

UNIVERSIDADE FEDERAL DO ESPÍRITO SANTO
CENTRO DE CIÊNCIAS AGRÁRIAS E ENGENHARIAS
PROGRAMA DE PÓS-GRADUAÇÃO EM GENÉTICA E MELHORAMENTO

DRIELLI CANAL

The draft genome assembly of *Psidium guajava* L.

ALEGRE, ES

2018

DRIELLI CANAL

The draft genome assembly of *Psidium guajava* L.

Dissertation presented to the
UniversidadeFederal do Espírito Santo, as part
of therequirements of the Genetics and
BreedingGraduate Program, for the attainment
of the titleMagister Scientiae.

Advisor: DSc. Dr. Adésio Ferreira

ALEGRE, ES

2018

Modelo de ficha catalográfica fornecido pelo Sistema Integrado de Bibliotecas da Ufes
para ser confeccionada pelo autor

DRIELLI CANAL

Dissertação apresentada à Universidade Federal do Espírito Santo como requisito parcial para obtenção do Título de Mestre pelo Programa de Pós-Graduação em Genética e Melhoramento.

APROVADA: 26 de fevereiro de 2019

Comissão Examinadora:

Dr. Adésio Ferreira

Universidade Federal do Espírito
Orientador

Dra. Marcia Flores da Silva Ferreira

Universidade Federal do Espírito Santo
Coorientadora

Dr. Otavio J. B. Brustolini

Universidade Federal do Espírito Santo
Coorientador

Dra. Paola Carpinetti-Oliveira

Universidade Federal do Espírito Santo
Coorientadora

Dr. Pedro Henrique Dias Dos Santos

Universidade Estadual do Norte Fluminense
Membroexterno

Dr. Wellington Ronildo Clarindo

Universidade Federal do Espírito Santo
Membroexterno

dedicated to my beloved family

Acknowledgements

First, I would like to thank God for grace, health and mercy with me. Second, to my advisors Dr. Adésio Ferreira and Marcia Flores for their patience, help, encouragement and guidance during my dissertation work. I would also like to express my sincere gratitude to my friends Otávio Brustolin for his mentorship and support and to Miquéias Fernandes for his patience, teachings and advice. I also sincerely appreciate the contributions of Pedro, Paola and Well in this work. I'm very grateful to all my colleagues at the Biometry Lab and Laboratório de Genética e Melhoramento for their time, the jokes, respect, affection, confidence, always supporting and assisting in working hours, even on weekends and late at night. I would like to thank the Laboratório de Biotecnologia and Laboratório de Citogenética, for the reagents, time, space provided and support. The main sources of funding Conselho Nacional de Desenvolvimento Científico e Tecnológico (CNPQ), also, Fundação de Apoio à Pesquisa do Espírito Santo (FAPES), Coordenação de Aperfeiçoamento de Pessoal de Nível Superior (CAPES) and Universidade Federal do Espírito Santo (UFES) for the opportunities. Finally, and most importantly, my utmost gratitude goes to my family, none of this would have been possible without their unconditional support, trust, sacrifice and love.

List of Abbreviations

AAdenine

BLAST Basic Local Alignment Search Tool

bp Base pair

BUSCO Benchmarking Universal Single-Copy Orthologs

C Cytosine

CDS Coding sequence

CTAB Cetyl Trimethyl Ammonium Bromide protocol

DAPI 40 ,60 - diamidino-2-phenylindole

de novo starting from the beginning; anew

DBG *de Bruijn* graphs

DIRS *Dictyostelium* intermediate repeat sequence

DNA deoxyribonucleic acid

DTTdithiothreitol

FastQ format for short-read sequence data along with quality

FI fluorescence intensity

FMCFlow cytometry analysis

G Guanine

GB giga bytes

GC Guanine-Cytosine, used as shorthand for GC-content

GFF3 Genome Feature Format, version 3

GO Gene ontology

HSPs High-scoring segment pairs

kbpkilo base pairs

LARDs Large Retrotransposon Derivates

LINE Long interspersed nuclear element

LTR Long Terminal Repeat

N50 Size of contig or scaffold in an assembly such that 50% of the assembly is incontigs of that size or larger

