

---

## The complete chloroplast genome of *Ensete glaucum* (Roxb.) Cheesman

---

Nguyen, H. D.<sup>1</sup>, Le, V. M.<sup>2</sup>, Nguyen, P. A. T.<sup>3</sup>, Dinh, H. N.<sup>1</sup>, Phan, N. H.<sup>1</sup>, Vu, M. T.<sup>1</sup> and Do, H. D. K.<sup>1\*</sup>

<sup>1</sup>NTT Hi-Tech Institute, Nguyen Tat Thanh University, Ho Chi Minh City, Vietnam; <sup>2</sup>Research Center of Ginseng and Medicinal Materials, National Institute of Medicinal Materials, Ho Chi Minh City, Vietnam; <sup>3</sup> Department of Molecular Biotechnology, Biotechnology Research and Development Institute, Can Tho University, Can Tho City, Vietnam.

Nguyen, H. D., Le, V. M., Nguyen, P. A. T., Dinh, H. N., Phan, N. H., Vu, M. T. and Do, H. D. K. (2022). The complete chloroplast genome of *Ensete glaucum* (Roxb.) Cheesman. International Journal of Agricultural Technology 18(5):2109-2122.

**Abstract** *Ensete glaucum* (Roxb.) Cheesman, also called snow banana, originated in Asia and has ornamental and medicinal value. Results found its complete chloroplast genome, which is 168,483 bp in length and composed of a large single-copy region (LSC; 88,233 bp), a small single-copy region (SSC; 11,138 bp), and two inverted repeat regions (IR; 34,636 bp). The completely sequenced genome includes 135 coding regions of 87 protein-coding genes, 40 tRNAs, and 8 rRNAs. An analysis of repeat composition identified 31 simple sequence repeats and 44 long repeats, mostly in non-coding regions. Notably, the *ycf1* and *ycf2* genes contain various repeats within the coding sequences. A maximum likelihood phylogenetic analysis revealed a close relationship between *E. glaucum* and *Musella lasiocarpa* rather than *Musa* species. Within *E. glaucum*, the Vietnam sample had a chloroplast genome more similar to a sample from Taiwan than the Indian variety.

**Keywords:** Banana evolution, Musaceae, Phylogenetic relationship, Snow banana

### Introduction

Most land plants possess three genomes, the nuclear, mitochondria and chloroplast (cpDNA) genomes (Dobrogojski *et al.*, 2020). The cpDNA is typically quadripartite and ranges from 16 to over 200 kb and contains approximately 130 genes that perform photosynthesis and related metabolism in the chloroplast (Daniell *et al.*, 2016). The 1000 Plant Transcriptome and 10,000 Plant Genome Projects were recently initiated (Carpenter *et al.* 2019; Chen *et al.* 2011; Cheng *et al.* 2018). Data have been used to investigate various aspects of plants, such as phylogeny (Gitzendanner *et al.*, 2018; One Thousand Plant Transcriptomes Initiative, 2019) and metabolite pathways (Chakraborty, 2018; He *et al.* 2016; Liu *et al.*, 2020). Specifically, a one-

---

\*Corresponding Author: Do, H. D. K.; Email: [dhdtkhoa@ntt.edu.vn](mailto:dhdtkhoa@ntt.edu.vn)

billion-year history of green plants was proposed from the genetic sequencing of representative plants, to help understand how plants dominate on Earth (Gitzendanner *et al.*, 2018). The genomic data revealed the synthesis pathways of chemical compounds in medicinal plants and unknown metabolomics (Chakraborty, 2018; He *et al.* 2016). Another application of genomic data was to develop molecular markers for trait mapping, plant breeding and identification (Garrido-Cardenas *et al.*, 2018; Grover and Sharma, 2016; Henry, 2012; Semagn *et al.*, 2006; Vu *et al.*, 2021). Such studies have revealed the essential role of genomic data for exploring the evolution of land plants and other living organisms on Earth.

Musaceae Juss. 1789 is a monocotyledonous family that includes three genera: *Musa* L. 1753 (82 species), *Musella* (Franch.) C. Y. Wu ex H. W. Li 1978 (one species), and *Ensete* Bruce ex Horan.1862 (seven species) (POWO 2021). Various genomic studies based on the internal transcribed spacer (ITS) and chloroplast gene sequences (*rps16*, *atpB-rbcL*, and *trnL-F*) have revealed the phylogenetic relationships among these three genera of Musaceae, with *Musella* closer to *Ensete* than to *Musa* (Li *et al.*, 2010; Christelov á *et al.*, 2011; Wu *et al.*, 2021; Feng *et al.*, 2022). A divergence time analysis indicated that Musaceae arose approximately 69 Mya at the Cretaceous–Tertiary boundary (Christelov á *et al.*, 2011). The chloroplast genomes of *Musa*, *Musella*, and *Ensete* species have also been reported (Liu *et al.*, 2018; Yemataw *et al.*, 2018; Zhang *et al.*, 2018; Feng *et al.*, 2020; Song *et al.*, 2022). The entire banana genome was sequenced and provided useful information for the hybridization of banana cultivars (Martin *et al.*, 2016; Martin *et al.*, 2020).

Within Musaceae, *Ensete glaucum* (Roxb.) Cheesman 1948, also called the snow banana or elephant foot banana by the Vietnamese, has a native range from central Nepal to Papuasia (POWO, 2021). This species is being used in various traditional remedies and as an ornamental plant (Inta *et al.*, 2013; Joga *et al.*, 2021; Ochiai 2012). In addition, the microcrystalline cellulose in *E. glaucum* is a potential biomaterial for drug delivery (Pachauu *et al.*, 2019). Such works revealed potential medical applications of *E. glaucum*. In this study, the complete chloroplast genome of *E. glaucum* was sequenced using Oxford Nanopore Technology (ONT) to enlarge the genomic data for further studies of *Ensete* in particular and Musaceae. The *E. glaucum* cpDNA was characterized by size, gene content, and repeats. The genomic data extracted from cpDNA sequences of *Ensete*, *Musa* and *Musella* were also used to reconstruct the phylogenetic relationships among species of Musaceae.

