



ELSEVIER

Contents lists available at ScienceDirect

Data in brief

journal homepage: www.elsevier.com/locate/dib

Data Article

Genome assembly and phylogenomic data analyses using plastid data: Contrasting species tree estimation methods

D.J.P. Gonçalves^{a,*}, B.B. Simpson^a, G.H. Shimizu^b,
R.K. Jansen^{a,c}, E.M. Ortiz^d^a Department of Integrative Biology, The University of Texas at Austin, 2415 Speedway #C0930, Austin, TX, 78713, USA^b Department of Plant Biology, University of Campinas, 13083-970, Campinas, SP, Brazil^c Genomics and Biotechnology Research Group, Department of Biological Sciences, Faculty of Science, King Abdulaziz University, Jeddah, 21589, Saudi Arabia^d Department of Ecology & Ecosystem Management, Plant Biodiversity Research, Technical University of Munich, Emil-Ramann Strasse 2, Freising, D-85354, Germany

ARTICLE INFO

Article history:

Received 14 June 2019

Received in revised form 4 July 2019

Accepted 9 July 2019

Available online 27 July 2019

Keywords:

Data processing

Genome assembly

Phylogenetic analyses

Phylogenetic signal

Tree space

ABSTRACT

Phylogenomics has become increasingly popular in recent years mostly due to the increased affordability of next generation sequencing techniques. Phylogenomics has sparked interest in multiple fields of research, including systematics, ecology, epidemiology, and even personalized medicine, agriculture and pharmacy. Despite this trend, it is usually difficult to learn and understand how the analyses were done, how the results were obtained, and most importantly, how to replicate the study. Here we present the data and all of the code utilized to perform phylogenomic inferences using plastome data: from raw data to extensive phylogenetic inference and accuracy assessment. The data presented here utilizes plastome sequences available on GenBank (accession numbers of 94 species are available below) and the code is also available at <https://github.com/deisejpg/rosids>. Gonçalves et al. is the research article associated with the data analyses presented here.

© 2019 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

DOI of original article: <https://doi.org/10.1016/j.jympev.2019.05.022>.

* Corresponding author.

E-mail address: deisejpg@gmail.com (D.J.P. Gonçalves).

<https://doi.org/10.1016/j.dib.2019.104271>

2352-3409/© 2019 The Author(s). Published by Elsevier Inc. This is an open access article under the CC BY license (<http://creativecommons.org/licenses/by/4.0/>).

Specifications table

Subject area	Biology
More specific subject area	Systematics of angiosperms, plastome evolution
Type of data	Table, PHYLIP, NEXUS, and MARKDOWN formatted files
How data was acquired	Illumina HiSeq. 2500 and Illumina HiSeq 4000
Data format	Raw and Analyzed
Experimental factors	Total genomic DNA isolated and sequenced from fresh or silica dried leaf tissue of sampled species of rosids and genome sequences from NCBI
Experimental features	Illumina reads preprocessing, genome assembly, phylogenetic inference, topology tests, and phylogenetic signal assessment
Data accessibility	Within this article and at https://github.com/deisejpg/rosids
Related research article	Deise J.P. Gonçalves, Beryl B. Simpson, Edgardo M. Ortiz, Gustavo H. Shimizu, Robert K. Jansen Incongruence between gene trees and species trees and phylogenetic signal variation in plastid genes Molecular Phylogenetics and Evolution 10.1016/j.ympev.2019.05.022

Value of the data

- The present data provides details about phylogenomic analysis using a set of well-documented pipelines covering analyses from preprocessing Illumina reads to inferring and testing phylogenies using multiple methods of phylogenetic inference
- The data introduce the practical use of multispecies coalescent methods using plastid protein-coding genes and could be adjusted and used with molecular data from different molecular markers and organisms
- Accessibility to scripts utilized and data files containing the alignments and trees will enhance the replication of the analyses presented

1. Data

A dataset comprising 78 plastid protein-coding genes of 94 species of rosids is presented in Table 1. Here we present all the code used in the analysis of this dataset [1], including the scripts used to quality filter, assemble, extract regions of interest, and perform phylogenomic analysis, in a series of tutorial-like files: I. Genome assembly; II. Phylogenetic Analysis; III. Tree space; IV. Phylogenetic Signal. Part of the data was obtained from GenBank (<http://www.ncbi.nlm.nih.gov/genbank>). Data for and 27 species from groups of rosids that lacked the information on the database were generated using Illumina HiSeq. A total of 657,471,631 million paired-end reads with an average length of 150 bp was generated (Table 2). Despite the interest on extracting and using only the genes from the plastome, the pipeline for genome assembly presented here also separates contigs from the three cellular genomic compartments with a potential for use used in studies that target not only the plastome, but also mitochondria, nuclear ribosomal DNA, and other nuclear markers. The next set of tutorial-like markdown files present the code utilized for preparing alignments and for inferring phylogenies using an array of strategies of data partition and methods of phylogenetic inference. The code used to explore the similarities/dissimilarities between topologies and for the phylogenetic signal calculation is also presented.

