, Misawa K, Kuma Ki, Miyata T: **MAFFT: a novel method for rapid multiple sequence alignment based on fast Fourier transform**. *Nucleic acids research* 2002, **30**(14):3059-3066.
- 55. Stamatakis A: **The RAxML v8. 2. X Manual**. Heidleberg Institute for Theoretical Studies Available at: <u>https://cme</u> h-its org/exelixis/resource/download/NewManual pdf 2016.
- 56. Konishi T, Sasaki Y: **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* 1994, **91**(9):3598-3601.
- 57. Wu S, Chen J, Li Y, Liu A, Li A, Yin M, Shrestha N, Liu J, Ren G: Extensive genomic rearrangements mediated by repetitive sequences in plastomes of Medicago and its relatives. *BMC Plant Biol* 2021, **21**(1):421.
- 58. Li J, Su Y, Wang T: The Repeat Sequences and Elevated Substitution Rates of the Chloroplast accD Gene in Cupressophytes. *Front Plant Sci* 2018, **9**:533.
- 59. Gurdon C, Maliga P: Two distinct plastid genome configurations and unprecedented intraspecies length variation in the accD coding region in Medicago truncatula. DNA Res 2014, 21(4):417-427.
- 60. Richardson AO, Palmer JD: Horizontal gene transfer in plants. J Exp Bot 2007, 58(1):1-9.
- 61. de Vries J, Sousa FL, Bolter B, Soll J, Gould SB: **YCF1: A Green TIC?** *Plant Cell* 2015, **27**(7):1827-1833.
- 62. Filip E, Skuza L: Horizontal Gene Transfer Involving Chloroplasts. Int J Mol Sci 2021, **22**(9).
- 63. Zhong Q, Yang S, Sun X, Wang L, Li Y: The complete chloroplast genome of the Jerusalem artichoke (Helianthus tuberosus L.) and an adaptive evolutionary analysis of the ycf2 gene. *PeerJ* 2019, 7:e7596.
- 64. Antil S, Abraham JS, Sripoorna S, Maurya S, Dagar J, Makhija S, Bhagat P, Gupta R, Sood U, Lal R *et al*: **DNA barcoding, an effective tool for species identification: a review**. *Mol Biol Rep* 2023, **50**(1):761-775.