## Materials and methods

Fresh *E. glaucum* leaves were collected in Ninh Thuan Province, Vietnam (108.767135E, 11.9903180N) and stored in liquid nitrogen. A specimen was deposited at the Research Center of Ginseng and Medicinal Materials under voucher number TTS-MK 231. Total genomic DNA was isolated using the modified CTAB method (Doyle and Doyle 1987). The qualified DNA extract was used to prepare a sequencing library with the SQK-LSK109 ligation sequencing kit (ONT), following the manufacturer's instructions and then sequenced using a single FLO-MIN106 (R9.4) flow cell on a MinION Mk1B device (ONT) for 30 h. The MinKNOW interface was used to monitor the sequencing process and raw reads were base called using Guppy 5 (ONT). Then the processed reads were assembled using Canu and Geneious Prime 2021.1 to complete the chloroplast genome sequence (Koren *et al.* 2017). The obtained chloroplast genome was annotated using Geneious Prime 2021.1 with the reference chloroplast genomes of *Musa ingens* (NCBI accession number MW864253), *Ensete glaucum* (NCBI accession number LC610748), and *Musella lasiocarpa* (NCBI accession number KY807173). The newly sequenced *E. glaucum* chloroplast genome was deposited to NCBI under accession number MZ856381. A map of the chloroplast genome was illustrated using OGDRAW (Greiner *et al.*, 2019).

To identify repeats in the *E. glaucum* chloroplast genome, Phobos embedded in Geneious Prime was used to find simple sequence repeats (SSRs) and REPuter was used to locate long repeats (Kurtz *et al.*, 2001; Mayer, 2006-2010). For the SSRs, the features were set to a minimum length of 10 bp for mononucleotides, 12 bp for dinucleotides, 15 bp for trinucleotides, 16 bp for tetranucleotides, 20 bp for pentanucleotides, and 24 bp for hexanucleotides. For long repeats, a minimum length of 20 bp was selected to identify forward (direct), reverse, complementary, and palindromic repeats.

For phylogenetic analysis, 79 protein-coding regions in 15 chloroplast genomes of 12 Musaceae species (*Musa laterita* [NCBI accession number MW864255], *M. acuminata* subsp. *malaccensis* [HF677508], *M. yunnanensis* [MW864261], *M. itinerans* [NC\_035723], *M. rubinea* [MW864259], *M. nagensium* [MW864258], *M. balbisiana* [NC\_028439], *M. troglodytarum* [MW864260], *M. beccarii* [MK012089], *M. ingens* [MW864253], *Musella lasiocarpa* [KY807173 and LC610747] and *Ensete glaucum* [MZ286962, MZ856381 and LC610748]) were used together with that of *Ravenala madagascariensis* Sonn. 1782 (Strelitziaceae Hutch. 1934) as the outgroup. The sequences were aligned using MUSCLE in Geneious Prime (Edgar, 2004). jModeltest 2.0 was used to estimate the best model for the data matrix, which

was identified as the TVM+G model (Darriba *et al.*, 2012). The IQ-TREE package was used to construct a phylogenetic tree using the maximum likelihood method with 1000 bootstrap replicates (Minh *et al.*, 2020). The phylogenetic tree was illustrated manually using Figtree (<http://tree.bio.ed.ac.uk/software/figtree/>).

## Results

### *Features of the chloroplast genome*

The ONT device produced 1,277,260 raw reads, which ranged from 41 to 27,789 bp. After polishing the raw data, 9,729 reads remained, with lengths from 310 to 19,378 bp. Among the corrected reads, 3,514 (range 338 to 19,024 bp) were assembled to complete the *E. glaucum* chloroplast genome with 91× coverage. The results revealed a typical quadripartite *E. glaucum* chloroplast genome (168,483 bp) that includes large (LSC; 88,233 bp) and small (SSC; 11,138 bp) single-copy regions separated by two inverted repeat regions (IR; 34,636 bp) (Figure 1). This genome sequence had 37% GC content and 135 coding regions, including 87 protein-coding genes, 40 tRNAs, and 8 rRNAs (Table 1). Of the coding regions, 19 parts were duplicated in the IR regions (*rpl2*, *rpl23*, *rps12*, *rps15*, *ndhB*, *ycf1*, *ycf2*, *trnH-GUG*, *trnI-CAU*, *trnL-CAA*, *trnV-GAC*, *trnI-Gau*, *trnR-ACG*, *trnN-GUU*, *rrn4.5*, *rrn5*, *rrn16*, and *rrn23*). The chloroplast genome newly sequenced in this study showed high similarity to the two *E. glaucum* genomes available on NCBI with accession numbers MZ283962 (99.1%) and LC610748 (99.9%). The junction between the LSC and IR regions of *E. glaucum* was located in the *trnH-GUG/rps19* intergenic space (IGS), whereas these junctions in *Musa* and *Musella* extended to the *rps19-rpl22* IGS. However, the LSC/IR junction at *trnH-GUG/rps19* IGS was also found in some *Musa* species, such as *M. beccarii* and *M. troglodytarum*.

The SSR analysis revealed that mononucleotide repeats accounted for 87% of the SSRs in the *E. glaucum* cpDNA (Figure 2A); dinucleotides and trinucleotides comprised 10% and 3%, respectively. There were no tetra-, penta-, or hexanucleotides. Most of the SSRs were located in non-coding regions (77%; Figure 2B). Of the SSRs, 23% were in the *ycf1*, *rps14*, and *rpoC2* coding sequences (Table 2). The *E. glaucum* SSRs ranged from 10 to 17 bp and were made of A and T nucleotides (Table 2).