2. Experimental design, materials and methods*2.1. Data preprocessing and genome assembly*

For the 27 samples for which data were generated, leaf tissue was ground and total genomic DNA was isolated using DNeasy Plant Mini Kit (Qiagen) according to the manufacturer's protocol or a modified version of [2] described in Ref. [3]. The DNA was quantified using a Qubit Fluorometric Quantitation (Thermo Fisher) instrument and was sequenced at the Genome Sequencing and Analysis

Table 1

Classification according to APG IV (2016) of samples used in the study (89 samples of rosids, considered here as fabids + malvids, and five of outgroup), voucher information of newly sequenced plastomes and GenBank accession numbers. Bold font indicates plastid genome sequences generated in this study.

Order	Family	Species	Voucher ID	GenBank Accession Numbers
Ingroup				
Brassicales	Brassicaceae	<i>Brassica napus</i>		NC_016734
Brassicales	Caricaceae	<i>Carica papaya</i>		NC_010323
Brassicales	Moringaceae	<i>Moringa oleifera</i>	CONN-129179	MK726020
Brassicales	Salvadoraceae	<i>Azima tetraacantha</i>	CONN00225893	MK726028
Celastrales	Celastraceae	<i>Euonymus japonicus</i>		NC_028067.1
Cucurbitales	Cucurbitaceae	<i>Citrullus lanatus</i>		NC_032008
Cucurbitales	Cucurbitaceae	<i>Cucumis hystrix</i>		NC_023544
Cucurbitales	Cucurbitaceae	<i>Gynostemma pentaphyllum</i>		NC_029484
Fabales	Fabaceae	<i>Cicer arietinum</i>		NC_011163
Fabales	Fabaceae	<i>Inga leiocalycina</i>		NC_028732
Fabales	Fabaceae	<i>Lupinus luteus</i>		NC_023090
Fabales	Polygalaceae	<i>Polygala alba</i>	TEX-DJPG731	MK726019
Fagales	Betulaceae	<i>Ostrya rehderiana</i>		NC_028349
Fagales	Fagaceae	<i>Castanea mollissima</i>		NC_014674
Fagales	Juglandaceae	<i>Juglans regia</i>		NC_028617
Geraniales	Francoaceae	<i>Francoa sonchifolia</i>		NC_021101
Geraniales	Francoaceae	<i>Melianthus villosus</i>		NC_023256
Geraniales	Francoaceae	<i>Viviania marifolia</i>		NC_023259
Geraniales	Geraniaceae	<i>Erodium rupestre</i>		NC_030719
Geraniales	Geraniaceae	<i>Geranium palmatum</i>		NC_014573
Geraniales	Geraniaceae	<i>Hypseocharis bilobata</i>		NC_023260
Geraniales	Geraniaceae	<i>Monsonia speciosa</i>		NC_014582
Geraniales	Geraniaceae	<i>Pelargonium alternans</i>		NC_023261
Malpighiales	Chrysobalanaceae	<i>Chrysobalanus icaco</i>		NC_024061
Malpighiales	Chrysobalanaceae	<i>Hirtella racemosa</i>		NC_024060
Malpighiales	Erythroxylaceae	<i>Erythroxylum novogranatense</i>		NC_030601
Malpighiales	Euphorbiaceae	<i>Ricinus communis</i>		NC_016736
Malpighiales	Malpighiaceae	<i>Galphimia angustifolia</i>	TEX-DJPG803	MK726010
Malpighiales	Salicaceae	<i>Salix babylonica</i>		NC_028350
Malvales	Malvaceae	<i>Gossypium turneri</i>		NC_026835
Malvales	Malvaceae	<i>Hibiscus syriacus</i>		NC_026909
Malvales	Malvaceae	<i>Theobroma cacao</i>		NC_014676
Malvales	Thymelaeaceae	<i>Aquilaria sinensis</i>		NC_029243
Myrtales	Alzateaceae	<i>Alzatea verticillata</i>	K-TNV548	MK726006
Myrtales	Combretaceae	<i>Laguncularia racemosa</i>	CONN00225898	MK726017
Myrtales	Combretaceae	<i>Terminalia guyanensis</i>	UEC-GHS1070	MK726027
Myrtales	Lythraceae	<i>Heimia apetala</i>	CONN00225896	MK726012
Myrtales	Lythraceae	<i>Lagerstroemia guilinensis</i>		NC_029885
Myrtales	Lythraceae	<i>Lagerstroemia fauriei</i>		NC_029808
Myrtales	Lythraceae	<i>Lagerstroemia indica</i>		NC_030484
Myrtales	Melastomataceae	<i>Blakea schlimii</i>		NC_031877
Myrtales	Melastomataceae	<i>Henriettea barkeri</i>		NC_031880
Myrtales	Melastomataceae	<i>Memecylon pauciflorum</i>	K-TNV679	MK726029
Myrtales	Melastomataceae	<i>Miconia dodecandra</i>		NC_031882
Myrtales	Melastomataceae	<i>Rhexia virginica</i>		NC_031886
Myrtales	Melastomataceae	<i>Tibouchina urvilleana</i>	CONN00225897	MK726030
Myrtales	Myrtaceae	<i>Allosyncarpia ternata</i>		NC_022413
Myrtales	Myrtaceae	<i>Corymbia eximia</i>		NC_022409
Myrtales	Myrtaceae	<i>Eucalyptus globulus</i>		NC_008115
Myrtales	Myrtaceae	<i>Eugenia uniflora</i>		NC_027744
Myrtales	Myrtaceae	<i>Heteropyxis natalensis</i>	K-MFF s.n.	MK726014
Myrtales	Myrtaceae	<i>Psidium guajava</i>		NC_033355
Myrtales	Myrtaceae	<i>Stockwellia quadrifida</i>		NC_022414
Myrtales	Myrtaceae	<i>Xanthostemon chrysanthus</i>	K-TNV684	MK726024
Myrtales	Onagraceae	<i>Ludwigia octovalvis</i>		NC_031385
Myrtales	Onagraceae	<i>Oenothera argillicola</i>		NC_010358