- 65. Li J, Wang Y, Cui J, Wang W, Liu X, Chang Y, Yao D, Cui J: **Removal effects of** aquatic plants on high-concentration phosphorus in wastewater during summer. J Environ Manage 2022, **324**:116434.
- 66. Soliman AT, Hamdy RS, Hamed AB: Ludwigia stolonifera (Guill. & Perr.) PH Raven, insight into its phenotypic plasticity, habitat diversity and associated species. *Egyptian Journal of Botany* 2018, **58**(3):605-626.
- 67. Kamoshita A, Ikeda H, Yamagishi J, Lor B, Ouk M: Residual effects of cultivation methods on weed seed banks and weeds in Cambodia. Weed Biology and Management 2016, 16(3):93-107.
- 68. Anita VPD, Matra DD, Siregar UJ: Chloroplast genome draft assembly of Falcataria moluccana using hybrid sequencing technology. *BMC Res Notes* 2023, **16**(1):31.
- 69. Xu S, Teng K, Zhang H, Gao K, Wu J, Duan L, Yue Y, Fan X: Chloroplast genomes of four Carex species: Long repetitive sequences trigger dramatic changes in chloroplast genome structure. *Front Plant Sci* 2023, **14**:1100876.
- 70. Guo YY, Yang JX, Li HK, Zhao HS: Chloroplast Genomes of Two Species of Cypripedium: Expanded Genome Size and Proliferation of AT-Biased Repeat Sequences. *Front Plant Sci* 2021, **12**:609729.
- 71. Liu S-H, Edwards C, Hoch PC, Raven PH, Barber JC: **Complete plastome sequence** of Ludwigia octovalvis (Onagraceae), a globally distributed wetland plant. *Genome* announcements 2016, 4(6):e01274-01216.
- Bateman RM, Rudall PJ, Murphy ARM, Cowan RS, Devey DS, Perez-Escobar OA: Whole plastomes are not enough: phylogenomic and morphometric exploration at multiple demographic levels of the bee orchid clade Ophrys sect. Sphegodes. J Exp Bot 2021, 72(2):654-681.
- 73. Lin Z, Zhou P, Ma X, Deng Y, Liao Z, Li R, Ming R: Comparative analysis of chloroplast genomes in Vasconcellea pubescens A.DC. and Carica papaya L. Sci Rep 2020, 10(1):15799.
- 74. Lihodeevskiy GA, Shanina EP: The Use of Long-Read Sequencing to Study the Phylogenetic Diversity of the Potato Varieties Plastome of the Ural Selection. Agronomy 2022, 12(4):846.
- 75. Nath O, Fletcher SJ, Hayward A, Shaw LM, Masouleh AK, Furtado A, Henry RJ, Mitter N: A haplotype resolved chromosomal level avocado genome allows analysis of novel avocado genes. *Hortic Res* 2022, **9**:uhac157.
- 76. Wanichthanarak K, Nookaew I, Pasookhush P, Wongsurawat T, Jenjaroenpun P, Leeratsuwan N, Wattanachaisaereekul S, Visessanguan W, Sirivatanauksorn Y, Nuntasaen N: Revisiting chloroplast genomic landscape and annotation towards comparative chloroplast genomes of Rhamnaceae. BMC Plant Biology 2023, 23(1):59.
- 77. Luo Y, He J, Lyu R, Xiao J, Li W, Yao M, Pei L, Cheng J, Li J, Xie L: Comparative Analysis of Complete Chloroplast Genomes of 13 Species in Epilobium, Circaea, and Chamaenerion and Insights Into Phylogenetic Relationships of Onagraceae. *Front Genet* 2021, 12:730495.
- 78. Zhang XF, Landis JB, Wang HX, Zhu ZX, Wang HF: Comparative analysis of chloroplast genome structure and molecular dating in Myrtales. *BMC Plant Biol* 2021, **21**(1):219.
- 79. Xu J, Shen X, Liao B, Xu J, Hou D: Comparing and phylogenetic analysis chloroplast genome of three Achyranthes species. *Sci Rep* 2020, **10**(1):10818.
- 80. Lian C, Yang H, Lan J, Zhang X, Zhang F, Yang J, Chen S: **Comparative analysis of chloroplast genomes reveals phylogenetic relationships and intraspecific variation in the medicinal plant Isodon rubescens**. *PLoS One* 2022, **17**(4):e0266546.

