

Genome assembly of the pink Ipê (*Handroanthus impetiginosus*, Bignoniaceae), a highly-valued ecologically keystone Neotropical timber forest tree --Manuscript Draft--

Manuscript Number:	GIGA-D-17-00159R2	
Full Title:	Genome assembly of the pink Ipê (<i>Handroanthus impetiginosus</i> , Bignoniaceae), a highly-valued ecologically keystone Neotropical timber forest tree	
Article Type:	Data Note	
Funding Information:	CNPq (471366/2007-2)	Professor Rosane Garcia Collevatti
	CNPq (457406/2012-7)	Professor Rosane Garcia Collevatti
	CNPq (476709/2012-1)	Dr Evandro Novaes
	FAP-DF (193.000.570/2009)	Dr Dario Grattapaglia
Abstract:	<p>Background: <i>Handroanthus impetiginosus</i> (Mart. ex DC.) Mattos is a keystone Neotropical hardwood tree widely distributed in seasonally dry tropical forests of South and Mesoamerica. Regarded as the "new mahogany", it is the second most expensive timber, the most logged species in Brazil, and currently under significant illegal trading pressure. The plant produces large amounts of quinoids, specialized metabolites with documented antitumorous and antibiotic effects. The development of genomic resources is needed to better understand and conserve the diversity of the species, to empower forensic identification of the origin of timber and to identify genes for important metabolic compounds.</p> <p>Findings: The genome assembly covers 503.7Mb (N50=81,316 bp), 90.4% of the 557 Mbp genome, with 13,206 scaffolds. A repeat database with 1,508 sequences was developed allowing masking ~31% of the assembly. Depth of coverage indicated that consensus determination adequately removed haplotypes assembled separately due to the extensive heterozygosity of the species. Automatic gene prediction provided 31,688 structures and 35,479 mRNA transcripts, while external evidence supported a well-curated set of 28,603 high-confidence models (90% of total). Finally, we used the genomic sequence and the comprehensive gene content annotation to identify genes related to the production of specialized metabolites.</p> <p>Conclusions: This genome assembly is the first well-curated resource for a Neotropical forest tree and the first one for a member of the Bignoniaceae family, opening exceptional opportunities to empower molecular, phytochemical and breeding studies. This work should inspire the development of similar genomic resources for the largely neglected forest trees of the mega-diverse tropical biomes.</p>	
Corresponding Author:	Rosane Garcia Collevatti, PhD Universidade Federal de Goias Goiania, GO BRAZIL	
Corresponding Author Secondary Information:		
Corresponding Author's Institution:	Universidade Federal de Goias	
Corresponding Author's Secondary Institution:		
First Author:	orzenil.silva@embrapa.br B Silva-Junior, PhD	
First Author Secondary Information:		
Order of Authors:	orzenil.silva@embrapa.br B Silva-Junior, PhD	
	Dario Grattapaglia, PhD	

	Evandro Novaes, PhD
	Rosane Garcia Collevatti, PhD
Order of Authors Secondary Information:	
Response to Reviewers:	<p>Dr Hans Zauner Assistant Editor GigaScience</p> <p>We double checked the manuscript and perform minor corrections. We also moved Figure S1 (ipê photo) to the main-text (it is now Figure 1) and change figure and supplemental material numbers.</p> <p>Hope our changes are satisfactory.</p> <p>All the best Rosane Collevatti</p>
Additional Information:	
Question	Response
Are you submitting this manuscript to a special series or article collection?	No
<p>Experimental design and statistics</p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our Minimum Standards Reporting Checklist. Information essential to interpreting the data presented should be made available in the figure legends.</p> <p>Have you included all the information requested in your manuscript?</p>	Yes
<p>Resources</p> <p>A description of all resources used, including antibodies, cell lines, animals and software tools, with enough information to allow them to be uniquely identified, should be included in the Methods section. Authors are strongly encouraged to cite Research Resource Identifiers (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our Minimum Standards Reporting Checklist?</p>	Yes
<p>Availability of data and materials</p> <p>All datasets and code on which the conclusions of the paper rely must be</p>	Yes

either included in your submission or deposited in [publicly available repositories](#) (where available and ethically appropriate), referencing such data using a unique identifier in the references and in the “Availability of Data and Materials” section of your manuscript.

Have you have met the above requirement as detailed in our [Minimum Standards Reporting Checklist](#)?

1
2
3
4 1 **Genome assembly of the Pink Ipê (*Handroanthus impetiginosus*, *Bignoniaceae*), a highly-**
5
6 2 **valued ecologically keystone Neotropical timber forest tree**
7
8 3

9
10 4 Orzenil Bonfim da Silva-Junior^{1,2}, Dario Grattapaglia^{1,2}, Evandro Novaes³, Rosane G. Collevatti^{4*}
11
12 5

13 6 ¹*EMBRAPA Recursos Genéticos e Biotecnologia, EPqB, Brasília, DF. 70770-910. Brazil.*
14

15 7 ²*Programa de Ciências Genômicas e Biotecnologia – Universidade Católica de Brasília, SGAN 916*
16
17 8 *Modulo B, Brasilia, DF 70790-160. Brazil*

18 9 ³*Escola de Agronomia, Universidade Federal de Goiás, CP 131. Goiânia, GO. 74001-970. Brazil.*
19

20 10 ⁴*Laboratório de Genética & Biodiversidade, Instituto de Ciências Biológicas, Universidade Federal*
21
22 11 *de Goiás. Goiânia, GO. 74001-970. Brazil.*
23
24 12

25 13 ***Corresponding author:** Rosane Garcia Collevatti, Instituto de Ciências Biológicas, Universidade
26
27 14 Federal de Goiás, 74001-970, Goiânia, GO, Brasil.
28
29 15

30
31 16 E-mail: rosanegc68@hotmail.com. Phone: +55 62 3521-1729.
32
33 17
34
35 18
36
37 19
38
39 20
40
41 21
42
43 22
44
45 23
46
47 24
48
49 25
50
51 26
52
53 27
54
55 28
56
57 29
58
59 30
60
61 31
62
63 32
64
65

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 34 **Abstract**

5
6 35

7
8 36 **Background:** *Handroanthus impetiginosus* (Mart. ex DC.) Mattos is a keystone Neotropical
9
10 37 hardwood tree widely distributed in seasonally dry tropical forests of South and Mesoamerica.
11 38 Regarded as the "new mahogany", it is the second most expensive timber, the most logged
12 39 species in Brazil, and currently under significant illegal trading pressure. The plant produces
13 40 large amounts of quinoids, specialized metabolites with documented antitumorous and
14 41 antibiotic effects. The development of genomic resources is needed to better understand and
15 42 conserve the diversity of the species, to empower forensic identification of the origin of timber
16 43 and to identify genes for important metabolic compounds.
17
18
19
20
21

22 44

23 45 **Findings:** The genome assembly covers 503.7Mb (N50=81,316 bp), 90.4% of the 557 Mbp
24 46 genome, with 13,206 scaffolds. A repeat database with 1,508 sequences was developed
25 47 allowing masking ~31% of the assembly. Depth of coverage indicated that consensus
26 48 determination adequately removed haplotypes assembled separately due to the extensive
27 49 heterozygosity of the species. Automatic gene prediction provided 31,688 structures and 35,479
28 50 mRNA transcripts, while external evidence supported a well-curated set of 28,603 high-
29 51 confidence models (90% of total). Finally, we used the genomic sequence and the
30 52 comprehensive gene content annotation to identify genes related to the production of
31 53 specialized metabolites.
32
33
34
35
36
37
38
39
40
41

42 54

43 55 **Conclusions:** This genome assembly is the first well-curated resource for a Neotropical forest
44 56 tree and the first one for a member of the *Bignoniaceae* family, opening exceptional
45 57 opportunities to empower molecular, phytochemical and breeding studies. This work should
46 58 inspire the development of similar genomic resources for the largely neglected forest trees of
47 59 the mega-diverse tropical biomes.
48
49

50 60

51 61

52 62 **Keywords:** heterozygous genome, RNA-seq, transposable elements, quinoids, *Bignoniaceae*
53
54
55
56
57
58
59
60
61
62
63
64
65

63 63

64 64

65 65

66 DATA DESCRIPTION

67

68 **Context.** The generation of plant genome assemblies is a key driver to develop powerful
69 genomic resources that allow gaining detailed insights into the evolutionary history of species
70 while enabling breeding and conservation efforts [1, 2]. Such advances took place first in model
71 plant species [3] followed by the mainstream [4] and minor crops [5], and some major forest
72 trees [6-9]. Genome sequences have also driven important advances in the description and
73 understanding of essential plant metabolic processes that underlie survival across distinct
74 lineages. Research on the functional roles of specialized metabolites, many of them
75 phylogenetically restricted [10], has recently addressed the gap in the species-specific
76 knowledge of specialized plant metabolism by sequencing the genome of key medicinal plants
77 [11, 12]. Innovation in this field has relied on a combination of high-throughput genomics,
78 including massive parallel sequencing and arrays with animal and clinical studies to elucidate the
79 mechanisms of target compounds as adjuvant therapies, to demonstrate the necessary
80 formulations for its biological effects and to determine which substances are beneficial or toxic.
81 Apart from recent reports of shallow transcriptome characterization using 454 pyrosequencing
82 [13] and a low-coverage (11X) fragmented genome assembly [14], essentially no well-curated
83 genome assembly and gene content annotation exist for Neotropical forest trees, despite their
84 recognized value by indigenous communities for the healing properties of their special
85 metabolites, increasingly exploited by large pharmaceutical corporations [15, 16]. An example of
86 such tree is the species *Handroanthus impetiginosus* (Mart. ex DC.) Mattos (syn. *Tabebuia*
87 *impetiginosa*, Bignoniaceae), popularly known as Pink Ipê, Lapacho or Pau d'arco, a source of
88 both high value timber and traditional medicine.