(continued on next page)

Table 1 (continued)

Order	Family	Species	Voucher ID	GenBank Accession Numbers
Myrtales	Onagraceae	<i>Oenothera grandiflora</i>		NC_029211
Myrtales	Onagraceae	<i>Oenothera oakesiana</i>		NC_029212
Myrtales	Onagraceae	<i>Oenothera villaricae</i>		NC_030532
Myrtales	Penaeaceae	<i>Saltera sarcocolla</i>	CONN00225892	MK726025
Myrtales	Vochysiaceae	<i>Callisthene erythroclada</i>	UEC-DG439	MK726008
Myrtales	Vochysiaceae	<i>Erisma bracteosum</i>	UEC-GHS937	MK726009
Myrtales	Vochysiaceae	<i>Erismadelphus exsul</i>	BRLU-SB751	MK726007
Myrtales	Vochysiaceae	<i>Korupodendron songweanum</i>	BRLU-SB3491	MK726013
Myrtales	Vochysiaceae	<i>Qualea grandiflora</i>	UEC-DG382	MK726022
Myrtales	Vochysiaceae	<i>Ruizterania albiflora</i>	UEC-TM258	MK726023
Myrtales	Vochysiaceae	<i>Salvertia convallariodora</i>	TEX-DJPG569	MK726026
Myrtales	Vochysiaceae	<i>Vochysia acuminata</i>	TEX-DJPG150	MK726031
Oxalidales	Oxalidaceae	<i>Averrhoa carambola</i>		NC_033350
Oxalidales	Oxalidaceae	<i>Oxalis drummondii</i>	TEX-DJPG722	MK726021
Rosales	Cannabaceae	<i>Cannabis sativa</i>		NC_026562
Rosales	Elaeagnaceae	<i>Elaeagnus macrophylla</i>		NC_028066
Rosales	Moraceae	<i>Ficus racemosa</i>		NC_028185
Rosales	Moraceae	<i>Morus indica</i>		NC_008359
Sapindales	Sapindaceae	<i>Acer morrisonense</i>		NC_029371
Sapindales	Anacardiaceae	<i>Anacardium occidentale</i>		NC_035235.1
Sapindales	Anacardiaceae	<i>Mangifera indica</i>		NC_035239.1
Sapindales	Burseraceae	<i>Boswellia sacra</i>		NC_029420
Sapindales	Meliaceae	<i>Azadirachta indica</i>		NC_023792
Sapindales	Citrusaceae	<i>Citrus aurantiifolia</i>		NC_024929
Sapindales	Rutaceae	<i>Zanthoxylum piperitum</i>		NC_027939
Sapindales	Sapindaceae	<i>Dimocarpus longan</i>	UCONN-201400014	MK726005
Sapindales	Sapindaceae	<i>Dipteronia sinensis</i>		NC_029338
Sapindales	Sapindaceae	<i>Litchi chinensis</i>		NC_035238.1
Sapindales	Sapindaceae	<i>Sapindus mukorossi</i>		NC_025554
Vitales	Vitaceae	<i>Tetrastigma hemsleyanum</i>		NC_029339
Vitales	Vitaceae	<i>Vitis vinifera</i>		NC_007957
Zygophyllales	Krameriaceae	<i>Krameria bicolor</i>	TEX-JN14-09-03-1	MK726015
Zygophyllales	Krameriaceae	<i>Krameria lanceolata</i>	TEX-JN6-IV-2015-1	MK726016
Zygophyllales	Zygophyllaceae	<i>Guaicum angustifolium</i>	TEX-BBS20-IV-2015-1	MK726011
Zygophyllales	Zygophyllaceae	<i>Larrea tridentata</i>	TEX-JN14-08-24-1	MK726018
Outgroup				
Caryophyllales	Caryophyllaceae	<i>Silene capitata</i>		NC_035226
Apiales	Araliaceae	<i>Aralia undulata</i>		NC_022810
Solanales	Solanaceae	<i>Nicotiana tabacum</i>		NC_001879.2