- 81. Mohanta TK, Mishra AK, Khan A, Hashem A, Abd Allah EF, Al-Harrasi A: Gene Loss and Evolution of the Plastome. *Genes (Basel)* 2020, **11**(10).
- 82. Barthet MM, Hilu KW: Expression of matK: functional and evolutionary implications. *American Journal of Botany* 2007, **94**(8):1402-1412.
- 83. Li L, Liu C, Hou K, Liu W: Comparative Analyses of Plastomes of Four Anubias (Araceae) Taxa, Tropical Aquatic Plants Endemic to Africa. *Genes (Basel)* 2022, 13(11).
- 84. Zeb U, Wang X, AzizUllah A, Fiaz S, Khan H, Ullah S, Ali H, Shahzad K: Comparative genome sequence and phylogenetic analysis of chloroplast for evolutionary relationship among Pinus species. Saudi J Biol Sci 2022, 29(3):1618-1627.
- 85. Wu Z, Liao R, Yang T, Dong X, Lan D, Qin R, Liu H: Analysis of six chloroplast genomes provides insight into the evolution of Chrysosplenium (Saxifragaceae). *BMC Genomics* 2020, **21**(1):621.
- 86. Kode V, Mudd EA, Iamtham S, Day A: The tobacco plastid accD gene is essential and is required for leaf development. *Plant J* 2005, 44(2):237-244.
- 87. Madoka Y, Tomizawa K, Mizoi J, Nishida I, Nagano Y, Sasaki Y: 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* 2002, **43**(12):1518-1525.
- 88. Gu H, Wang Y, Xie H, Qiu C, Zhang S, Xiao J, Li H, Chen L, Li X, Ding Z: Drought stress triggers proteomic changes involving lignin, flavonoids and fatty acids in tea plants. *Sci Rep* 2020, **10**(1):15504.
- Bharadwaj B, Mishegyan A, Nagalingam S, Guenther A, Joshee N, Sherman SH, Basu C: Physiological and genetic responses of lentil (Lens culinaris) under flood stress. Plant Stress 2023:100130.
- 90. Kikuchi S, Asakura Y, Imai M, Nakahira Y, Kotani Y, Hashiguchi Y, Nakai Y, Takafuji K, Bedard J, Hirabayashi-Ishioka Y *et al*: A Ycf2-FtsHi Heteromeric AAA-ATPase Complex Is Required for Chloroplast Protein Import. *Plant Cell* 2018, **30**(11):2677-2703.
- 91. Schreier TB, Clery A, Schlafli M, Galbier F, Stadler M, Demarsy E, Albertini D, Maier BA, Kessler F, Hortensteiner S *et al*: Plastidial NAD-Dependent Malate Dehydrogenase: A Moonlighting Protein Involved in Early Chloroplast Development through Its Interaction with an FtsH12-FtsHi Protease Complex. *Plant Cell* 2018, **30**(8):1745-1769.
- 92. Drescher A, Ruf S, Calsa T, Jr., Carrer H, Bock R: The two largest chloroplast genome-encoded open reading frames of higher plants are essential genes. *Plant J* 2000, **22**(2):97-104.
- 93. Xing J, Pan J, Yi H, Lv K, Gan Q, Wang M, Ge H, Huang X, Huang F, Wang Y *et al*: **The plastid-encoded protein Orf2971 is required for protein translocation and chloroplast quality control**. *Plant Cell* 2022, **34**(9):3383-3399.
- 94. Chen J, Zang Y, Shang S, Yang Z, Liang S, Xue S, Wang Y, Tang X: Chloroplast genomic comparison provides insights into the evolution of seagrasses. *BMC Plant Biol* 2023, **23**(1):104.
- 95. Xie Z, Merchant S: The plastid-encoded ccsA gene is required for heme attachment to chloroplast c-type cytochromes. *J Biol Chem* 1996, **271**(9):4632-4639.
- 96. Kranz R, Lill R, Goldman B, Bonnard G, Merchant S: Molecular mechanisms of cytochrome c biogenesis: three distinct systems. *Mol Microbiol* 1998, **29**(2):383-396.

- 97. Maheswari P, Kunhikannan C, Yasodha R: Chloroplast genome analysis of Angiosperms and phylogenetic relationships among Lamiaceae members with particular reference to teak (Tectona grandis L.f). *J Biosci* 2021, 46.
- 98. Zhang Y, Du L, Liu A, Chen J, Wu L, Hu W, Zhang W, Kim K, Lee SC, Yang TJ et al: The Complete Chloroplast Genome Sequences of Five Epimedium Species: Lights into Phylogenetic and Taxonomic Analyses. Front Plant Sci 2016, 7:306.
- 99. Huang LS, Sun YQ, Jin Y, Gao Q, Hu XG, Gao FL, Yang XL, Zhu JJ, El-Kassaby YA, Mao JF: Development of high transferability cpSSR markers for individual identification and genetic investigation in Cupressaceae species. Ecol Evol 2018, 8(10):4967-4977.
- 100. Leontaritou P, Lamari FN, Papasotiropoulos V, Iatrou G: Exploration of genetic, morphological and essential oil variation reveals tools for the authentication and breeding of Salvia pomifera subsp. calycina (Sm.) Hayek. *Phytochemistry* 2021, 191:112900.
- 101. Snoussi M, Riahi L, Ben Romdhane M, Mliki A, Zoghlami N: Chloroplast DNA Diversity of Tunisian Barley Landraces as Revealed by cpSSRs Molecular Markers and Implication for Conservation Strategies. Genet Res (Camb) 2022, 2022:3905957.
- 102. Song SL, Lim PE, Phang SM, Lee WW, Hong DD, Prathep A: Development of chloroplast simple sequence repeats (cpSSRs) for the intraspecific study of Gracilaria tenuistipitata (Gracilariales, Rhodophyta) from different populations. *BMC Res Notes* 2014, 7:77.
- 103. Wheeler GL, Dorman HE, Buchanan A, Challagundla L, Wallace LE: A review of the prevalence, utility, and caveats of using chloroplast simple sequence repeats for studies of plant biology. *Appl Plant Sci* 2014, **2**(12).
- 104. Feng JL, Wu LW, Wang Q, Pan YJ, Li BL, Lin YL, Yao H: Comparison Analysis Based on Complete Chloroplast Genomes and Insights into Plastid Phylogenomic of Four Iris Species. *Biomed Res Int* 2022, 2022:2194021.