Similar to the SSRs, the long repeats were mainly in non-coding regions (59%), and forward was the major type (Figure 3A and 3B). The coding regions containing long repeats include *ycf1*, *ycf2*, *psaA*, *psaB*, *accD*, *trnG-UCC*, *trnG-GCC*, *trnS-GCU*, *trnS-UGA*, *trnS-GGA*, *trnM-CAU*, *trnP-UGG*, *trnF-GAA*,

*trnV-UAC*, *trnV-UAC*, and *trnA-UGC* (Table 2). The long repeats ranged from 20 to 48 bp (Table 2).

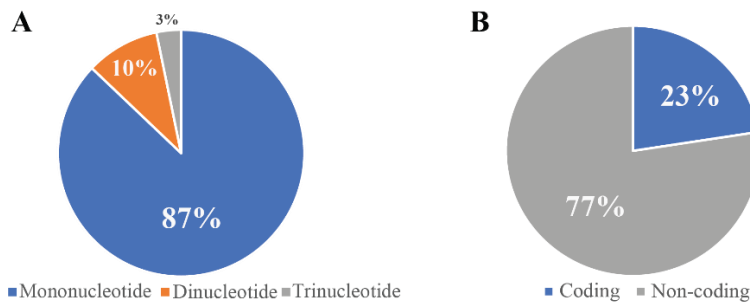
**Table 1.** Gene content of *Ensete glaucum* chloroplast genome

Groups of genes	Names of genes
Ribosomal RNAs	<i>rrn4.5(2x)</i> , <i>rrn5(2x)</i> , <i>rrn16(2x)</i> , <i>rrn23(2x)</i> <i>trnA-UGC*(2x)</i> , <i>trnC-GCA</i> , <i>trnD-GUC</i> , <i>trnE-UUC</i> , <i>trnF-GAA</i> , <i>trnG_UCC*</i> , <i>trnG-GCC</i> , <i>trnH-GUG(2x)</i> , <i>trnI-CAU(2x)</i> , <i>trnI-GAU*(2x)</i> , <i>trnK-UUU*</i> , <i>trnL-UAA*</i> , <i>trnL-UAG</i> , <i>trnL-CAA(2x)</i> , <i>trnM-CAU</i> , <i>trnM-CAU</i> , <i>trnN-GUU(2x)</i> , <i>trnP-UGG</i> , <i>trnQ-UUG</i> , <i>trnR-UCU</i> , <i>trnR-ACG(2x)</i> , <i>trnS-GCU</i> , <i>trnS-UGA</i> , <i>trnS-GGA</i> , <i>trnT-GGU</i> , <i>trnT-UGU</i> , <i>trnV-UAC*</i> , <i>trnV-GAC(2x)</i> , <i>trnW-CCA</i> , <i>trnY-GUA</i>
Transfer RNAs	
Photosystem I	<i>psaA</i> , <i>psaB</i> , <i>psaC</i> , <i>psaI</i> , <i>psaJ</i>
Photosystem II	<i>psbA</i> , <i>psbB</i> , <i>psbC</i> , <i>psbD</i> , <i>psbE</i> , <i>psbF</i> , <i>psbH</i> , <i>psbI</i> , <i>psbJ</i> , <i>psbK</i> , <i>psbL</i> , <i>psbM</i> , <i>psbN</i> , <i>psbT</i> , <i>psbZ</i>
Cytochrome	<i>petA</i> , <i>petB*</i> , <i>petD*</i> , <i>petG</i> , <i>petL</i> , <i>petN</i>
ATP synthases	<i>atpA</i> , <i>atpB</i> , <i>atpE</i> , <i>atpF*</i> , <i>atpH</i> , <i>atpI</i>
Large unit of Rubisco	<i>rbcL</i>
NADH dehydrogenase	<i>ndhA*</i> , <i>ndhB*(2x)</i> , <i>ndhC</i> , <i>ndhD</i> , <i>ndhE</i> , <i>ndhF</i> , <i>ndhG</i> , <i>ndhH</i> , <i>ndhI</i> , <i>ndhJ</i> , <i>ndhK</i>
ATP-dependent protease subunit P	<i>clpP*</i>
Envelope membrane protein	<i>cemA</i>
Large units of ribosome	<i>rpl2*(2x)</i> , <i>rpl14</i> , <i>rpl16*</i> , <i>rpl20</i> , <i>rpl22</i> , <i>rpl23(2x)</i> , <i>rpl32</i> , <i>rpl33</i> , <i>rpl36</i>
Small units of ribosome	<i>rps2</i> , <i>rps3</i> , <i>rps4</i> , <i>rps7(2x)</i> , <i>rps8</i> , <i>rps11</i> , <i>rps12*(2x)</i> , <i>rps14</i> , <i>rps15</i> , <i>rps16*</i> , <i>rps18</i> , <i>rps19</i>
RNA polymerase	<i>rpoA</i> , <i>rpoB</i> , <i>rpoC1*</i> , <i>rpoC2</i>
Initiation factor	<i>infA</i>
Miscellaneous protein	<i>accD</i> , <i>ccsA</i> , <i>matK</i>
Hypothetical proteins conserved reading frames	and <i>ycf1</i> , <i>ycf2(2x)</i> , <i>ycf3*</i> , <i>ycf4</i> , <i>ycf15(2x)</i>

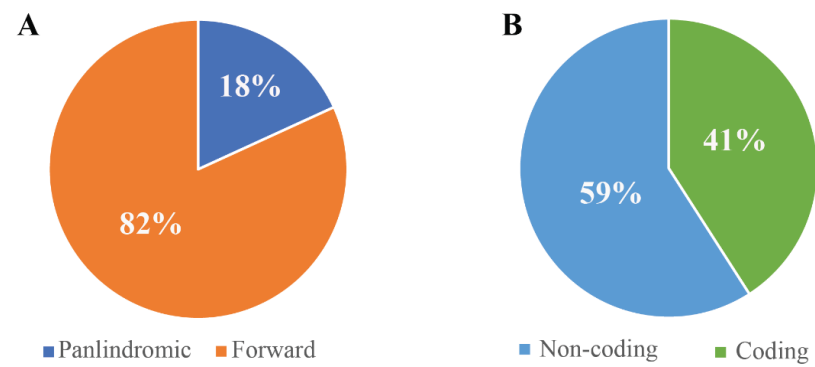
\*- genes with introns; 2x-duplicated genes;  $\Psi$ -pseudogenes.