89

90 Species of *Handroanthus* and *Tabebuia* have virtually no genomic tools and resources, beyond a
91 handful of 21 microsatellites [17] with their known caveats for more sophisticated genetic
92 analyses in the areas of population genomics and evolution [18]. Whole-genome sequencing has
93 now become accessible to a point that efforts to develop improved genomic resources for such
94 species are possible and warranted. We built a preliminary assembly of the nuclear genome of a
95 single individual of *Handroanthus impetiginosus* based on short-reads and longer mate-pair DNA
96 sequence data to provide the necessary framework for the development of genomic resources
97 to support multiple genomic and genetic analyses of this keystone Neotropical hardwood tree

1
2
3
4 98 regarded as the "new mahogany". It is the second most expensive timber and the most logged
5
6 99 species in Brazil [19], exported largely to North America for residential decking and currently
7
8 100 under significant illegal trading pressure. Additionally, the tree produces large amounts of
9
10 101 natural products such as those of quinoid systems (1,4-anthraquinones, 1,4-naphthoquinones,
11
12 102 and 1,2-furanonaphthoquinones), specialized metabolites with promising antitumor, anti-
13
14 103 inflammatory and antibiotic effects [20, 21]. The high pressure of logging and illegal trading on
15
16 104 this species with a notable ecological keystone status urges conservation efforts of existing
17
18 105 populations.
19

20 107 **METHODS**

21
22 108
23
24 109 **Sample collection and sequencing.** DNA of a single adult tree of *H. impetiginosus* (UFG-1)
25
26 110 (Figure 1) was extracted using Qiagen DNeasy Plant Mini kit (Qiagen, DK). Flow cytometry was
27
28 111 used to check the genome size of tree UFG-1 indicating a genome size of (557 ±39) Mb /1C
29
30 112 (Figure S1) consistent with published estimates [22]. Total RNA from shoots of five seedlings and
31
32 113 from the differentiating xylem of the adult tree (UFG-1) was extracted using Qiagen RNeasy
33
34 114 Plant Mini kit (Qiagen, DK) and pooled for RNA sequencing. DNA and RNA sequencing was
35
36 115 performed at the High-Throughput Sequencing and Genotyping Center of the University of
37
38 116 Illinois Urbana-Champaign, USA. The following libraries were generated for sequencing: (1) two
39
40 117 shotgun genomic libraries of short fragments (300bp and 600bp) from tree UFG-1 (2) one
41
42 118 shotgun library from combined pools of five RNA samples tagged with a single index sequence.
43
44 119 Paired-end sequencing, 2x150 nt, was performed in two lanes of an Illumina HiSeq 2500
45
46 120 instrument (Illumina, CA, USA). Three additional mate-pair libraries (fragment lengths of 4kb to
47
48 121 5.5kb, 8kb to 10kb and 15kb to 20kb) for UFG-1 were also sequenced in two lanes of an Illumina
49
50 122 HiSeq 2000 instrument (2x101 bp). This long-range sequence resource was used to generate the
51
52 123 final genome assembly for annotation. A complete overview of the genome assembly and
53
54 124 annotation pipeline is provided (Figure S2).
55

56 125
57 126 *[Insert Figure 1 here]*
58

59 127
60 128 **Genome assembly using short paired-end and mate pair sequencing data.** Short reads and
61
62 129 mate-pair reads were stripped of sequencing adapters using *Fastq-mcf* [23]. Reads that mapped
63
64
65

1
2
3
4 130 to a database containing mitochondrial and chloroplast genomes of plants with *Bowtie1* [24]
5
6 131 (option `-v 3 -a -m 1`) were discarded. Mate-pair reads were inspected using a *Perl* script
7
8 132 (*TrimAdaptor.pl*), and sequences that did not contain the circularization adaptor were
9
10 133 discarded. By using the filtered short reads, Jellyfish2 (Jellyfish, RRID:SCR_005491) [25] and
11
12 134 GenomeScope [26] were applied to obtain estimates of the *H. impetiginosus* genome size,
13
14 135 repeat fraction and heterozygosity prior to the assembly. *ALLPATHS-LG* (*ALLPATHS-LG*,
15
16 136 *RRID:SCR_010742*) [27] was used for *de novo* assembly of the sequence data from both paired-
17
18 137 end and mate-pair data, with default options, in a stepwise strategy for error correction of
19
20 138 reads, handling of repetitive sequences and use of mate-pair libraries.

21 139

22 140 **Transposable elements and repetitive DNA.** Repetitive elements were detected and annotated
23
24 141 on the genome assembly with the RepeatModeler *de novo* repeat family identification and
25
26 142 modeling package (RepeatModeler, RRID:SCR_015027) [28]. Using RECON, RepeatScout and
27
28 143 Tandem Repeat Finder, repetitive sequences were detected in the scaffolds longer than 10 kb
29
30 144 using a combination of similarity-based and *de novo* approaches. The TE sequences were
31
32 145 evaluated using modeling capabilities of the RepeatModeler program, with default settings, to
33
34 146 compare the TE library against the entire assembled sequences and to refine and classify
35
36 147 consensus models of putative interspersed repeats. A complementary analysis intended to
37
38 148 augment the number of TE sequences classified according to current criteria [29] was performed
39
40 149 using the PASTEC program [30]. RepeatMasker Open-4.0 (RepeatMasker, RRID:SCR_012954)
41
42 150 [31] was used with the sequences from the *de novo* repetitive element library to annotate the
43
44 151 interspersed repeats and to detect simple sequence repeats (SSRs) on the genome assembly.

45 152

46 153 **Protein-coding genes annotation.** Protein-coding genes annotation was performed with a
47
48 154 pipeline that combines RNA-seq assembled transcript and protein alignments to the reference
49
50 155 with *de novo* predictions methods (Figure S2). RNA-Seq reads were screened for the presence of
51
52 156 adapters, which were removed using *Fastq-mcf* [23]. *Trimmomatic* (*Trimmomatic*,
53
54 157 *RRID:SCR_011848*) [32] was used to (1) remove low quality, no base called segments (N's) from
55
56 158 sequencing reads; (2) scan the read with a 4-base sliding window, cutting when the average
57
58 159 quality per base dropped below 15; and (3) remove reads shorter than 32 bp after trimming.
59
60 160 Trimmed reads mapped to mitochondrial, chloroplast and ribosomal sequences from plants with
61
62 161 *Bowtie1* [24] (options `-v 3 -a -m 1`) were also removed. Transcript *de novo* assemblies were
63
64
65

1
2
3
4 162 performed using *SOAP-Transdenovo* [33] and *Trinity de-novo* [34] from the processed reads. The
5
6 163 assemblies were concatenated and used as input to *EvidentialGene* [35], a comprehensive
7
8 164 transcriptome pipeline to identify likely complete coding regions and their proteins in the final,
9
10 165 combined, transcriptome assembly. Gene modeling was carried out using standard procedures
11
12 166 and tools described, for instance, in [36]. In summary, a genome-guided transcriptome assembly
13
14 167 of *H. impetiginosus* was performed with the JGI PERTRAN RNA-seq Read Assembler pipeline [37]
15
16 168 using both the RNA-Seq trimmed reads and sequences from the *de novo* transcript assembly.
17
18 169 Loci were identified by the assembled transcript alignments using BLASTX [38] and EXONERATE
19
20 170 [39] alignments of peptide sequences to the repeat-soft-masked genome using RepeatMasker
21
22 172 Known peptide sequences included manually curated data sets for plant species available from
23
24 173 UniProtKB/Swiss-Prot [41] and sequences available from Phytozome [1] version 11 for
25
26 174 *Arabidopsis thaliana*, *Oryza sativa*, *Erythranthe guttata*, *Solanum lycopersicum*, *Solanum*
27
28 175 *tuberosum*, *Populus trichocarpa* and *Vitis vinifera*. Gene structure were predicted by homology-
29
30 176 based predictors, FGENESH++, FGENESH_EST [42, 43] and GenomeScan [44]. Gene predictions
31
32 177 were improved by PASA (PASA, RRID:SCR_014656) [45], including adding UTRs, correcting
33
34 178 splicing and adding alternative transcripts. PASA-improved gene model peptides were subjected
35
36 179 to peptide homology analysis with the above-mentioned proteomes to obtain Cscore values and
37
38 180 peptide coverage. Cscore is the ratio of the peptide BLASTP score to the mutual best hit BLASTP
39
40 181 score, and peptide coverage is the highest percentage of peptide aligned to the best homolog. A
41
42 182 transcript was selected if its Cscore value was greater than or equal to 0.5 and its peptide
43
44 183 coverage was greater than or equal to 0.5 or if it had transcript coverage but the proportion of
45
46 184 its coding sequence overlapping repeats was less than 20%. For gene models where greater than
47
48 185 20% of the coding sequence overlapped with repeats, the Cscore value was required to be at
49
50 186 least 0.9 and homology coverage was required to be at least 70% to be selected. Selected gene
51
52 187 models were then subjected to classification analysis using *InterProScan 5* (*InterProScan*,
53
54 188 *RRID:SCR_005829*) [46] for PFAM domains, PANTHER, Enzyme Comissioned Number (EC) and
55
56 189 KEGG categories. Gene ontology annotation was obtained, where possible, from Interpro2GO
57
58 190 and EC2GO mappings.