Figures



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



Figure 2: Comparative results of *L. grandiflora* 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.



Figure 3: Circular representation of annotations plastomes in *Ludwigia octovalis*, *Ludwigia grandiflora* and *Ludwigia peploides* 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* is represented as haplotype 2 only diverge by the orientation of the SSC segment.



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* and *L. grandiflora* (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* and *L. grandiflora* haplotype 1) and haplotype 2 (*L. octovalvis* and *L. grandiflora* haplotype 2) plastomes.



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.



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). 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.



Figure 7: The Ka/Ks ratios of the 80 protein-coding genes of *Ludwigia* sp. plastomes. The blue curve represents *L. grandiflora* versus *L. peploides*, purple curve denotes *L. grandiflora* versus *L. octovalvis* and green curve *L. peploides* 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 the *Ludwigia matK* marker. Only four sequences are complete (499 aa, yellow star), the others correspond to barcode amplification ranging from 128 to 386 aa, with an average of 244 aa.

Additional Figures

r	٦.

В

	ABySS		MEGAHIT		VELVET		SPAdes	
	not corrected	corrected	not corrected	corrected	not corrected	corrected	not corrected	corrected
				Using all c	ontigs		ļ	
Genome fraction (%) Duplication ratio Largest alignment	86.868 1.047 56 588	85.279 1.042 41262	86.428 1.796 30 904	85.158 1.041 90352	91.927 2.002 3531	86.796 1.128 17235	84.682 1.042 90 399	84.483 1 90 272
			Us	sing contigs	s > 200 nt			
Genome fraction (%) Duplication ratio Largest alignment	86.419 1.028 56 588	85.279 1.042 41262	86.377 1.681 30 904	85.057 1.029 90 352	76.589 1.177 3531	86.181 1.11 17235	84.682 1.042 90 399	84.483 1 90 272
			Us	ing contigs	> 500 nt			
Genome fraction (%) Duplication ratio Largest alignment	85.564 1.009 56588	84.517 1.012 41262	83.503 1.041 30 904	84.774 1.004 90 352	45.468 1.015 3531	79.279 1.054 17235	84.682 1.042 90 399	84.483 1 90272
			Us	ing contigs	> 1000 nt			
Genome fraction (%) Duplication ratio Largest alignment	83.701 1 56 588	84.199 1.002 41262	81.256 1.007 30 904	84.545 1.001 90352	22.194 1 3531	66.438 1.011 17235	84.563 1.026 90 399	84.483 1 90272
	ABySS MEGAHIT		AHIT	VELVET		SPAdes		
	not corrected	corrected	not corrected corrected		not corrected corrected		not corrected corrected	
				Using all c	ontigs			
NGA50 LGA50	15 215 3	26 577 3	19 986 3	90 352 1	469 93	2796 9	90 399 1	90 272 1
Misassemblies								
# misassemblies Misassembled contigs length	0 0	0 0	4 1595	0 0	0 0	0 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 # N's per 100 kbp	12.4 0	10.58 0	62.99 0	16.26 0	27.92 0	74.88 6.1	19.17 0	0
			Us	ing contigs	> 500 nt			
NGA50 LGA50	15 215 3	26 577 3	19 986 3	90 352 1	-	2796 9	90 399 1	90272 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 # indels per 100 kbp # N's per 100 kbp	62.39 2.9 0	46.17 2.2 0	123.32 4.33 0	8.1 1.47 0	221.33 28.51 0	148.48 63.74 7.22	96.57 19.17 0	0 0 0
Using contigs > 100								
NGA50 LGA50	15 215 3	26 577 3	19 986 3	90 352 1	-	2796 9	90 399 1	90 272 1
Misassemblies								
# misassemblies Misassembled contigs length	0 0	0 0	0 0	0 0	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 # N's per 100 kbp	1.5 0	0.74 0	0	0.74 0	25.41 0	56.93 9.25	19.49 0	0

Add. Figure 1: QUAST 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.



Add. Figure 2: 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



Add. Figure 3: 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.



Add. Figure 4: 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, *Epibolium* 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.



Add. Figure 5: 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.