**Figure 1.** Map of the *Ensete glaucum* chloroplast genome: The genes inside (outside) the circle are transcribed clockwise (counterclockwise); The grey inside the circle indicates the GC content; LSC, large single copy; SSC, small single copy; IRA-IRB, inverted repeat regions



**Figure 2.** Percentage of SSRs in the *Ensete glaucum* chloroplast genome. A) Types and B) Locations of the SSRs



**Figure 3.** Percentage of long repeats in the *Ensete glaucum* chloroplast genome. A) Types and B) Locations of the repeats

**Table 2.** Features of repeats in chloroplast genome of *Ensete glaucum*

Type	Length	Location	Sequence
Dinucleotide	19	<i>rps12-clpP</i>	TATATATATATATATATAT
Dinucleotide	17	<i>trnF-GAA-ndhJ</i>	TATATATATATATATAT
Dinucleotide	12	<i>matK-trnK-UUU</i>	ATATATATATAT
Mononucleotide	16	<i>ycf1*</i>	AAAAAAAAAAAAAAAAAAAA
Mononucleotide	15	<i>rps19-trnH-GUG</i>	TTTTTTTTTTTTTTTTT
Mononucleotide	14	<i>ndhA intron</i>	AAAAAAAAAAAAAAAAAAAA
Mononucleotide	13	<i>trnE-UUC-trnT-GGU</i>	TTTTTTTTTTTTTTTTT
Mononucleotide	13	<i>ycf1*</i>	TTTTTTTTTTTTTTTTT
Mononucleotide	12	<i>atpB-rbcL</i>	TTTTTTTTTTTTTTTTT
Mononucleotide	12	<i>ycf1*</i>	AAAAAAAAAAAAAA
Mononucleotide	11	<i>rps16-trnQ-UUG</i>	AAAAAAAAAAAAAA
Mononucleotide	11	<i>rps2-rpoC2</i>	TTTTTTTTTTTTT
Mononucleotide	11	<i>rps14*</i>	TTTTTTTTTTTTT
Mononucleotide	11	<i>ycf3 intron</i>	AAAAAAAAAAAAAA
Mononucleotide	11	<i>trnF-GAA-ndhJ</i>	AAAAAAAAAAAAAA
Mononucleotide	11	<i>accD-psaI</i>	TTTTTTTTTTTTT
Mononucleotide	10	<i>matK-trnK-UUU</i>	AAAAAAAAAAAAAA
Mononucleotide	10	<i>rps16 intron</i>	AAAAAAAAAAAAAA

**Table 2. (Con.)**

Type	Length	Location	Sequence
Mononucleotide	10	<i>atpH-atpI</i>	TTTTTTTTTT
Mononucleotide	10	<i>rpoC2*</i>	TTTTTTTTTT
Mononucleotide	10	<i>rpoC1 intron</i>	TTTTTTTTTT
Mononucleotide	10	<i>petN-psbM</i>	TTTTTTTTTT
Mononucleotide	10	<i>ycf3-trnS-GGA</i>	AAAAAAAAAA
Mononucleotide	10	<i>trnF-GAA-ndhJ</i>	AAAAAAAAAA
Mononucleotide	10	<i>trnF-GAA-ndhJ</i>	AAAAAAAAAA
Mononucleotide	10	<i>ycf4-cemA</i>	AAAAAAAAAA
Mononucleotide	10	<i>rpl22-rps19</i>	TTTTTTTTTT
Mononucleotide	10	<i>ycf1*</i>	AAAAAAAAAA
Mononucleotide	10	<i>rpl32-trnL-UAG</i>	TTTTTTTTTT
Mononucleotide	10	<i>ycf1*</i>	TTTTTTTTTT
Trinucleotide	15	<i>atpH-atpI</i>	TAATAATAATAATA
Forward	39	<i>ycf2*</i>	TGATAGTGACGATATCGATATTGATGATAGTGACGAT AT
Forward	38	<i>rps12-trnV-GAC</i>	ACTATGAAATTAATATTTCTATAACTATGAAATTAAT A
Forward	36	<i>ycf1*</i>	GCGATGTAGAAAAGTGAGGAAGAAAGCGATGTAGAAA
Forward	34	<i>psaB-psaA*</i>	GTTCTATACATATGACCAGCGATCAGGAAAAGAA
Forward	29	<i>trnF-GAA-ndhJ</i>	AAATAATAATAATAAGTTAAAAAAAAAA
Forward	28	<i>ycf2*</i>	GCTAACTATGACGAATGCGCTAACTATG
Forward	26	<i>trnF-GAA-ndhJ</i>	TAAAAAAAAATAATAATAAATAAGT
Forward	26	<i>ycf1*</i>	TCAAGGCATAAGAAAATTAATAAGTC
Forward	26	<i>ycf1*</i>	ATATTGAGAACAAATAGTAATATTGAG
Forward	25	<i>trnT-UGU-trnL-UAA</i>	AATTATTTCTTAAAACTAACTATTT
Forward	25	<i>ycf4-cemA</i>	CCGAAAAGGACTCTTATTCTATGTC
Forward	25	<i>trnN-GUU-ycf1</i>	ATATAATAATCAAGAAATTGCAATA
Forward	24	<i>trnQ-UUG-psbK</i>	AACACATGTAGATTAGATATAGAA
Forward	24	<i>ycf2*</i>	CCGAAATCTGATTCAAATCCAATA
Forward	24	<i>ycf1*</i>	TGAAAATAATATTAGTGAAAAATA
Forward	23	<i>rps16-trnQ-UUG</i>	ATAATCTAGTTTATCTATTTTCA
Forward	22	<i>trnK-UUU intron</i>	AAGAATAGTTAGGATTCATTAA