191

192 **DATA VALIDATION AND QUALITY CONTROL**

193

1
2
3
4 194 **Global properties of the *H. impetiginosus* tree genome from the unassembled reads.**

5
6 195 Sequencing of the *H. impetiginosus* tree genome generated c. 599 million reads, comprising 73
7
8 196 Gbp of sequence data. This represents nearly 132× the expected sequence coverage. After
9
10 197 removal of adaptors, followed by standard error correction and trimming with ALLPATHS-LG,
11
12 198 with default options, c. 46 Gbp of data was found useful for the assembly process, yielding
13
14 199 sequencing coverage of 82x (63x from the fragments libraries and 19x from the mate pair
15
16 200 libraries). The estimated physical coverage was 400x based on the observed fragment size
17
18 201 distributions (Table S1). ALLPATHS-LG k-mer spectrum frequency analysis (at K=25) on useful
19
20 202 reads, error corrected reads, estimated a haploid genome size of 540,968,531 bp, a repeat
21
22 203 fraction of 38.0%, and a SNP rate of 1/88 bp (1.14%). An alternative analysis of the k-mer
23
24 204 frequencies using GenomeScope [26] produced a haploid genome size estimate of 503,748,072
25
26 205 bp, repetitive content of 36.6% and SNP rate of 1/60 bp (1.65%). Both estimates (Figure 2A) are
27
28 206 consistent with the flow cytometry estimates and in line with the expectations regarding the
29
30 207 heterozygous content of the *H. impetiginosus* genome, a predominantly outcrossed tree [47].
31
32 208 Sequencing errors caused an extreme peak at k = 1 in the k-mer frequency distribution. Both k-
33
34 209 mer histograms display two distinct peaks comprising the largest area of each histogram at
35
36 210 depths 27 and 55. The bimodal distributions characterize the expected behavior for k-mer
37
38 211 frequencies of a heterozygous diploid genome as seen, for example, in the recently reported
39
40 212 Oak genome [48]. In the right homozygous peak (at K=55), k-mers are shared between the two
41
42 213 homologous chromosomes. The left or heterozygous peak, with half the k-mer depth of the
43
44 214 homozygous peak, contains k-mers that are unique to each haplotype due to heterozygosity.
45
46 215 The difference in height between these peaks (heterozygous/homozygous ratio) is a measure of
47
48 216 the heterozygosity within the genome, which is 1.65% according to the GenomeScope modeling
49
50 217 equation.

51 218

52 219 *[Insert Figure 2 here]*

53 220

54 221 **Genome assembly.** State-of-the-art haploid genome assembler pipelines from short-reads
55
56 222 ALLPATHS-LG [27] and SOAPdenovo2 (SOAPdenovo2, RRID:SCR_014986) [49] were considered
57
58 223 for an initial evaluation on the dataset of reads. Two relatively new algorithms specifically
59
60 224 developed for de novo assembly of heterozygous genomes, MaSuRCA (MaSuRCA,
61
62 225 RRID:SCR_010691) [50] and PLATANUS (PLATANUS, RRID:SCR_015531) [51], were also
63
64
65

1
2
3
4 226 attempted as alternatives to the other two assemblers designed for genomes of low
5
6 227 heterozygosity. Reads were first preprocessed and error corrected using the algorithms
7
8 228 provided by each assembler. PLATANUS was set to run but after 10 weeks it did not produce any
9
10 229 result in an Intel(R) Xeon(R) server with 64 X7560 2.27GHz CPUs, 256 GB RAM, except for the k-
11
12 230 mer count table on the input trimmed reads. After 9 week-long runtimes in an Intel(R) Xeon(R)
13
14 231 server with 64 X7560 2.27GHz CPUs, 512 GB RAM, MaSuRCA successfully completed the
15
16 232 generation of the super-reads from the trimmed reads but the process was aborted on the
17
18 233 overlap-correction process in the Celera Assembler due to excessive CPU usage. SOAPdenovo2
19
20 234 ran very fast (3 days) but produced an assembly with total scaffold size of 860 Mbp. Analysis
21
22 235 with SOAPdenovo2 was run with different k-mer sizes, from 31 to 71, step of 10, but none of
23
24 236 them produced a reasonable assembly size in view of the expected size estimated by flow
25
26 237 cytometry and the k-mer frequency. ALLPATHS-LG was therefore used to assemble the genome
27
28 238 with default options. The short reads from fragmented libraries were error-corrected using
29
30 239 default settings (K-mer size of 24, ploidy of 2), fragment-filled and assembled into initial
31
32 240 unipaths (k-mer size of 96, ploidy of 2). Jumping reads from the mate-pair libraries were then
33
34 241 aligned to the unipaths and all alignments were processed in a seed-extension strategy with
35
36 242 junction point recognition within the read aimed to remove invalid and duplicate fragments to
37
38 243 perform error correction and initial scaffolding. This initial process produced an assembly graph
39
40 244 that was turned into scaffolds by analyzing branch points in the graph topology. This late
41
42 245 process converted single-base mismatches into ambiguous base codes at branch. It also
43
44 246 flattened some other structural features of the assembly including short indels. The contig
45
46 247 assembly comprised 109,064 sequences of length 500 bp or longer with total length of
47
48 248 466,314,780 bp. Genome assembly after scaffolding comprised 57,815 scaffolds of length 1 kbp
49
50 249 or longer with total length of 610,091,865 bp and N50 of 57 Kbp. The fraction of bases captured
51
52 250 in gaps was 23.9% and the rate of ambiguous bases for all bases captured in the assembly was
53
54 251 0.24%. This assembly was only slightly larger in size (<10%) than the empirically determined
55
56 252 genome size using flow cytometry [22].
57
58 253

59
60 254 **Alternative scaffold and gap-filling.** Although the ALLPATHS-LG performance was good in
61
62 255 recovering the expected genome size in the assembled contigs there was a high fraction of the
63
64 256 bases captured in gaps in the scaffolds ($\sim \frac{1}{4}$ of the total genome assembly). *De novo* assembly
65
66 257 algorithms applied to moderate-to-high levels of heterozygosity cannot match the performance

1
2
3
4 258 achieved in assemblies of homozygous genomes, especially at the contig assembly level [52]. We
5
6 259 thus used the assembled contigs to perform an alternative scaffolding step with SSPACE
7
8 260 (SSPACE, RRID:SCR_005056) [53] using the error-corrected short fragment reads and the
9
10 261 jumping reads. In this approach, genome assembly comprised 16,090 scaffolds of length 1 kbp
11
12 262 or longer with total length of 577,446,088 bp and N50 of 95 Kbp, respectively. The fraction of
13
14 263 bases captured in gaps dropped from 23.9% to 18.9% in contrast to ALLPATHS-LG scaffolding,
15
16 264 totaling 109,533,288 bp. The rate of ambiguous bases for all bases captured in the assembly
17
18 265 dropped from 0.24% to 0.13%. All preprocessed reads were reused in an attempt to close the
19
20 266 intra-scaffold gaps using the GapCloser (GapCloser, RRID:SCR_015026) [54] algorithm. Genome
21
22 267 assembly after gap-filling was 586,206,884 bp in 15,671 scaffolds of length 1 kbp or longer and
23
24 268 only 20,583,469 bp (3.51% of the genome assembly) remained in 24,907 gaps. N50 of scaffolds
25
26 269 of length 1 kbp or longer, with gaps, was 97,344 Kb (L50 = 1,792). Sequences longer than 20 kb
27
28 270 were assembled in only 6,791 scaffolds totaling 538,102,146 bp, ~97% of the genome size
29
30 271 estimated from flow cytometry (557 Mb).

31 272
32 273 ***Evaluation of accuracy of the genome assembly.*** A subset of fragments and jumping read pairs
33
34 274 (~15x sequencing coverage each) were used to uncover inaccuracies in the genome assembly.
35
36 275 Scaffolds with identified errors were broken or flagged for inspection. REAPR [55] was used to
37
38 276 test each base of the genome assembly looking for small local errors (such as a single base
39
40 277 substitutions, and short insertions or deletions) and structural errors (such as scaffolding errors)
41
42 278 located by means of changes to the expected distribution of inferred sequencing fragments
43
44 279 from the mapped reads using SMALT v0.7.6 [56]. REAPR reported that only 343,588,027 (~60%)
45
46 280 bases in the assembly should be free of errors, with 5,476 reported (1,658 within contigs, 3,818
47
48 281 over gaps) in the remaining 242,618,857 bp. The most frequent (~92%) type of inaccuracy
49
50 282 reported was *Perfect_cov* and *Link*. *Perfect_cov* means low coverage of perfect uniquely
51
52 283 mapping reads while *Link* describes situations in which reads map elsewhere in the assembly.
53
54 284 The recognition of this inaccuracy at the base pair level should thus reflect the repetitive nature
55
56 285 of the genome as inferred from the k-mer frequency spectra analysis (~36-38% of repeats).
57
58 286 Besides the base pair inaccurate calls due to repeats, other structural problems in the assembly
59
60 287 were identified based on sequence-coverage differences from the expected fragment size
61
62 288 distribution and the program used this information to break these. Given the high
63
64 289 heterozygosity and divergence between haplotypes on this diploid genome sequence,
65

1
2
3
4 290 homologous sequences can assemble separately or merge. Moreover, unresolved repeat
5
6 291 structures in the assembly might also contribute heavily to this issue. Structural errors in REAPR
7
8 292 were likely called at the boundaries of these regions. The final genome assembly after REAPR
9
10 293 breaks had 19,319 sequences of length 1 kbp or longer, with 576,829,188 bp. N50 size of
11
12 294 scaffolds dropped from 97,344 Kb (L50 = 1,792) to 71,491 bp (L50 = 2,379). The number of
13
14 295 remaining gaps in the assembly was 21,417 totaling 30,066,113 bp (5.05%).

15 296
16
17 297 Paired-end reads from the short fragment libraries were aligned back independently to this
18
19 298 genome assembly using SMALT (map -r 0 -x -y 0.5; default alignment penalty scores). Per-
20
21 299 scaffold depth of coverage was computed, regardless of mapping quality, using GATK
22
23 300 DepthofCoverage. The mean read depth across the scaffolds resulted in 66.45x. The mean read
24
25 301 length of the mapped reads was 139.8 bp and the corresponding k-mer coverage for size of 25
26
27 302 was 55.04x which matches with the homozygous peak computed from the k-mer frequency
28
29 303 distribution from the unassembled reads. The read depth frequencies are shown in Figure 2B.
30
31 304 The heterozygous/homozygous peak height (> 1) in the distribution suggests that the assembly
32
33 305 contains redundant copies of unmerged haplotypes due to the structural heterozygosity of the
34
35 306 diploid genome of the species. To specifically deal with the heterozygosity we introduced a step
36
37 307 to, leniently, recognize and remove alternative heterozygous sequences. Sequences of scaffolds
38
39 308 were aligned one versus all using BLAT (BLAT, RRID:SCR_011919) [57] and results were
40
41 309 concatenated in a single file of alignments and sorted. Similar sequences were identified on the
42
43 310 base of pairwise similarity using filterPSL utility from AUGUSTUS [58] with default parameters,
44
45 311 and retaining all best matches to each single sequence queried against all others that satisfy
46
47 312 minimal percentage of identity (minId=92%) and minimal percentage of coverage of the query
48
49 313 read (minCover=80%). We considered as heterozygous redundant those scaffolds that showed
50
51 314 pairwise similarity to exactly another sequence and their depth of coverage fell in a Poisson
52
53 315 distribution with parameters given by the heterozygous peak of the read depth distribution over
54
55 316 all scaffolds ($\lambda = 34$; Figure 2B). The final step was to keep only one copy – the largest one
56
57 317 – of the heterozygous scaffolds among pairs with high similarity.