В



Α

Add. Figure 6: Diagram showing the position of tandem repeats in the *accD* gene. *L. octovalis* (in red) and *L. peploides* and *L. grandiflora* (in green). We also observe the consequences of these repetitions on the insertion of amino acids, also repeated.



Add. Figure 7: 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.



В



Add. Figure 8: 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.

Tables

	L.octovalvis*	<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				

Table 1. The general characteristics of the three Ludwigia plastomes

* KX827312 (ref)

Table 2 : Genes present in the plastome of Ludwigia sp.

Function	Name						
Photosynthesis							
Rubisco	rbcL						
Photosystem I (PSI)	psaA; psaB; psaC; psal; psaJ						
PSI assembly factors	ycf3 [#] (pafl); ycf4 (pafll)						
Photosystem II	psbA; psbB; psbC; psbD; psbE; psbF; psbH; psbJ; psbJ; psbK; psbL; psbM; pbf1 (psbN) psbT; psbZ						
ATP synthase	atpA; atpB; atpE; atpF [#] ; atpH; atpI						
Cytochrome <i>b6f</i>	petA; petB [#] ; petD [#] ; petG; petL; petN						
Cytochrome biogenesis	ccsA						
NADPH dehydrogenase	ndhA [#] ; ndhB** [#] ; ndhC; ndhD; ndhE; ndhF; ndhG; ndhH; ndhI; ndhJ						
	Transcription and translation						
Transcription	rpoA; rpoB; rpoC1 [#] ; rpoC2						
Small ribosomal proteins	rps2; rps3; rps4; rps7**; rps8; rps11; rps12** [#] ; rps14; rps15; rps16 [#] ; rps18; rps19						
Large ribosomal proteins	rpl2** [#] ; rpl14; rpl16 [#] ; rpl20; rpl22; rpl23**; rpl32; rpl33; rpl36						
Translation initiation	infA						
Ribosomal RNA	m5**; m4;5**; m16**; m23**						
Transfer RNA	tmA-UGC** [#] ;tmC-GCA;tmD-GUC;tmE-UUC;tmF-GAA;tmfM-CAU;tmG-GCC;tmG-UCC [#] ;tmH-GUG;;tmI-CAU**;tmI-GAU** [#] ;tmK-UUU [#] ;tmL-CAA**;tmL-UAA [#] ;tmL-UAG;tmM-CAU;tmN-GUU**;tmP-UGG;tmQ-UUG;tmR-ACG**;tmR-UCU;tmS-GCU;tmS-GGA;tmS-UGA;tmT-GGU;tmT-UGU;tmV-GAC**;tmV-UAC [#] ;tmW-CCA;tmY-GUA						
Other functions							
Group II intron splicing	matK						
Inorganic carbon uptake	cemA						
Protease	clpP1 [#]						
Fatty acid synthesis/Heat tolerance	accD						
TIC machinery (protein import)	ycf1 (Tic214); ycf2**						
Unknown function pseudogene	ycf15**						
	** duplicated in IR region; # spliced genes						