**Table 2. (Con.)**

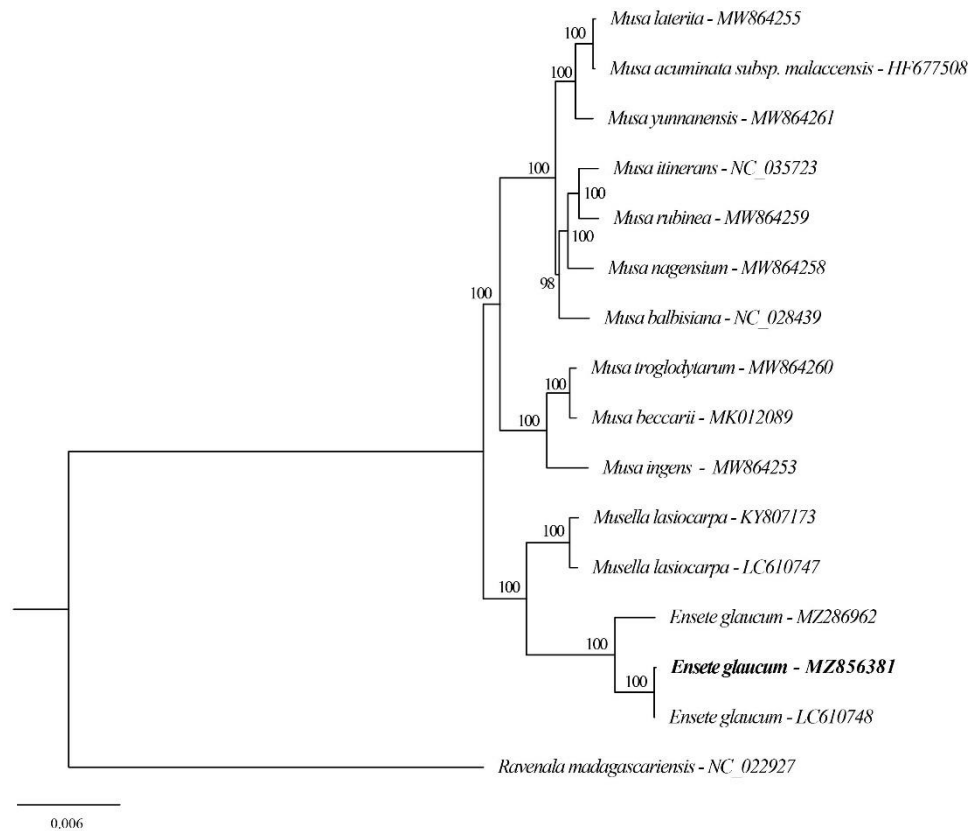
Type	Length	Location	Sequence
Forward	22	<i>trnG-UCC-trnG-GCC*</i>	GATGCGGGTTCGATCCCGCTA
Forward	22	<i>psaB-psaA*</i>	GCAATATCGGTCAGCCATAAAC
Forward	22	<i>ndhC-trnV-UAC</i>	ATCCTATAATTAATACTATGA
Forward	22	<i>rps19-trnH-GUG</i>	AGAAAATCCTTTAGCTAGAAAA
Forward	21	<i>trnK-UUU intron</i>	ACTTACATGAGCATTTTCAGAA
Forward	21	<i>trnS-GCU-trnS-UGA-trnS-GGA*</i>	AGAGAGGGGATTCGAACCCCTCG
Forward	21	<i>petN-psbM</i>	TTTCATTTTCATTTTTTCATTT
Forward	21	<i>trnJ-M-CAU-trnP-UGG*</i>	GACAGGATTTGAACCCGTGAC
Forward	36	<i>accD-psaI</i>	TTACTTATAAAATAAATAATAT
Forward	21	<i>rps8-rpl14</i>	TTGATTATAAAATTATTTATT
Forward	21	<i>rpl16-rps3</i>	TTTATATAGTTATTAAGTTTA
Forward	21	<i>ycf1*</i>	AGTTCTAGTTCTAAAACGAAT
Forward	20	<i>trnH-GUG-psbA</i>	AACAATATTGTATCAAACAA
Forward	20	<i>rps16 intron</i>	AACCTAAGACAAATTAGATT
Forward	20	<i>atpI-rpoC2</i>	TTTTGTTTTTCTTTTTTTT
Forward	20	<i>petA-psbJ</i>	TCTTTTCTTGTTTCTTCGTA
Forward	48	<i>accD*</i>	CAATGATTCCATGAAGAAGTAGAGTCTGATTTCTAT GAAGAAGTAGA
Forward	48	<i>ycf2*</i>	CTTTTTGTCCAAGTCACTTCCCTTTTTGTCCAAGTCAC TCCCTTTTT
Forward	47	<i>ycf1*</i>	AGTCAAAAGAAAAATGAAAATAATATCAGTCAAAAG AAAAATGAAAA
Palindromic	46	<i>rpl32-trnL-UAG</i>	TCTACTTTTCAATAAGAAAAATATTTTCTATTGTGA AAAGTAGA
Palindromic	44	<i>psaC-ndhE</i>	ATATACATTAATAATGTATTATATAATACATTATTAAT GTATAT
Palindromic	31	<i>petA-psbJ</i>	AAGAGTAAGAAAAGAAGTCAACGGGACCTTA
Palindromic	26	<i>petD-rpoA</i>	ATTAATGTATCTAAGAATAGTGACTT
Palindromic	20	<i>trnH-GUG-psbA</i>	AAACAAAGTAGCAATACCCC
Palindromic	20	<i>trnF-GAA-trnV-UAC*</i>	ATAGCTCAGTTGGTAGAGCA
Palindromic	20	<i>trnVUAC-trnA-UGC*</i>	GCTCTACCAACTGAGCTATA
Palindromic	23	<i>petN-psbM</i>	ATAGTATGGTAGAAAAGATTATATAATTCTTTCTA CCATACTAT

Asterisks indicate the repeats in the coding regions

### Phylogenetic relationships

The phylogenetic analysis showed the monophyly of *Musa*, *Musella*, and *Ensete* with high support (Figure 4). In Musaceae, *Ensete* was closer to *Musella* than to *Musa*. Within *Ensete*, the sample from Vietnam was more closely

related to a Taiwan sample (accession no. LC610748) than an Indian sample (accession no. MZ286962).