58
59 318
60
61 319 **A preliminary assembly of the *H. impetiginosus* genome.** At the end of the accuracy evaluation
62
63 320 processes, the genome assembly had a total size of 503,308,897 bp, with gaps, in 13,206
64
65 321 scaffolds. The N50 of scaffolds of 1 kbp or longer was 80,946 bp (L50 = 1,906), the average size

1
2
3
4 322 of the sequences was 38,118 bp. Using 20 kbp as an approximate value of longest plant gene
5
6 323 length [59, 60], the percentage of scaffolds that equaled or surpassed this value in relation to
7
8 324 the empirically determined genome size is 83%, which corresponds to over 92% of the assembly
9
10 325 total size. Contigs generated by cutting scaffolds at each gap (of at least 25 base pair, i.e. 25 or
11
12 326 more Ns) produced N50 of 40,064 bp (L50 = 3,551) with average sequence size of 19,765 bp. The
13
14 327 remaining gaps comprised 26,447,057 bp (5.25% of the genome assembly) in 11,094 segments,
15
16 328 with size of $2,384 \pm 3,167$ bp. The total assembly size represents over 90% of the flow cytometry
17
18 329 genome estimate (557 Mb) and should provide a good start to build a further improved
19
20 330 reference genome assembly of the species using long-range scaffolding techniques such as
21
22 331 whole genome maps using either imaging methods [61] or contact maps of chromosomes based
23
24 332 on chromatin interactions [62]. Table 1 summarizes the main statistics of the *Handroanthus*
25
26 333 *impetiginosus* genome assembly with respect to the decisions made in the assembly process.

334

27 335 A reassessment of the assembly accuracy was carried out using REAPR on the final genome
28
29 336 assembly. A total of 121 errors within a contig were still recognized, a much smaller number
30
31 337 than previously annotated (1,658 errors). Figure 3A shows the frequency distribution for the
32
33 338 read depth computed from the paired-end read alignment to the scaffolds sequences. It
34
35 339 indicates the expected effect on the distribution in comparison to the previous more redundant
36
37 340 assembly. The height of the heterozygous peak was successfully lowered by removing unmerged
38
39 341 copies of the same heterozygous loci. Figure 3B shows the relation between the observed
40
41 342 number of scaffolds in the final assembly and their read coverage in comparison to a Poisson
42
43 343 approximation with $\lambda = 63$ which was the observed average sequencing coverage for
44
45 344 reads set from short fragment libraries. Loss of information due to repeat sequences is clearly a
46
47 345 limitation of this *H. impetiginosus* assembly. Given the high rate of non-classified consensus
48
49 346 sequences we can infer that most families/subfamilies of repeats might be underrepresented.

347

348 *[Insert Figure 3 here]*

349

50
51
52
53
54 350 To complement the depth of read coverage analyses, we performed additional analyses to
55
56 351 identify the most probable causes of breaks in the assembly. We inspected contig termini
57
58 352 defining the positions of the terminal nucleotides of each contig from the genome assembly
59
60 353 created by cutting at each gap (of at least one base pair, i.e. one or more Ns). This analysis was

1
2
3
4 354 developed using a protocol described elsewhere [63] and results are summarized in Figure 4.
5
6 355 Contig termini overlap most prominently (~50%) with regions that do not encompass any
7
8 356 annotated feature or regions that have no depth of coverage (~15%) based on mapped reads to
9
10 357 the assembly. It suggests that contigs end in large repeats not yet resolved given the inherent
11
12 358 limitations of short-read sequence data. Another possibility is that these regions can contain
13
14 359 low-copy young euchromatic segmental duplication with higher sequence similarity to the
15
16 360 consensus sequence. Annotated interspersed repeats (~18%) and short tandem repeats (~9%)
17
18 361 were the most prominently annotated features with overlap to contigs ends. Less than 8%
19
20 362 (2,473 of 31,668) of annotated gene models were found to overlap contigs ends, indicating that
21
22 363 very few are likely to be interrupted in this unfinished assembly. It is a trend that was confirmed
23
24 364 using BUSCO analysis which reported only 3% of fragmented genes. Based on variant
25
26 365 identification analysis with FreeBayes (FreeBayes, RRID:SCR_010761) using read data mapped to
27
28 366 the genome assembly, we found virtually no allelic variants located at contigs end, suggesting
29
30 367 that interruption of continuity and contiguity in the assembly is not related to differences
31
32 368 between haplotypes.

31 369

32 370 *[Insert Figure 4 here]*

33 371

34
35
36 372 **Repetitive DNA.** A total of 1,608 consensus sequences (average length = 773 bp and totaling
37
38 373 1,281,536 bp) representing interspersed repeats in the genome assembly were found. Search
39
40 374 for domains in these sequences with similarity to known large families of genes that could
41
42 375 confound the identification of true repeats indicated 85 false positives in the consensus library
43
44 376 of repeats. Further 50 sequences were annotated with predicted protein domains frequently
45
46 377 associated with protein coding genes. These 135 sequences were wiped out from the consensus
47
48 378 library. Most of the remaining 1,473 sequences (71.1%) could not find classification in the
49
50 379 hierarchical well-known classes of Transposable Elements [64] but 16.6% could be classified as
51
52 380 Class I (retrotransposons) including three orders: LTR (12.8%), LINE (1.6%) and SINE (2.2%); 8.4%
53
54 381 are Class II (DNA transposons). Other categories comprised non-autonomous TEs: TRIM (0.4%)
55
56 382 and MITE (3.5%). Unknown non-classified sequences in the consensus library cover a wide range
57
58 383 of sequence sizes from 42 bp up to 5,987 bp (average = 345 bp, median = 503 bp). The 1,473
59
60 384 sequences representing interspersed repeats in the consensus repeat library were used to mask
61
62 385 the genome with RepeatMasker. The masked fraction of the genome assembly comprised

1
2
3
4 386 155,348,349 bp, i.e. 30.9% of the total assembled genome of 503 Mbp. Remarkably, if we add to
5
6 387 these ~155 Mbp the 54 Mbp of non-captured base pairs in the assembly when considering the
7
8 388 empirically determined genome size (= 557–503), the repetitive fraction of the genome
9
10 389 approximates 37.5% (209 Mbp out 557 Mbp). This is within the expected range (36.6% - 38.0%)
11
12 390 for the repetitive fraction of the genome estimated from the reads set using k-mer profiling
13
14 391 approaches.

15 392
16 393 More than 50% of the masked bases in the assembly, or 80 Mbp, came from non-classified
17
18 394 sequences in the consensus library. In the well-known repeats, retrotransposons are the most
19
20 395 abundant class in the assembly comprising 50 Mbp (~1/3 of the masked bases) with prominence
21
22 396 of LTR/Gypsy (~23 Mb) and LTR/Copy (18 Mb) families of repeats. DNA transposons and non-
23
24 397 autonomous orders of transposons masked 12 Mbp and 11 Mbp (~1/6 of the masked bases),
25
26 398 respectively, highlighting the prominence of DNA/hAT families of class II and MITE (Figure 5).
27
28 399 Simple sequence repeats (SSRs) detection using RepeatMasker identified a total of 182,115
29
30 400 microsatellites with a density of 2.76 kb per SSR in the genome assembly. This density
31
32 401 corroborates the general finding that the overall frequency of microsatellites is inversely related
33
34 402 to genome size in plant genomes [65]. This SSRs density in *H. impetiginosus* (genome size of 557
35
36 403 Mbp/SSR density of 362 per Mbp) is higher than in larger plant genomes such as those of maize
37
38 404 (1,115 Mbp/163 SSRs per Mbp), *S. bicolor* (738 Mbp/175 per Mbp), *G. raimondii* (761 Mbp/74.8
39
40 405 per Mbp) [66] but lower than densities in smaller genomes such as those of *A. thaliana* (120
41
42 406 Mbp/ 418 per Mbp), *Medicago truncatula* (307 Mbp/ 495 per Mbp) and *C. sativus* (367 Mbp/
43
44 407 552 per Mbp) [67]. Different SSR motifs ranging from 1 to 6 bp showed that the di-nucleotide
45
46 408 repeats were the most abundant repeats followed by the mono- (Figure S3A). The frequency of
47
48 409 SSR decreased with increase in motif length (Supplementary Figure S3B), which is a trend usually
49
50 410 observed both in monocots and dicots [67].

51 411

52 412 *[Insert Figure 5 here]*

53 413

54 414 **Transcriptome assembly and gene content annotation and analysis.** A single run of Illumina
55
56 415 HiSeq 2500 sequencing, from a pool of RNA samples, generated nearly 148 million of paired end
57
58 416 reads. After adapter removal, trimming and coverage normalization, 55.2 million high-quality
59
60 417 reads (38%) were used to assemble the transcriptome using *de novo* (Trinity and SOAP-Trans-

1
2
3
4 418 denovo transcripts combined with the EvidentialGene pipeline) and genome guided methods
5
6 419 (PERTRAN). The PASA pipeline was used to integrate transcripts alignments to the genome
7
8 420 assembly from these set of sequences, generating 54,320 EST assemblies representing putative
9
10 421 protein-coding loci in the genome assembly. Loci were identified by the assembled transcript
11
12 422 alignments using BLASTX [36] and EXONERATE [37] alignments of plant peptides to the repeat-
13
14 423 soft-masked genome using RepeatMasker. After gene model prediction and refinements, a total
15
16 424 of 36,262 gene models were found in the genome assembly and 31,668 of them were retained
17
18 425 after quality assessment based on Cscore, protein coverage, and overlap to repeats as described
19
20 426 in Methods. The number of predicted mRNA transcripts was 35,479.