Facility (GSAF) at The University of Texas at Austin. Two species were kindly provided by The Royal Botanic Gardens, Kew, DNA bank (<https://www.kew.org/data/dnaBank/>).

Once the reads were available, the genome assembly pipeline was used to remove adaptors and PHIX, for quality trimming, and for genome assembly.

2.2. Phylogenetic inference

After gathering sequences of plastid protein-coding genes, the alignments and phylogenetic inference were performed. The code used to prepare the alignments using MAFFT [4] and MACSE [5] as well as the scripts used to infer phylogenies using Maximum Likelihood (ML), IQ-TREE [6], and Multispecies Coalescent (MSC) methods, SVDquartets [7], and ASTRAL-II [8] is presented in phylogenetic analysis pipeline.

2.3. Calculating distances of tree topologies and phylogenetic signal

Commented scripts present how the inferred phylogenies were further explored. First, Robinson-Foulds and Kendall-Colijn algorithms implemented in the R package TREESPACE [9] were used to

Table 2

Summary of output for Illumina sequencing of the 23 complete and the 4 draft plastomes.

Species	Total # Reads	Plastid reads	Average fold coverage ^a
<i>Laguncularia racemosa</i>	26,332,422	992,869	758.64
<i>Terminalia guyanensis</i>	28,526,166	335,902	41.137
<i>Heimia apetala</i>	27,057,568	3480,902	3630.02
<i>Memecylon pauciflorum</i>	29,404,628	1,776,470	2078.67
<i>Tibouchina urvilleana</i>	51,945,408	1,565,235	1321.25
<i>Heteropyxis natalensis</i>	23,511,880	437,572	620.57
<i>Saltera sarcocolla</i>	42,509,058	1,975,497	1589.48
<i>Callisthene erythroclada</i>	13,523,069	1,475,969	1679.93
<i>Erism bracteosum</i>	28,714,784	740,005	872.25
<i>Qualea grandiflora</i>	14,026,674	315,666	361.56
<i>Korupodendron songweanum</i>	26,380,860	621,934	697.51
<i>Ruizterania albiflora</i>	13,416,345	346,923	395.21
<i>Salvertia convallariodora</i>	13,719,886	1,131,571	1262.57
<i>Vochysia acuminata</i>	14,190,315	1,404,610	1588.55
<i>Azima tetracantha</i>	19,768,664	588,345	676.25
<i>Moringa oleifera</i>	35,924,836	3,914,845	3054.87
<i>Dimorcarpus longan</i>	39,914,336	1,125,881	870.47
<i>Galphimia angustifolia</i>	20,773,804	1,235,549	1449.65
<i>Oxalis drummondii</i>	30,898,958	1,231,010	1698.62
<i>Krameria bicolor</i>	20,882,876	1,032,768	1236.57
<i>Krameria lanceolata</i>	20,840,984	191,073	206.57
<i>Guaiacum angustifolium</i>	27,224,600	131,532	143.81
<i>Larrea tridentata</i>	23,104,094	2,094,664	2797.51
<i>Alzatea verticillata</i> (reads mapped to <i>S. sarcocolla</i>)	11,665,088	51,656	61.91
<i>Xanthostemon chrysanthus</i> (reads mapped to <i>H. natalensis</i>)	9,809,306	51,304	58.28
<i>Polygala alba</i> (reads mapped to <i>G. angustifolia</i>)	30,134,340	125,019	101.66
<i>Erismadelphus exsul</i> (reads mapped to <i>K. songweanum</i>)	13,270,682	286,659	368.03