**Figure 4.** Phylogenetic tree of Musaceae species inferred from 79 protein-coding regions of chloroplast genomes. Numbers at the nodes are the bootstrap values. The chloroplast genome sequenced in this study is in bold italics

## Discussion

The newly sequenced *E. glaucum* chloroplast genome has a quadripartite structure, as reported in the cpDNAs of other land plants and other Musaceae (Daniell *et al.*, 2016; Liu *et al.*, 2018; Yemataw *et al.*, 2018; Zhang *et al.*, 2018; Feng *et al.*, 2020; Song *et al.*, 2022). The gene content and order were similar among Musaceae species, with approximately 135 coding regions, including 87 protein-coding genes, 40 tRNAs, and 8 rRNAs. However, IR region expansion increased the gene number. There were three categories of junctions between IR and LSC regions in *Musa* species, in which the IR expanded to include *rpl2*, *rps19*, and partial *rps19/rpl22* IGS (Wu *et al.*, 2021). Previously, the junctions

of IR and LSC regions of monocots were classified into five types based on the expansion of IR to *rpl2/trnH\_GUG* IGS (type I), *trnH\_GUG/rps19* IGS (type II), *rps19/rpl22* (type III), *rpl22/rps3* IGS (type IV), and from *rps3* (type V) (Do *et al.*, 2020). In Musaceae, the *Musa* species have type I to III IR/LSC borders, whereas *E. glaucum* and *M. lasiocarpa* both have type II. The variable junctions in *Musa* species revealed the dynamic evolution of cpDNA among *Musa* species. Although *Ensete* includes seven species, only the *E. glaucum* cpDNA was characterized here. Therefore, further study should examine all members of *Ensete* to investigate the evolutionary history of the chloroplast genome within *Ensete*.

In the chloroplast genome, repeat sequences play an essential role in structural evolution. Extensive cpDNA rearrangements have been found in different species, such as Pinaceae and Campanulaceae (Cosner *et al.*, 1997; Haberle *et al.*, 2008; Wu *et al.*, 2011). The SSRs can also be used to study genetic populations and molecular markers of plants (Powell *et al.*, 1995; Godwin *et al.*, 1997). In Musaceae, the cpDNA repeats have been analyzed in *Musa* species (Martin *et al.*, 2013; Novák *et al.*, 2014; Song *et al.*, 2022). This study is the first to examine the repeat makeup in *Ensete glaucum* cpDNA. Although no cpDNA structural rearrangement has been reported in Musaceae, repeats provide useful information for further studies of the genetic evolution of Musaceae.

In this study, the phylogenetic analysis showed a close relationship between *Musella* and *Ensete* with high support values, as noted previously (Christelová *et al.*, 2011; Janssens *et al.*, 2016; Wu *et al.*, 2021). Of the seven *Ensete* species, previous studies have used only one *E. glaucum* sample, while the current study included three *E. glaucum* cpDNAs to reconstruct the phylogeny of Musaceae. It showed that the sample from Vietnam (MZ856381) was more closely related to that from Taiwan (LC610748) than that from India (MZ286962). It is suggested that variability among *E. glaucum* chloroplast genomes in its native distribution from Central Nepal to Papuasias and provided valuable information for further genomic studies of *Ensete* species in particular and Musaceae in general.

This study reports the detailed features of the chloroplast genome of *E. glaucum*, including its genome size, gene content and order, SSRs, and long repeats. This information is essential for further studies of the evolutionary history of *Ensete* chloroplast genomes, and those of other Musaceae.

## Acknowledgements

This work was funded by NTTU Foundation for Science & Technology, Nguyen Tat Thanh University under grant number 2021.01.161/HD-KHCN; and by Vingroup Joint Stock

Company through the Domestic Master/ PhD Scholarship Programme of Vingroup Innovation Foundation (VINIF), Vingroup Big Data Institute (VINBIGDATA) to Nguyen Hoang Danh under grant number VINIF.2021.ThS.58; and Department of Science and Technology of Ninh Thuan Province under grant number 11/2020/HĐ-SKHCHN.