21
22 427
23
24 428 Structural features of the gene content are shown in Table 2 and Table 3. The average number
25
26 429 of exons per gene was ~5 and its average length was 285 bp. The average number of introns per
27
28 430 gene was ~4 and its average length was 445 bp. The GC content is significantly different
29
30 431 between exons and introns (t-test p-value < 0.0001). Coding sequences have ~43% of GC, while
31
32 432 introns have less with ~33% (Table 2). GC content tends to be higher in coding (exonic) than in
33
34 433 non-coding regions [68], which may be related to gene architecture and alternative splicing [69-
35
36 434 71]. A comparison of the gene features parameters, such as number and length (Figure S4A),
37
38 435 was carried out between *H. impetiginosus* and *Erythranthe guttata*, another plant in the order
39
40 436 Lamiales (Asterids), the model plant *A. thaliana* and the model tree *P. trichocarpa* (Rosids). As
41
42 437 depicted in the frequency histograms, the exons parameters are stable among these species
43
44 438 (Figure S4B). For the introns (Figure S4C), frequency histograms have a sharp peak around 90 bp
45
46 439 and a larger peak that is much lower in density. There is a small intron-size variability from
47
48 440 species to species in the distributions, especially for larger introns, which rarely go beyond than
49
50 441 10,000 bp. The intron length distributions in these four species is similar to those observed in
51
52 442 lineages that are late in the evolutionary time scale, such as plants and vertebrates [72]. The
53
54 443 sharp peak in the distributions at their “minimal intron” size is supposed to affect function by
55
56 444 enhancing the rate at which mRNA is exported from the cell nucleus [73, 74]. In the model plant
57
58 445 *A. thaliana*, a minimal intron group was previously defined [73] as anything that lies within three
59
60 446 standard deviations of the optimum peak at 89±12 bp (53 bp – 125 bp). According to this
61
62 447 definition, Table 3 summarizes the distribution of the minimal intron among genes of *H.*
63
64 448 *impetiginosus* and other selected plant species in the Asterids and Rosids lineages. We have
65
66 449 calculated the percentages of minimal introns out of the total introns and the fraction of

1
2
3
4 450 minimal-intron-containing genes with at least one minimal intron. Computed values were
5
6 451 similar between *H. impetiginosus* and those of selected species with higher number of large
7
8 452 introns (smaller minimal intron peak) but were more distinctive with those species such as *A.*
9
10 453 *thaliana* and *E. guttata* in which the number of large introns was lower (larger minimal intron
11
12 454 peak). This is thought as a general trend and was also observed in previous work [73]. These
13
14 455 comparative analyses about the structural properties of the predicted genes indicate that the
15
16 456 genome assembly of *H. impetiginosus* contains highly accurate gene structures.

17 457

18 458 *[Insert Figure 6 here]*

19
20 459

21
22 460 To further validate the gene content annotation, we used the transcript assemblies and selected
23
24 461 plant proteomes to inspect if these sequences could align in its entirety to the genomic
25
26 462 sequence. Out of the 31,668 primary mRNA transcripts (considering only the longest one when
27
28 463 isoforms were predicted) in the genome, 11,488 have 100% of their CDS covered by EST
29
30 464 assemblies. The remaining 20,054 transcripts have either a minimum of 80% of their CDS
31
32 465 covered by EST assemblies or a cscore ≥ 0.5 . From these latter, the encoded putative peptides
33
34 466 have excellent sequence similarity support from BLASTP comparisons with dicot species
35
36 467 *Erythranthe guttata* (5,224 genes), *Sesamum indicum* (4,625 genes), potato or tomato (2,777
37
38 468 genes), soybean (1,484 genes) and the poplar tree (1,424 genes) reflecting the taxonomic
39
40 469 relationship between *H. impetiginosus* and these other related dicots. Gene models support was
41
42 470 also found from more distantly related dicots (1,826 genes) and monocots (1,042 genes).
43
44 471 Altogether, 31,048 gene models (98%) show well-supported similarity hits to other known plant
45
46 472 protein sequences. Additional 517 predicted protein sequences did not produce hits and 103
47
48 473 sequences produced ambiguous hits from non-target species or represent possible
49
50 474 contaminants in the assembly such as endophytic fungi (ascomycetes, 42 sequences;
51
52 475 basidiomycetes, 17 sequences). Figure 6A summarizes the main finding regarding the similarity
53
54 476 analyses with known proteins.

55
56 477

57
58 478 BUSCO (BUSCO, RRID:SCR_015008) [75] single-copy genes plant profiles were used to estimate
59
60 479 completeness of the expected gene space as well as the duplicate fraction of the genome
61
62 480 assembly. Out of the 956 profiles searched on the assembly, 59 (6.1%) were reported missing
63
64 481 and 30 (3.1%) returned fragmented. From the profiles with complete match to the assembly,
65

1
2
3
4 482 867 (90.7%) were reported as single-copy and 247 (25.8%) were found completely duplicated.
5
6 483 We benchmarked our results by searching the BUSCO profiles on the genomes of other lamids,
7
8 484 *Erythranthe guttata* and *Olea europaea*. In *E. guttata* the analysis reported completeness level
9
10 485 of 88% (848 single-copy profiles with complete match) while fragmented genes were 52 (5.4%).
11
12 486 In *O. europaea*, the completeness level was 94% (905 complete single-copy profiles) and
13
14 487 fragmented genes were only 14 (1.4%). Summary of BUSCO analysis is presented in Figure 6B.

15 488
16
17 489 Databases for gene ontology (GO) annotation are rich resources to describe functional
18
19 490 properties of experimentally derived gene sets. To explore relationships between the GO terms
20
21 491 in the *H. impetiginosus* and related, well-curated, genomes we used WEGO [76] to perform a
22
23 492 genome-wide comparative analyses among broad functional GO terms with other lamids. The P-
24
25 493 value of Pearson Chi-Square test was considered to indicate significant relationships between
26
27 494 the proportions of genes of each GO term in these two datasets and to suggest patterns of
28
29 495 enrichment (Figures S5 and S6). These analyses revealed several GO terms in which the
30
31 496 proportion of genes in the two compared species were related. For the terms in which the
32
33 497 comparison did not indicate a significant relationship of gene proportions between the two
34
35 498 datasets, the compared GO terms suggested enrichments in *H. impetiginosus* for GO terms
36
37 499 involved in metabolic processes and catalytic activity in comparison to *E. guttata* and *O.*
38
39 500 *europaea*.

38 501

40 502 [Insert Figure 7 here]

41 503

43 504 The central role of enzymes as biological catalysts is a well-studied issue related to the chemistry
44
45 505 of cells [77]. An important feature of most enzymes is that their activities can be regulated to
46
47 506 function properly to comply with physiological needs of the organism. We observed that GO
48
49 507 term for enzyme regulatory activity encompass a higher proportion of genes in *H. impetiginosus*
50
51 508 than in the two other lamids, albeit the difference did not reach significance in *E. guttata*.
52
53 509 Research in *Arabidopsis*, an herbaceous plant, has found little connectivity between metabolites
54
55 510 and enzyme activity [78]. In comparison to *Arabidopsis* broader GO terms, *H. impetiginosus*
56
57 511 showed, as discussed above, enrichment for the proportion genes assigned to metabolic process
58
59 512 (49.1% > 47.4%; p-value 0.002) and catalytic activity (46.2% > 42.9%; p-value = 0). The
60
61 513 proportion of genes for enzyme regulatory activity was also higher in *H. impetiginosus* than *A.*

1
2
3
4 514 *thaliana*, though not statistically significant (p -value = 0.083). Investigations into whether and
5
6 515 how metabolic process and enzyme activities relate and how it could influence the known
7
8 516 richness of metabolites for forest trees of the mega diverse tropical biomes, particularly in the
9
10 517 genus *Tabebuia* and *Handroanthus*, shall be an interesting issue for future molecular and
11
12 518 chemistry studies.

13 519

14
15 520 **Benchmarking the genome assembly of *H. impetiginosus*.** Based on current standards for plant
16
17 521 genome sequence assembly [60, 79, 80] we have provided a quality assembly of high future
18
19 522 utility. To support functional analyses we classified the gene models into high-confidence and
20
21 523 low-confidence groups. Out of the 31,688 protein-coding loci annotated in the genome
22
23 524 assembly, 28,603 (90%) produced high-confidence gene models (Supplementary File S1). This
24
25 525 subset contains approximately the same number of genes reported in less fragmented genome
26
27 526 assemblies for other lamids. *E. guttata* ($2n=28$) reports 28,140 protein-coding genes [81]; *O.*
28
29 527 *europaea* ($2n = 46$) has 56,349 protein-coding genes [82] but its genome has likely undergone a
30
31 528 whole genome duplication event. Most of *Tabebuia* and *Handroanthus* species studied so far
32
33 529 have $2n = 40$ [22]. The fraction of gene duplicates in the BUSCO analysis (see Figure 5B) was
34
35 530 intended to estimate the level of redundancy in the genome assembly. We benchmarked our
36
37 531 results by searching the completed duplicated BUSCO profiles in the genomes of *E. guttata* and
38
39 532 *O. europaea*. In the first, we found them to be 15% (150 out 956), while in the latter the
40
41 533 duplicated profiles were 38% (364 out of 956). In these three lamids, it appears that the
42
43 534 frequency of small- and large-scale duplications, such as (paleo)polyploidy, can explain the
44
45 535 differences in the number of annotated genes and levels of gene duplication (*E. guttata* \leq *H.*
46
47 536 *impetiginosus* \ll *O. europaea*). It suggests that the *H. impetiginosus* genome has not undergone
48
49 537 a recent whole-genome duplication event, although a deeper analysis of this question, beyond
50
51 538 the scope of this study, remains open.