NCBI National Center for Biotechnology Information

NGS Next-Generation Sequencing

MbpMegabase pairs ($\times 10^6$ bp)

MITE Miniature inverted-repeat transposable elements

MYA Million years ago

NGS Next generation of sequencing

OLC Overlap-Layout-Consensus

PE Paired-End

PI propidiumiodide

pg picograms

RNA ribonucleic acid

SINE Short interspersed nuclear element

SMRT Single molecule real time (sequencing)

SSR Simple sequence repeats

SRA Sequence Read Archive database

T Thymine

TR tandem repeat

TE Transposable element

TIR Terminal Inverted Repeats

TPS terpene synthase genes

WGS Whole genome sequencing

ZMWs zero-mode waveguide

List of Figures

Figure 1. Total nucleotide distributions within the assembly.....	14
Figure 2. Validation of draft genome using BUSCO	15
Figure 3. Flowcytometry methods for genome size estimation and GC content in <i>Psidium guajava</i> ‘Cortibel RM’	16
Figure 4. Main transposable elements (TE) orders in the <i>Psidium guajava</i> genome	19
Figure 5. Percentage of microsatellites found in <i>P. guajava</i> genome and in coding DNA sequences.....	17
Figure 6. Percentage of the most frequent motifs in each class of microsatellites (SSRs) found in all genome and in coding DNA sequences of <i>Psidium guajava</i>	17
Figure 7. Features of <i>Psidiumguajava</i> genome	18
Figure S1. Pipeline for TE identification (de novo).....	40
Figure S2. Pipeline for TE annotation.....	40

List of Tables

Table 1. Summary of the guava genome assembly.	14
Table 2. Gene content of the <i>P. guajava</i> genome and structural annotation.	25
Table S1. Summary of sequencing data used in the <i>de novo</i> assembly of the draft <i>Psidium guajava</i> genome.	37
Table S2. NCBI Reference sequence numbers for the mitochondrion genomes used in the preprocessing step.	37
Table S3. Main statistics of draft of sequenced plant genomes.	38
Table S4. Annotation of transposable elements.	39

Table of Contents

1. Introduction	11
2.0 Results and discussion	12
2.1 Genome sequencing and assembly	12
2.2 Evaluation of assembly results	14
2.3 Repeat content	16
2.4 Genome Structure	24
3.0 Conclusions and future directions	25
4.0 Methods	26
4.1 Plant Material and DNA sequencing	26
4.2 Genome assembly	26
4.3 Assembly quality	27
4.4 Repeat content	29
4.5 Gene prediction	30
5.0 Bibliography	30
6.0 Appendix	38

ABSTRACT

Psidium guajava is a major commercial fruit of Myrtaceae family, mainly cropped for fresh fruit consumption, valued for its taste and high C vitamin content. Also, is used for industrialized production and has medicinal applications. Is widely distributed in all tropical and subtropical regions. However, guava conventional breeding could be a difficult and slow process and itshampered by limited molecular resources.This dissertation focuses on establishing the genome sequence for the *Psidium guajava*, that comprises all inheritable traits of an organism and provides essential information required to determine phenotypic traits and accelerate breeding programns. In this dissertation, I describe the process of creating a *de novo* assembly of *Psidium guajava* that was generated using a hybrid assembly strategy with Illumina short reads and long reads produced by Pacbio platforms.The genome assembly comprised 2,881 scaffolds with a total length of 412.5 Mb and a scaffold N50 of 312 Kb. Genome completeness indicated that 95,15% of the expected eudicotyledons genes were present in the genome assembly and cytogenetic analysis confirmed the genome size of 465 Mband GC content estimation around 40%. To detect the genomic features, an automatic annotation was performed, which revealed a large proportion of repetitive elements (57%), which suggests that they have played important roles in the evolution of *P. guajava*.The richness and distribution of microsatellites was also evaluated, revealing 210,037 potential SSR s allowing the fast and cost-effective development of molecular markersfor genetic research.A structural sequence annotation resulted in the prediction of about 61,089 putative genes, that are possible related to protein coding regions. The next stage of this project included the annotation of these proteins.

Keywords: whole genome sequencing, repetitive elements, microsatellites, PacBio sequencing, Illumina, bioinformatics.