## References

- Carpenter, E. J., Matasci, N., Ayyampalayam, S., Wu, S., Sun, J., Yu, J., Vieira, F. R. J., Bowler, C., Dorrell, R. G., Gitzendanner, M. A., Li, L., Du, W., Ullrich, K. K., Wickett, N. J., Barkmann, T. J., Barker, M. S., Leebens-Mack, J. H. and Wong, G. K-S. (2019). Access to RNA-Sequencing Data from 1,173 Plant Species: The 1000 Plant Transcriptomes Initiative (1KP). *GigaScience*, 8:1-7.
- Cosner, M., Jansen, R., Palmer, J. and Downie, S. R. (1997). The highly rearranged chloroplast genome of *Trachelium caeruleum* (Campanulaceae): multiple inversions, inverted repeat expansion and contraction, transposition, insertions/deletions, and several repeat families. *Curr Genet*, 31:419-429.
- Chakraborty, P. (2018). Herbal Genomics as Tools for Dissecting New Metabolic Pathways of Unexplored Medicinal Plants and Drug Discovery. *Biochimie Open*, 6:9-16.
- Chen, S., Xiang, L., Guo, X. and Li, Q. (2011). An Introduction to the Medicinal Plant Genome Project. *Frontiers of Medicine*, 5:178-84.
- Cheng, Shifeng *et al.* Cheng, S., Melkonian, M., Smith, S. A., Brockington, S., Archibald, J. M., Delaux, P-M., Li, F-W., Melkonian, B., Mavrodiev, E.V., Sun, W., Fu, Y., Yang, H., Soltis, D. E., Graham, S. W., Soltis, P. S., Liu, X., Xu, X. and Wong, G. K-S. (2018). 10KP: A Phylodiverse Genome Sequencing Plan. *GigaScience*, 7:1-9.
- Christelov P., Valarik, M., Hribova, E., Langhe E. D. and Dolezel, J. (2011). A multi gene sequence-based phylogeny of the Musaceae (banana) family. *BMC Ecology and Evolution*, 11:103.
- Daniell, H., Lin, C. S., Yu, M. and Chang W. J. (2016). Chloroplast genomes: diversity, evolution, and applications in genetic engineering. *Genome Biology*, 17:134.
- Darriba, D., Taboada, G. L., Doallo, R. and Posada, D. (2012). JModelTest 2: More Models, New Heuristics and Parallel Computing. *Nature Methods*, 9:772-772.
- Do, H. D. K., Kim, C., Chase, M. W. and Kim, J. H. (2020). Implications of plastome evolution in the true lilies (monocot order Liliales). *Molecular Phylogenetics and Evolution*, 148:106818.
- Dobrogojski, J., Adamiec, M. and Lucinski, R. (2020). The chloroplast genome: a review. *Acta Physiologiae Plantarum*, 42:98.
- Doyle, J. J. and Doyle, J. L. (1987). A Rapid DNA Isolation Procedure for Small Quantities of Fresh Leaf Tissue. *Phytochemical bulletin*, 19:11-15.
- Edgar, R. C. (2004). MUSCLE: A Multiple Sequence Alignment Method with Reduced Time and Space Complexity. *BMC Bioinformatics*, 5:113.
- Feng, H., Chen, Y., Li, C., Xu, X., Luo, H. and He, C. (2022). Organelle DNA sequence data provide new insights into the maternal and paternal lineages of *Musa* species germplasms. *Genetic Resources and Crop Evolution*, 69:737-754.
- Feng, H., Chen, Y., Xu, X., Luo, H., Wu, Y. and He, C. (2020) The complete chloroplast genome of *Musa beccarii*. *Mitochondrial DNA Part B*. 5:2384-2385.
- Garrido-Cardenas, J.A., Mesa-Valle, C. and Manzano-Agugliaro, F. (2018). Trends in Plant Research Using Molecular Markers. *Planta*, 247:543-57.
- Godwin, I. D., Aitken, E. A. B. and Smith, L. W. (1997), Application of inter simple sequence

- repeat (ISSR) markers to plant genetics. *ELECTROPHORESIS*, 18:1524-1528.
- Greiner, S., Lehwark, P. and Bock, R. (2019). OrganellarGenomeDRAW (OGDRAW) version 1.3.1: expanded toolkit for the graphical visualization of organellar genomes. *Nucleic Acids Research*, 47:W59-W64
- Grover, A. and Sharma, P. C. (2016). Development and Use of Molecular Markers: Past and Present. *Critical Reviews in Biotechnology*, 36:290-302.
- Gitzendanner, M. A., Soltis, P. S., Wong, G. K., Ruhfel, B. R. and Soltis, D. E. (2018). Plastid Phylogenomic Analysis of Green Plants: A Billion Years of Evolutionary History. *American Journal of Botany*, 105:291-301.
- Haberle, R. C., Fourcade, H. M., Boore, J. L. and Jansen, R. K. (2008). Extensive Rearrangements in the Chloroplast Genome of *Trachelium caeruleum* Are Associated with Repeats and tRNA Genes. *Journal of Molecular Evolution*, 66:350-361.
- He, Y., Xiao, H., Deng, C., Xiong, L., Nie, H. and Peng, C. (2016). Survey of the Genome of *Pogostemon Cablin* Provides Insights into Its Evolutionary History and Sesquiterpenoid Biosynthesis. *Scientific Reports*, 6:26405.
- Henry, R. J. (2012). Evolution of DNA Marker Technology in Plants. In: Robert J. Henry ed. *Molecular Markers in Plants*, Oxford, Blackwell Publishing Ltd., pp. 1-19.
- Inta, A., Trisonthi, P. and Trisonthi, C. (2013). Analysis of Traditional Knowledge in Medicinal Plants Used by Yuan in Thailand. *Journal of Ethnopharmacology*, 149:344-51.
- Janssens, S. B., Vandeloock, F., De Langhe, E., Verstraete, B., Smets, E., Vandenhoutte, I. and Swennen, R. (2016). Evolutionary dynamics and biogeography of Musaceae reveal a correlation between the diversification of the banana family and the geological and climatic history of Southeast Asia. *New Phytol*, 210:1453-1465.
- Joga, R. J., Sangma, E., Karmakar, B., Lyngdoh, V. and Aochen, C. (2021). Phytochemical Investigations on the Therapeutic Properties of *Ensete Glaucum* (Roxb.) Cheesman. *Indian Journal of Traditional Knowledge*, 20:68-73.
- Koren, S., Walenz, B. P., Berlin, K., Miller, J. R., Bergman, N. H. and Phillippy, A. A. (2017). Canu: Scalable and Accurate Long-Read Assembly via Adaptive k-Mer Weighting and Repeat Separation. *Genome Research*, 27:722-36.
- Kurtz, S., Choudhuri, J. V., Ohlebusch, E., Schleiermacher, C., Stoye, J. and Giegerich, R. (2001). REPuter: The Manifold Applications of Repeat Analysis on a Genomic Scale. *Nucleic Acids Research*, 29:4633-4642.
- Li, L. F., Häkkinen, M., Yuan, Y. M., Hao, G. and Ge, X. J. (2010). Molecular phylogeny and systematics of the banana family (Musaceae) inferred from multiple nuclear and chloroplast DNA fragments, with a special reference to the genus *Musa*. *Molecular Phylogenetics and Evolution*, 57:1-10.
- Liu, J., Gao, C. W. and Niu, Y. F. (2018). The complete chloroplast genome sequence of flowering banana, *Musa ornata*. *Mitochondrial DNA Part B Resources*, 3:960-961.
- Liu, X., Zhu, X., Wang, H., Liu, T., Cheng, J. and Jiang, H. (2020). Discovery and Modification of Cytochrome P450 for Plant Natural Products Biosynthesis. *Synthetic and Systems Biotechnology*, 5:187-99.
- Martin, G., Baurens, F. C., Cardi, C., Aury, J. M. and D'Hont, A. (2013). The Complete Chloroplast Genome of Banana (*Musa acuminata*, Zingiberales): Insight into Plastid Monocotyledon Evolution. *PLOS ONE*, 8:e67350.
- Martin, G., Baurens, F. C., Droc, G., Rouard, M., Cenci, A., Kilian, A., Hastie, A., Doležel, J., Aury, J.-M., Alberti, A., Carreel, F. and Angélique D'Hont, A. (2016). Improvement of the banana "*Musa acuminata*" reference sequence using NGS data and semi-automated bioinformatics methods. *BMC Genomics*, 17:243.
- Martin, G., Cardi, C., Sarah, G., Ricci, S., Jenny, C., Fondi, E., Perrier, X., Glaszmann, J.-C.,