52 539

53
54 540 Our genome assembly metrics were benchmarked against comparable genome assemblies of
55
56 541 other highly heterozygous forest tree genomes (File S2 and Figure S7). The *H. impetiginosus*
57
58 542 assembly has 503 Mbp in 13,206 scaffolds ≥ 2 kbp, representing over 90% of the flow cytometry
59
60 543 estimated size (557 Mb). For *Quercus robur*, the assembly had 17,910 scaffolds ≥ 2 kbp with
61
62 544 scaffolds N50 of 260 kbp, but corresponding to 1.34 Gbp, i.e. 81% larger than the expected 740
63
64 545 Mbp genome, which is clearly an undesirable result [83]. For *Quercus lobata* with a genome size
65

1
2
3
4 546 of 730 Mbp two assemblies were provided: a haplotype-reduced assembly, with 40,158 contigs
5
6 547 totaling 760 Mb, N50 of 95 kbp and a more complete version for gene models, containing
7
8 548 94,394 scaffolds ≥ 2 kbp, totaling 1.15 Gbp, with an N50 of 278 kbp [48]. Despite our lower
9
10 549 NG50/N50 scaffold length < 100 kbp, the *H. impetiginosus* assembly has a large (60%)
11
12 550 percentage of scaffolds ≥ 20 kbp. This value is higher than the reported values for *Quercus lobata*
13
14 551 v0.5 (53%), *Quercus lobata v1.0* (51%) and *Quercus rubra* (48%), even if those assemblies had
15
16 552 higher NG50/N50 scaffold lengths. Finally, contigs termini analysis has found virtually no allelic
17
18 553 variants located at contigs ends, suggesting that interruption of continuity and contiguity in the
19
20 554 assembly is not related to differences between haplotypes. This genome assembly for
21
22 555 *Handroanthus impetiginosus* will thus be useful for variant calling, one of the main future
23
24 556 objectives for generating this resource.

25
26 557
27 558 **Genome-guided exploration of specialized metabolism genes of quinoid systems.** Aside from
28
29 559 its high valued wood, *H. impetiginosus* and other Ipê species are also known for their medicinal
30
31 560 effects. Extracts from its bark and wood have many ethnobotanical uses: against cancer,
32
33 561 malaria, fevers, trypanosomiasis, fungal and bacterial infections and stomach disorders [84, 85].
34
35 562 The wood extracts have also been demonstrated to have anti-inflammatory effects [86] [87].
36
37 563 The main bioactive components isolated from the Pink Ipê are Lapachol and its products [88],
38
39 564 which are naphthoquinones derived from the o-succinylbenzoate (OSB) pathway [89]. Lapachol
40
41 565 is also responsible for the well-known high resistance of the Ipê wood against rotting fungi and
42
43 566 insects [90]. In addition, naphthoquinones are aromatic substances with ecological importance
44
45 567 for the interaction of plants with other plants, insects and microbes [89]. Given their medicinal
46
47 568 and biological relevance, we have searched the *H. impetiginosus* annotated genes for the
48
49 569 enzymes involved in the biosynthesis of naphthoquinones. By searching for the KEGG identifiers
50
51 570 of these enzymes (e.g. K01851) in the InterPro annotation results, we found all the important
52
53 571 known enzymes that lead to the biosynthesis of lapachol (Figure 7). Unfortunately, however, the
54
55 572 last two steps of the lapachol biosynthesis pathway still constitute unidentified enzymes [89].
56
57 573 For comparative purposes, we downloaded the annotation file of five other species from the
58
59 574 Phytozome database. The number of *H. impetiginosus* genes encoding for the enzymes of each
60
61 575 step in the pathway is comparable to the numbers found in other species. However, three
62
63 576 exceptions were found. *H. impetiginosus* has five genes encoding the enzyme that converts
64
65 577 chorismate to isochorismate, the first step in the o-succinylbenzoate (OSB) pathway. Two other

1
2
3
4 578 steps found to have relatively more genes in *H. impetiginosus* are the ones that lead to the
5
6 579 synthesis of 1,4-Dihydroxy-2-naphtoyl-CoA and of 2-Phytyl-1,4-naphthoquinone. The availability
7
8 580 of sequences for these genes may open new avenues for biotechnological products and for a
9
10 581 better understanding of their ecological roles.

11 582

12
13 583 **RE-USE POTENTIAL**

14 584

15
16 585 We have reported a well-curated but still unfinished genome assembly for *Handroanthus*
17
18 586 *impetiginosus*, a highly valued, ecologically keystone tropical timber and a species rich in natural
19
20 587 products. The fragmentation of this preliminary assembly might be still be limiting for deeper
21
22 588 insights of whole-genome comparative analyses or studies of genome evolution [91], although
23
24 589 we think that such studies may be carried out using this assembly at least at the gene-level or
25
26 590 gene-family level. Nevertheless, the broad validation performed provides a useful genomic
27
28 591 resource for genetic and functional analysis including, but not limited to, downstream
29
30 592 applications such as variant calling, molecular markers development and functional studies.
31
32 593 Extensive documentation of quality throughout the assembly process was provided showing
33
34 594 that acceptable continuity was reached and that the fragmentation of the final sequence mostly
35
36 595 derives from loss of information on high-copy families of long interspersed repeats or the
37
38 596 presence of low-copy segmental duplications likely recently evolved with higher sequence
39
40 597 similarity to the consensus sequence. Certainly, there are still inaccuracies at the base and
41
42 598 assembly level but all efforts were made to deliver results to end user with the appropriate
43
44 599 documentation, making this initial read set, sequence and annotations as a primary and reliable
45
46 600 starting grounds for further improvement.

47 601

48
49 602 We have documented in detail the main features of the reported assembly. The total assembly
50
51 603 size of scaffolds with ≥ 2 kbp in length is 90% of the flow cytometry determined genome size, we
52
53 604 believe a remarkable accomplishment given the anticipated difficulties in assembling such a
54
55 605 repetitive and highly heterozygous diploid genome based exclusively on short-read sequencing.
56
57 606 The percentage of base pairs in scaffolds with ≥ 20 kbp is 83% (461 Mbp of 557 Mbp) of the
58
59 607 empirically determined genome size, which corresponds to 92% of the assembled total size (461
60
61 608 Mbp of 503 Mbp). Using 20 kbp as an approximate value of the longest plant gene length, this
62
63 609 result shows that 60% of the assembly is accessible for reliable gene annotation. Furthermore
64
65

1
2
3
4 610 the N50/NG50 (41 kbp/34 kbp) contig length is longer than 30 kbp, which has been suggested to
5
6 611 be an adequate minimum threshold for high utility of a genome assembly [79]. The percentage
7
8 612 of documented gaps in scaffolds is only 5.3% and the few misassembled signatures present in
9
10 613 the assembly were fully documented based on acceptable metrics such as fragment coverage
11
12 614 distribution error (FCD error). Less than 8% (2,473 of 31,668) of annotated gene models were
13
14 615 found to overlap contigs ends, indicating that very few are likely to be interrupted in this
15
16 616 unfinished assembly. No allelic variants were found at contigs ends, suggesting that interruption
17
18 617 of continuity and contiguity in the assembly is not related to differences between haplotypes,
19
20 618 therefore providing a valuable resource for variant calling and functional analysis. Over 86%
21
22 619 (27,380 of 31,668) of the gene models represented in the assembly have external evidential
23
24 620 support measured by Pasa-validated EST alignments from RNA-Seq or high-coverage alignments
25
26 621 with known plant proteins (>90% coverage). Furthermore, 80% (25,369 of 31,668) of transcripts
27
28 622 have conceptual translation that contain protein domain annotation, excluded those associated
29
30 623 to TEs. Finally, a summary of BUSCO analysis indicates that the detected number of plant single
31
32 624 copy orthologs represents 90% of the searched profiles (867 of 956) while only 6% is missing and
33
34 625 3% is fragmented.

35 626
36 627 This is the first well-curated genome for a Neotropical forest tree and the first one reported for
37
38 628 a member of the Bignoniaceae family. Besides expanding the opportunities for comparative
39
40 629 genomic studies by including an overlooked taxonomic family, the availability of this genome
41
42 630 assembly will foster functional studies with new targets and allow the development and
43
44 631 application of robust sets of genome-wide SNP genotyping tools to support multiple population
45
46 632 genomics analyses in *H. impetiginosus* and related species of the Tabebuia Alliance. This group
47
48 633 includes several of the most ecologically and economically important timber species of the
49
50 634 American tropics. Going beyond the species-specific significance of these results, this study
51
52 635 paves the way for developing similar genomic resources for other Neotropical forest trees of
53
54 636 equivalent relevance. This in turn will open exceptional prospects to empower a higher-level
55
56 637 understanding of the evolutionary history, species distribution and population demography of
57
58 638 the still largely neglected forest trees of the mega diverse tropical biomes. Furthermore, this
59
60 639 genome assembly provides a new resource for advances in the current integration between
61
62 640 genomics, transcriptomics and metabolomics approaches for exploration of the enormous
63
64 641 structural diversity and biological activities of plant-derived compounds.
65

642

643 AVAILABILITY OF SUPPORTING DATA

644

645 Sequences for the genome and assembly along with gene content annotation as well as the raw
646 sequencing reads have been deposited into GenBank, BioProject PRJNA324125. This Whole
647 Genome Shotgun (WGS) project has been deposited at DDBJ/ENA/GenBank under the accession
648 NKXS000000000. The version described in this paper is version NKXS01000000. BioSample for
649 WGS is SAMN05195323 and corresponding SRA run accessions are SRR3624821 - SRR3624825.
650 BioSample for RNA-Seq is SAMN07346903 with SRA run accession SRR5820886. Supporting data
651 and summary outputs for main analyses in this Data Note are available via the *GigaScience*
652 repository GigaDB [92]. The Perl script that automated the read set from mate-pair sequencing
653 preprocessing (TrimAdaptor.pl) was uploaded to GigaDB under permission of the original
654 authors at the High-Throughput Sequencing and Genotyping Center Unit of the University of
655 Illinois Urbana-Champaign.