Table 3A : Tandem repeats detected on *Ludwigia* sp. plastomes

Sequence	Lo	Lgh	Lpm	Length	Region	Locus	Comments
TTGTAGTCAGGGGTGTAGTACTAT				24	IRs	ycf2	
TAGAAGAGAGTGCAG		Х	Х	15	IRs	ycf2	15 nt deletion in Lgh and Lpm
ATGAAATATCGTATAATGAAGTACCACACGAGTGGATAT	Х	X		39	IRs	rpl2 intron	39 nt deletion in Lgh and Lo
AAAAATAGGATAGGAT		х	Х	16	LSC	ycf1-tmH-GUG	56 nt deletion in Lgh and Lpm
ΤΑΑΑΤΤΑΑΤΑΤΟΤΑΤΑΤΑ		х	х	18	LSC	psbZ-trnG-GCC	18 nt deletion in Lgh and Lpm
TTTTCTATCTATCTTATATCAA		х	х	22	LSC	tmK-UUU-rps16	22 nt deletion in Lgh and Lpm
AGATCCATAACATCATCAAA		х	х	20	LSC	rps16 intron	22 nt deletion in Lgh and Lpm
TATTAGTTATTAATATTATTAGA		х	х	23	LSC	trnP-UGG-psaJ	23 nt deletion in Lgh and Lpm
ΑΑΤΑΑΤΑΤΑΤΑΑΤΑΑCTTAAATA		х	х	23	LSC	rpl33-rps18	33 et 44 nt nt deletion in in Lgh et Lpm, respectively
TTTTTATTTAACATGCTATCAAATCAACAATGCCATACCGTAGGGCATCTGTT		х	х	53	LSC	rpl20-clpP1	107 nt deletion in Lgh and Lpm
ATATATTTCGATTCAATTC	Х		Х	19	LSC	tmH-GUG-psbA	3 copies in a 57 nt deletion in Lo and Lpm
ATAGAAATATCAGTATTTGAGTG	Х		Х	23	LSC	atpH-atpI	23 nt deletion in Lo and Lpm
TTAATTTTAATTGAAGAA	Х		Х	18	LSC	psbJ-psbL	17 and 24 nt deletion in Lo and Lpm, respectively
TTAAAGAATATTAATATTC	imperfect TR			19	LSC	tmR-UCU-atpA	A -> C mutation in second copy in Lo
ΤΑΤΤΑΤΤΑΤΤΑΤΤΑΑΤ	Х	Х		16	LSC	atpH-atpI	16 nt deletion in Lgh and Lo
TCTAAGGCTGAAATAAGG	Х	Х		18	LSC	pafl intron	18 nt deletion in Lgh and Lo
TGTGAATCTATCTAT			Х	15	LSC	tmS-UGA-psbZ	8 nt deletion in Lpm
TTTTTTCTAGTA				12	LSC	pafl intron	
CTAGTTATTGACATGG		imperfect TR	imperfect TR	16	LSC	psaJ-rpl33	G -> A mutation in second in Lpm et Lgh
ATTTTTATTAACTCT	Х		imperfect TR	15	SSC	ycf1	T->A mutation in first copy in Lpm, other sequence in first copy in Lo
AATCAAATAGTTGAT		Х	X	15	SSC	ycf1	other sequence in first copy of Lpm and Lgh
ΑΤΑΑΤΑΑΤΑΤΑΤΤΤΑΤΤΑΑΤΤΑΑΤΤΑΑΤΑ	X			28	SSC	ndhF-rpl32	160 nt deletion in Lo

Lo = Ludwigia octovalvis; Lgh = L. grandiflora subsp. Hexapetala; Lpm = L. peploides subsp. Montevidensis.

Sequence Lo Lgh Lpm Size (nt) Spacers (nt) Region Locus Comments TTCAATTGGAACGGACGATTCGTCAATCATCT 32 37 SSC ycf1 2 copies. In Lo, one mutation (G->A) in the second copie 35 28 - 22 - 11 LSC CATCGATGATGAAAGTGAAAACAGTAATGAAGAGG Х accD 3 perfects copies and 1 mutated (G->A) copie in Lgh and Lpm . Region of 174 AGATGGTGAAGAACCTTATGAAGATGGTGAAGAACCTTATG Х Х 41 22 LSC accD Region of 147 nt deleted in Lgh and Lpm TATCAAATCAACAATGCCATACCGTAGGGCAT Х Х 32 22 - 21 LSC rps12-clpP1 3 copies 408 in *L.p* , 406 in *L.g* LSC 2 copies. In Lgh, one mutation (C->T) in the second copie TTAAGAGCCGTACAGGCACCTTTTGATGCATACGG Х clpP1 TTAAGAGCCGTACAGGCACTTTTTGATGCATACGG Х 35 LSC clpP1 intron 1- intron 2 Х 811 41 TGCAATAGCCAAATGATGATGAGCAATATCAGTCAGCCATA 2178 LSC psaB-psaA

Table 3B : Direct repeats detected on Ludwigia sp. plastomes

Lo = Ludwigia octovalvis; Lgh = L. grandiflora subsp. Hexapetala; Lpm = L. peploides subsp. Montevidensis.