- D'Hont, A. and Yahiaoui, N. (2020). Genome ancestry mosaics reveal multiple and cryptic contributors to cultivated banana. *Plant Journal*, 102:1008-1025.
- Mayer, C. (2006-2010). Phobos 3.3.12. Retrieved from [http://www.rub.de/ecoevo/cm/cm\\_phobos.htm](http://www.rub.de/ecoevo/cm/cm_phobos.htm).
- Minh, B. Q., Schmidt, H. A., Chernomor, O., Schrempf, D., Woodhams, M. D., von Haeseler, A. and Lanfear, R. (2020). IQ-TREE 2: New Models and Efficient Methods for Phylogenetic Inference in the Genomic Era. *Molecular Biology and Evolution*, 37:1530-34.
- Novák, P., Hříbová, E., Neumann, P., Koblížková, A., Doležel, J. and Macas, J. (2014). Genome-Wide Analysis of Repeat Diversity across the Family Musaceae. *PLOS ONE*, 9:e98918.
- Ochiai, Y. (2012). From Forests to Homegardens: A Case Study of *Ensete Glaucum* in Myanmar and Laos. *Tropics*, 21:59-66.
- One Thousand Plant Transcriptomes Initiative. (2019). One Thousand Plant Transcriptomes and the Phylogenomics of Green Plants. *Nature*, 574:679-685.
- Pachua, L., Dutta, R. S., Hauzel, L., Devi, T. B. and Deka, D. (2019). Evaluation of Novel Microcrystalline Cellulose from *Ensete Glaucum* (Roxb.) Cheesman Biomass as Sustainable Drug Delivery Biomaterial. *Carbohydrate Polymers*, 206:336-643.
- Powell, W., Morgante, M., McDevitt, R., Vendramin, G. G. and Rafalski, J. A. (1995). Polymorphic simple sequence repeat regions in chloroplast genomes: Applications to the population genetics of pines. *The Proceedings of the National Academy of Sciences*, 92:7759-63.
- POWO. (2021). "Plants of the World Online. Facilitated by the Royal Botanic Gardens, Kew." Retrieved from <http://www.plantsoftheworldonline.org>.
- Semagn, K., Bjørnstad, A. and Ndjondjop, M. N. (2006). An Overview of Molecular Marker Methods for Plants. *African Journal of Biotechnology*, 5:2540-68.
- Song, W., Ji, C., Chen, Z., Cai, H., Wu, X., Shi, C. and Wang, S. (2022). Comparative Analysis the Complete Chloroplast Genomes of Nine *Musa* Species: Genomic Features, Comparative Analysis, and Phylogenetic Implications. *Frontiers in Plant Science*, 13:832884
- Vu, T. N., Pham, L. B. H., Nguyen, N. L., Luu, H. L., Huynh, T. T. H., Nguyen, H. H., Ha, H. H. and Le, T. T. H. (2021). Molecular Markers for Analysis of Plant Genetic Diversity. *Vietnam Journal of Biotechnology*, 18:589-608.
- Wu, C. S., Lin, C. P., Hsu, C. Y., Wang, R. J. and Chaw, S. M. (2011). Comparative Chloroplast Genomes of Pinaceae: Insights into the Mechanism of Diversified Genomic Organizations. *Genome Biology and Evolution*, 3:309-319.
- Wu, C. S., Sudianto, E., Chiu, H. L., Chao, C. P. and Chaw, S. M. (2021). Reassessing Banana Phylogeny and Organelle Inheritance Modes Using Genome Skimming Data. *Frontiers in Plant Science*, 12:713216.
- Yemataw, Z., Muzemil, S., Ambachew, D., Tripathi, L., Tesfaye, K., Chala, A., Farbos, A., O'Neill, P., Moore, K., Grant, M. and Studholme, D. J. (2018). Genome sequence data from 17 accessions of *Ensete ventricosum*, a staple food crop for millions in Ethiopia. *Data in Brief*, 18:285-293.
- Zhang, L., Guo, X., Wang, Z., Wang, M. and Hu, Q. (2018). Characterization of the complete chloroplast genome of *Musella lasiocarpa*. *Mitochondrial DNA Part B: Resources*, 3:728-729.

(Received: 13 April 2022, accepted: 30 July 2022)