656

657 List of abbreviations

658 BLASTP, Basic Local Alignment Search Tool for Proteins; BLAT, BLAST-like alignment tool; CDS,
659 coding DNA sequence; EC, Enzyme Commission Number; EST, Expressed Sequence Tag; GATK,
660 Genome Analysis Toolkit; GO, Gene Ontology; LINE, Long Interspersed Nuclear Elements; LTR,
661 Long Terminal Repeats; MBH, Mutual Best Hit; MITE, Miniature Inverted-Repeat Transposable
662 Elements; mRNA, messenger RNA; PASA, Program to Assemble Spliced Alignment; REAPR,
663 Recognition of Errors in Assemblies using Paired Reads; SINE, Short Interspersed Nuclear
664 Elements; SNP, Single Nucleotide Polymorphism; SSPACE, SSPACE-based Scaffolding of Pre-
665 Assembled Contigs after Extension; TE, transposable element.

666

667 Ethics approval

668 Not applicable

669

670 Consent for publication

671 Not applicable

672

673 Competing interests1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

674 The authors declare that they have no competing interests.

675

676 **Funding and acknowledgements**

677 This work was supported by competitive grants from CNPq to RGC (project no. 471366/2007-2
678 and Rede Cerrado CNPq/PPBio project no. 457406/2012-7), to EN (CNPq Proc. 476709/2012-1)
679 and to DG (PRONEX FAP-DF Project Grant "NEXTREE" 193.000.570/2009). RGC and DG have
680 been supported by productivity grants from CNPq, which we gratefully acknowledge. OBSJr has
681 been supported by an EMBRAPA doctoral fellowship and was an Affiliate Researcher at
682 Lawrence Berkeley National Laboratory (LBNL), Berkeley CA, at the time of this research. OBSJr
683 thanks to DM Goodstein and the members of the Phytozome team at the LBNL/Joint Genome
684 Institute (JGI) for their valuable help and support in working with the JGI pipelines for genomic
685 research. WE also thank Dr. Gabriela Ferreira Nogueira and André Luis X. de Souza for their help
686 with flow cytometry analysis.

687

688 **Authors' contributions**

689 OBSJr performed sequence data analysis and genome assembly and together with EN carried
690 out transcriptome and protein-coding gene annotation. RC and DG conceived the project,
691 collected samples, extracted genomic DNA and RNA, carried out flow cytometry analysis and
692 supervised the project. All authors were involved in discussions, writing and editing. All authors
693 read and approved the final manuscript.

694

695

REFERENCES

696

- 697 1. Goodstein DM, Shu SQ, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W,
698 Hellsten U, Putnam N *et al*: **Phytozome: a comparative platform for green plant**
699 **genomics**. *Nucleic Acids Research* 2012, **40**(D1):D1178-D1186.
- 700 2. Kang YJ, Lee T, Lee J, Shim S, Jeong H, Satyawani D, Kim MY, Lee SH: **Translational**
701 **genomics for plant breeding with the genome sequence explosion**. *Plant Biotechnology*
702 *Journal* 2016, **14**(4):1057-1069.
- 703 3. Bevan M, Walsh S: **The Arabidopsis genome: A foundation for plant research**. *Genome*
704 *Research* 2005, **15**(12):1632-1642.
- 705 4. Morrell PL, Buckler ES, Ross-Ibarra J: **Crop genomics: advances and applications**. *Nature*
706 *Reviews Genetics* 2012, **13**(2):85-96.
- 707 5. Varshney RK, Glaszmann JC, Leung H, Ribaut JM: **More genomic resources for less-**
708 **studied crops**. *Trends in Biotechnology* 2010, **28**(9):452-460.

1
2
3
4
5
6
7
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 709 6. Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J, Jenkins J,
5 710 Lindquist E, Tice H, Bauer D *et al*: **The genome of *Eucalyptus grandis***. *Nature* 2014,
6 711 **510**(7505):356-362.
- 7 712 7. Tuskan GA, DiFazio S, Jansson S, Bohlmann J, Grigoriev I, Hellsten U, Putnam N, Ralph S,
8 713 Rombauts S, Salamov A *et al*: **The genome of black cottonwood, *Populus trichocarpa***
9 714 **(Torr. & Gray)**. *Science* 2006, **313**(5793):1596-1604.
- 10 715 8. Neale DB, Wegrzyn JL, Stevens KA, Zimin AV, Puiu D, Crepeau MW, Cardeno C, Koriabine
11 716 M, Holtz-Morris AE, Liechty JD *et al*: **Decoding the massive genome of loblolly pine**
12 717 **using haploid DNA and novel assembly strategies**. *Genome Biology* 2014, **15**(3).
- 13 718 9. Nystedt B, Street NR, Wetterbom A, Zuccolo A, Lin YC, Scofield DG, Vezzi F, Delhomme
14 719 N, Giacomello S, Alexeyenko A *et al*: **The Norway spruce genome sequence and conifer**
15 720 **genome evolution**. *Nature* 2013, **497**(7451):579-584.
- 16 721 10. Moghe G, Last R: **Something old, something new: Conserved enzymes and the**
17 722 **evolution of novelty in plant specialized metabolism**. *Plant Physiology*
18 723 2015:pp.00994.02015.
- 19 724 11. Stone R: **Lifting the Veil on Traditional Chinese Medicine**. *Science* 2008, **319**(5864):709-
20 725 710.
- 21 726 12. Chappell J, DellaPenna D, O'Connor S: **Specific Aims for Medicinal Plant Genomics**
22 727 **Resource**. *Medicinal Plants Genomics Resource* 2017.
- 23 728 13. Brousseau L, Tinaut A, Duret C, Lang T, Garnier-Gere P, Scotti I: **High-throughput**
24 729 **transcriptome sequencing and preliminary functional analysis in four Neotropical tree**
25 730 **species**. *BMC Genomics* 2014, **15**(1):238.
- 26 731 14. Olsson S, Seoane-Zonjic P, Bautista Ro, Claros G, González-Martínez S, Scotti I, Scotti-
27 732 Saintagne C, Hardy O, Heuertz M: **Development of genomic tools in a widespread**
28 733 **tropical tree, *Symphonia globulifera* L.f.: a new low-coverage draft genome, SNP and**
29 734 **SSR markers**. *Molecular Ecology Resources* 2017, **17**(4):614-630.
- 30 735 15. Cadena-González A, Sorensen M, Theilade I: **Use and valuation of native and**
31 736 **introduced medicinal plant species in Campo Hermoso and Zetaquirá, Boyacá,**
32 737 **Colombia**. *Journal of Ethnobiology and Ethnomedicine* 2013, **9**(1):23.
- 33 738 16. Bodker G, Bhat KKS, Burley J, Vantomme P: **Medicinal plants for forest conservation**
34 739 **and health care**. *Food and Agriculture Organization of the United Nations* 1997.
- 35 740 17. Braga AC, Reis AMM, Leoi LT, Pereira RW, Collevatti RG: **Development and**
36 741 **characterization of microsatellite markers for the tropical tree species *Tabebuia aurea***
37 742 **(Bignoniaceae)**. *Molecular Ecology Notes* 2007, **7**(1):53-56.
- 38 743 18. Liu B, Shi Y, Yuan J, Hu X, Zhang H, Li N, Li Z, Chen Y, Mu D, Fan W: **Estimation of**
39 744 **genomic characteristics by analyzing k-mer frequency in de novo genome projects**.
40 745 *arXiv:13082012* 2013.
- 41 746 19. Schulze M, Grogan J, Uhl C, Lentini M, Vidal E: **Evaluating ipe (*Tabebuia*, Bignoniaceae)**
42 747 **logging in Amazonia: Sustainable management or catalyst for forest degradation?**
43 748 *Biological Conservation* 2008, **141**(8):2071-2085.
- 44 749 20. Inagaki R, Ninomiya M, Tanaka K, Watanabe K, Koketsu M: **Synthesis and Cytotoxicity**
45 750 **on Human Leukemia Cells of Furonaphthoquinones Isolated from *Tabebuia* Plants**.
46 751 *Chemical & Pharmaceutical Bulletin* 2013, **61**(6):670-673.
- 47 752 21. Park BS, Kim JR, Lee SE, Kim KS, Takeoka GR, Ahn YJ, Kim JH: **Selective growth-inhibiting**
48 753 **effects of compounds identified in *Tabebuia impetiginosa* inner bark on human**
49 754 **intestinal bacteria**. *Journal of Agricultural and Food Chemistry* 2005, **53**(4):1152-1157.