Table 3C : Palindromic repeats detected on *Ludwigia* sp. plastomes

Common perfect palidromic repeats		Locus	Comment
AGACTCTCATGAGAGTCT ATTAAATAGAATATTCTATTTAAT		trnC-GCA - petN trnE-UUC-trnT-GGI/	
TTGGTAAATTTACCAA		nshD	
TICATTICAATTICAATTIGAAATIGAAATGAA		tml-CALI-vct2	2 conies in IR
GAAAAAGGCCIIITIC		vcf2	2 copies in IR
TCTCAAATGATTAATCATTIGAGA		trol / / 44 intron	
GGATTACTAGTAATCC		tmD GUC tmV GUA	
		traC L/CC introp	
		1111G-0CC -1111R-0CU	
		tmG-GCC-tmtw-CAU	
CLAGIAIGUAIACIGG		nank	
Common palidromic repeats with covariation	· ·	Locus	
in L. octovalvis	in L. grandiflora et L. peploides		
ATAGAATCTATATTCTATTAGAATATAGATTCTAT	ATCGAATCTATATTCTATTAGAATATAGATTCGAT	ndhC-tmV-UAC	
ATGTATATATATCGAT	ATCTATATATAGAT	tmE-UUC-tmT-GGU	
Common palindromic and quasi-palidromic repeats			Comment
in L octavaluie	in Larandiflara and Lananlaides	LUCUS	Comment
		ter B. UQU, etc. A. ter B. UQU, etc. A	
I I I AACGAA IA I TAA IA I I TG I IAAA		tmk-UCU-atpA tmk-UCU-atpA	
ΤΤΑΑ ο GAATATTAATATTCTTTAA			
AATTGTA C TTACAATT	AATTGTAATTACAATT	ccsA	
AGGAAGATTGATCAATCTT t CT	AGGAAGATTGATCAATCTTCCT		
	ΤΤΑ CTA ΑΤΑ ΤΤΑ CTA Α		
ATATAGAATAT C CTATAT	ATATAGAATATTCTATAT		
ACATATCATGATA g GT	ACATATCATGATATGT	rpl22	
AATTACTAATTTCTATTACTATGTTCAATTGAACATAGTAATAGAAATTAGTAATT	AATTACTAATTTCTATTACT + TGTTCAATTGAACATAGTAATAGAAATTAGTAATT	atnH-atn/	
		aipi r aipi	
TAGTTAGAATTCTAACTA	TAGTT c GAATTCTAACTA	tmT-UGU-tmL-UAA	
TATTTTTTCTAGAAAAAATA	TATTTTTTCTAGAA g AAATA	ycf2	2 copies in IR
in L. octovalvis and L. peploides	in L. grandiflora		
	CCCATCAATCATGATTGATGGG	nshN_tmD_GUC	
000410410410410410	00041041041041041000	psuv-uno-000	
in L. octovalvis and L. grandiflora	in L. peploides		
ATGAAAAAAATCGATTTTTTCAT	ATGATAAAAATAGATTTTT a TCAT	tmK-UUU-rps16	
ATGAAAAAAATCGATTTTTTCAT- ATGATAAAAATCGATTTTTATCAT	ATGATAAAAATA gATTTTTATCAT	tmK-UUU-rps16	
Unique palidromic repeats			
l nenlnides			
ΤΤΑΤΑΤΑΤΑΤΑΤΑΤΑΤΑ		m/22 ndhE	Full deletion in / cotourluis & bases deletion in / grandiffers
		Tpi32-Hum	Fuil deletion in 2. Octovarvis, o bases deletion in 2. grandmora
L. octovalvis			
ATTGAAATTCGAATTTCAAT		psbZ-tmG-GCC	Full deletion in L. grandiflora and L. peploides
L. nenloides and L. grandiflora			
		tml 1/AC m/22	2 honors deleted and 2 honors mutated in / satavalvia
		1/1L-0AG-1002	5 bases deleted and 5 bases mutated in <i>L. Octovalvis</i>
AATATATATATATATATATAT		rpi32-ndni-	Full deletion in L. octovalvis
TATATTTATTATTAATTAATAATAATAATAA		rpl32-ndhF	Full deletion in L. octovalvis
L. octovalvis			
ATTGAAATTCGAATTTCAAT		nsh7 tmC CCC	Full deletion in / grandiflora and / panloides
		p302 1110 000	1 di decidi in E. grananora ana E. popolaco
I poplaidae and I amountiling			
L. pepiloles and L. granomora			
AAAAAATGGATCCATTTTT		tmL-UA G-rpi32	3 bases deleted and 3 bases mutated in L. octovalvis
AATATATTATTATAATAATAATATT		rpl32-ndhF	Full deletion in L. octovalvis
τατατττατταττααττααταατααταα		m/32-ndhF	Full deletion in / octovalvis
		ipice nom	