^a Coverage calculation was performed in bbmap (BBTools package, available at <https://jgi.doe.gov/data-and-tools/bbtools/>) with the option "covstats". Reference index was built with the 23 complete plastomes with each representing a scaffold. The average fold coverage was calculated by mapping reads to the reference index and values were taken from the scaffold correspondent to each species. For species with incomplete plastomes (marked in bold) we used the closely related species as the scaffold indicated within the parentheses.

visualize the distances of species trees and between species trees and gene trees inferred. The code used is available at tree space. Lastly, five taxa from different taxonomic levels that had alternative placements were selected for a set of measurements of gene-wise and site-wise log-likelihood support of alternative topologies phylogenetic signal.

Acknowledgements

We thank the Texas Advanced Computing Center for providing access to its supercomputer Lone-star5. This work was supported by the C. L. Lundell Professorship; the Texas Ecolab Program (DJPG); the National Science Foundation for the Doctoral Dissertation Improvement Grant (DJPG, grant number 1601522); the Plant Biology Graduate Program; the Department of Integrative Biology at UT Austin through the Linda Escobar, the Jean Andrews, and the Lorraine Stengl fellowships; the Science without Borders (DJPG, Capes/Laspau-1186-13-2); and the UT Graduate Continuing Fellowship (DJPG).

Conflict of interest

The authors declare that they have no known competing financial interests or personal relationships that could have appeared to influence the work reported in this paper.

References

- [1] D.J.P. Gonçalves, B.B. Simpson, E.M. Ortiz, G.H. Shimizu, R.K. Jansen, Incongruence between species tree and gene trees and phylogenetic signal variation in plastid genes, *Mol. Phylogenetics Evol.* 138 (2019) 219–232. <https://doi.org/10.1016/j.ympev.2019.05.022>.
- [2] J.J. Doyle, J.L. Doyle, DNA isolation from small amounts of plant tissue, *Phytochem. Bull.* 19 (1987) 11–15.
- [3] M.L. Weng, J.C. Blazier, M. Govindu, R.K. Jansen, Reconstruction of the ancestral plastid genome in Geraniaceae reveals a correlation between genome rearrangements, repeats, and nucleotide substitution rates, *Mol. Biol. Evol.* 31 (2013) 645–659. <https://doi.org/10.1093/molbev/mst257>.
- [4] K. Katoh, J. Rozewicki, K.D. Yamada, MAFFT online service: multiple sequence alignment, interactive sequence choice and visualization, *Briefings Bioinf.* (2017) 1–7. <https://doi.org/10.1093/bib/bbx108>.
- [5] V. Ranwez, S. Harispe, F. Delsuc, E.J.P. Douzery, MACSE: Multiple alignment of coding sequences accounting for frameshifts and stop codons, *PLoS One* 6 (2011) 1–10. <https://doi.org/10.1371/journal.pone.0022594.s001>.
- [6] L.T. Nguyen, H.A. Schmidt, A. von Haeseler, B.Q. Minh, IQ-TREE: a fast and effective stochastic algorithm for estimating maximum-likelihood phylogenies, *Mol. Biol. Evol.* 32 (2015) 268–274. <https://doi.org/10.1093/molbev/msu300>.
- [7] J. Chifman, L. Kubatko, Quartet inference from SNP data under the coalescent model, *Bioinformatics* 30 (2014) 3317–3324. <https://doi.org/10.1093/bioinformatics/btu530>.
- [8] S. Mirarab, T. Warnow, ASTRAL-II: coalescent-based species tree estimation with many hundreds of taxa and thousands of genes, *Bioinformatics* 31 (2015) i44–i52. <https://doi.org/10.1093/bioinformatics/btv234>.
- [9] T. Jombart, M. Kendall, J. Almagro-García, C. Colijn, Treespace: statistical exploration of landscapes of phylogenetic trees, *Mol. Ecol. Resour.* 17 (2017) 1385–1392. <https://doi.org/10.1111/1755-0998.12676>.