- 1
2
3
4 755 22. Collevatti RG, Dornelas MC: **Clues to the evolution of genome size and chromosome**
5 756 **number in *Tabebuia* alliance (Bignoniaceae)**. *Plant Systematics and Evolution* 2016,
6 757 **302(5):601-607**.
- 7 758 23. Aronesty E: **Comparison of sequencing utility programs**. *The Open Bioinformatics*
8 759 *Journal* 2013, **7:1-8**.
- 9 760 24. Langmead B, Trapnell C, Pop M, Salzberg SL: **Ultrafast and memory-efficient alignment**
10 761 **of short DNA sequences to the human genome**. *Genome Biology* 2009, **10(3)**.
- 11 762 25. Marçais G, Kingsford C: **A fast, lock-free approach for efficient parallel counting of**
12 763 **occurrences of k-mers**. *Bioinformatics* 2011, **27(6):764-770**.
- 13 764 26. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC:
14 765 **GenomeScope: fast reference-free genome profiling from short reads**. *Bioinformatics*
15 766 *btx153* 2017.
- 16 767 27. Gnerre S, MacCallum I, Przybylski D, Ribeiro F, Burton J, Walker B, Sharpe T, Hall G, Shea
17 768 T, Sykes S *et al*: **High-quality draft assemblies of mammalian genomes from massively**
18 769 **parallel sequence data**. *Proceedings of the National Academy of Sciences* 2011,
19 770 **108(4):1513-1518**.
- 20 771 28. Flutre T, Duprat E, Feuillet C, Quesneville H: **Considering Transposable Element**
21 772 **Diversification in De Novo Annotation Approaches**. *Plos One* 2011, **6(1)**.
- 22 773 29. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, Flavell A, Leroy P,
23 774 Morgante M, Panaud O *et al*: **A unified classification system for eukaryotic**
24 775 **transposable elements**. *Nature Reviews Genetics* 2007, **8(12):973-982**.
- 25 776 30. Hoede C, Arnoux S, Moisset M, Chaumier T, Inizan O, Jamilloux V, Quesneville H:
26 777 **PASTEC: An Automatic Transposable Element Classification Tool**. *Plos One* 2014, **9(5)**.
- 27 778 31. Smit AFA, Hubley R, Green P: **RepeatMasker Open-4.0 (2013-2015)**. 2015.
- 28 779 32. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina sequence**
29 780 **data**. *Bioinformatics* 2014, **30(15):2114-2120**.
- 30 781 33. Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, Huang W, He G, Gu S, Li S *et al*:
31 782 **SOAPdenovo-Trans: de novo transcriptome assembly with short RNA-Seq reads**.
32 783 *Bioinformatics* 2014, **30(12):1660-1666**.
- 33 784 34. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,
34 785 Raychowdhury R, Zeng QD *et al*: **Full-length transcriptome assembly from RNA-Seq**
35 786 **data without a reference genome**. *Nature Biotechnology* 2011, **29(7):644-U130**.
- 36 787 35. Gilbert D: **EvidentialGene: mRNA Transcript Assembly Software**. *EvidentialGene :*
37 788 *Evidence Directed Gene Construction for Eukaryotes* 2013.
- 38 789 36. Schmutz J, McClean P, Mamidi S, Wu A, Cannon S, Grimwood J, Jenkins J, Shu S, Song Q,
39 790 Chavarro C *et al*: **A reference genome for common bean and genome-wide analysis of**
40 791 **dual domestications**. *Nature Genetics* 2014, **46(7):707-713**.
- 41 792 37. Shu S, Goodstein DM, Hayes D, Mitros T, Rokhsar D: **JGI Plant Genomics Gene**
42 793 **Annotation Pipeline**. *SciTech Connect* 2017.
- 43 794 38. Gish W, States D: **Identification of protein coding regions by database similarity search**.
44 795 *Nature Genetics* 1993, **3(3):266-272**.
- 45 796 39. Slater GS, Birney E: **Automated generation of heuristics for biological sequence**
46 797 **comparison**. *Bmc Bioinformatics* 2005, **6**.
- 47 798 40. **RepeatMasker Open-4.0**. <http://www.repeatmasker.org>
48 799 [\[http://www.repeatmasker.org\]](http://www.repeatmasker.org)
- 49 800 41. UniProt Consortium: **UniProt: a hub for protein information**. *Nucleic Acids Research*
50 801 2014, **43(D1):D204-D212**.
- 51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 802 42. Salamov AA, Solovyev VV: **Ab initio gene finding in Drosophila genomic DNA.** *Genome Research* 2000, **10**(4):516-522.
5 803
6 804 43. Solovyev V, Kosarev P, Seledsov I, Vorobyev D: **Automatic annotation of eukaryotic**
7 805 **genes, pseudogenes and promoters.** *Genome Biology* 2006, **7**.
8 806 44. Yeh RF, Lim LP, Burge CB: **Computational inference of homologous gene structures in**
9 807 **the human genome.** *Genome Research* 2001, **11**(5):803-816.
10 808 45. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK, Hannick LI, Maiti R, Ronning CM,
11 809 Rusch DB, Town CD *et al*: **Improving the Arabidopsis genome annotation using**
12 810 **maximal transcript alignment assemblies.** *Nucleic Acids Research* 2003, **31**(19):5654-
13 811 5666.
14 812 46. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell
15 813 A, Nuka G *et al*: **InterProScan 5: genome-scale protein function classification.**
16 814 *Bioinformatics* 2014, **30**(9):1236-1240.
17 815 47. Braga AC, Collevatti RG: **Temporal variation in pollen dispersal and breeding structure**
18 816 **in a bee-pollinated Neotropical tree.** *Heredity* 2011, **106**(6):911-919.
19 817 48. Sork VL, Fitz-Gibbon ST, Puiu D, Crepeau M, Gugger PF, Sherman R, Stevens K, Langley
20 818 CH, Pellegrini M, Salzberg SL: **First Draft Assembly and Annotation of the Genome of a**
21 819 **California Endemic Oak Quercus lobata Nee (Fagaceae).** *G3-Genes Genomes Genetics*
22 820 2016, **6**(11):3485-3495.
23 821 49. Luo RB, Liu BH, Xie YL, Li ZY, Huang WH, Yuan JY, He GZ, Chen YX, Pan Q, Liu YJ *et al*:
24 822 **SOAPdenovo2: an empirically improved memory-efficient short-read de novo**
25 823 **assembler.** *Gigascience* 2012, **1**.
26 824 50. Zimin AV, Marcais G, Puiu D, Roberts M, Salzberg SL, Yorke JA: **The MaSuRCA genome**
27 825 **assembler.** *Bioinformatics* 2013, **29**(21):2669-2677.
28 826 51. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,
29 827 Nagayasu E, Maruyama H *et al*: **Efficient de novo assembly of highly heterozygous**
30 828 **genomes from whole-genome shotgun short reads.** *Genome Research* 2014,
31 829 **24**(8):1384-1395.
32 830 52. Malinsky M, Simpson JT, Durbin R: **Trio-sga: facilitating de novo assembly of highly**
33 831 **heterozygous genomes with parent-child trios.** *bioRxiv* 2016.
34 832 53. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W: **Scaffolding pre-assembled**
35 833 **contigs using SSPACE.** *Bioinformatics* 2011, **27**(4):578-579.
36 834 54. Nadalin F, Vezzi F, Policriti A: **GapFiller: a de novo assembly approach to fill the gap**
37 835 **within paired reads.** *Bmc Bioinformatics* 2012, **13**.
38 836 55. Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M, Otto TD: **REAPR: a universal**
39 837 **tool for genome assembly evaluation.** *Genome Biology* 2013, **14**(5):R47.
40 838 56. Ponstingl H, Ning ZM: **SMALT.** 2010 - 2015 *Genome Research Ltd* 2016.
41 839 57. Kent WJ: **BLAT--the BLAST-like alignment tool.** *Genome research* 2002, **12**(4):656-664.
42 840 58. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B: **AUGUSTUS: ab initio**
43 841 **prediction of alternative transcripts.** *Nucleic Acids Research* 2006, **34**:W435-W439.
44 842 59. Ramírez-Sánchez O, Pérez-Rodríguez P, Delaye L, Tiessen A: **Plant Proteins Are Smaller**
45 843 **Because They Are Encoded by Fewer Exons than Animal Proteins.** *Genomics,*
46 844 *Proteomics & Bioinformatics* 2016, **14**(6):357-370.
47 845 60. Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S, Chapman
48 846 JA, Chapuis G, Chikhi R *et al*: **Assemblathon 2: evaluating de novo methods of genome**
49 847 **assembly in three vertebrate species.** *GigaScience* 2013, **2**:10-10.
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 848 61. Lam ET, Hastie A, Lin C, Ehrlich D, Das SK, Austin MD, Deshpande P, Cao H, Nagarajan N,
5 849 Xiao M *et al*: **Genome mapping on nanochannel arrays for structural variation analysis**
6 850 **and sequence assembly**. *Nature Biotechnology* 2012, **30**(8):771-776.
7 851 62. Ay F, Noble WS: **Analysis methods for studying the 3D architecture of the genome**.
8 852 *Genome Biology* 2015, **16**.
9 853 63. Tørresen OK, Star B, Jentoft S, Reinart WB, Grove H, Miller JR, Walenz BP, Knight J,
10 854 Ekholm JM, Peluso P *et al*: **An improved genome assembly uncovers prolific tandem**
11 855 **repeats in Atlantic cod**. *BMC Genomics* 2017, **18**(1):95.
12 856 64. Wicker T, Sabot F, Hua-Van A, Bennetzen J, Capy P, Chalhoub B, Flavell A, Leroy P,
13 857 Morgante M, Panaud O *et al*: **A unified classification system for eukaryotic**
14 858 **transposable elements**. *Nature Reviews Genetics* 2007, **8**(12):973-982.
15 859 65. Morgante M, Hanafey M, Powell W: **Microsatellites are preferentially associated with**
16 860 **nonrepetitive DNA in plant genomes**. *Nature Genetics* 2002, **30**(2):194-200.
17 861 66. Wang Q, Fang L, Chen J, Hu Y, Si Z, Wang S, Chang L, Guo W, Zhang T: **Genome-Wide**
18 862 **Mining, Characterization, and Development of Microsatellite Markers in Gossypium**
19 863 **Species**. *Scientific Reports* 2015, **5**(1).
20 864 67. Sonah H, Deshmukh R, Sharma A, Singh V, Gupta D, Gacche R, Rana J, Singh N, Sharma T:
21 865 **Genome-Wide Distribution and Organization of Microsatellites in Plants: An Insight**
22 866 **into Marker Development in Brachypodium**. *PLOS ONE* 2011, **6**(6):e21298.
23 867 68. Bernardi G: **Isochores and the evolutionary genomics of vertebrates**. *Gene* 2000,
24 868 **241**(1):3-17.
25 869 69. Mizuno M, Kanehisa M: **Distribution profiles of GC content around the translation**
26 870 **initiation site in different species**. *FEBS letters* 1994, **352**(1):7-10.
27 871 70. Amit M, Donyo M, Hollander D, Goren A, Kim E, Gelfman S, Lev-Maor G, Burstein D,
28 872 Schwartz S, Postolsky B *et al*: **Differential GC content between exons and introns**
29 873 **establishes distinct strategies of splice-site recognition**. *Cell reports* 2012, **1**(5):543-556.
30 874 71. Wendel JF, Greilhuber J, Dolezel J, Leitch IJ: **Plant Genome Diversity Volume 1 - Plant**
31 875 **Genomes, their Residents, and their Evolutionary Dynamics**, vol. 1. Wien: Springer-
32 876 Verlag; 2012.
33 877 72. JiaYan W, JingFa X, LingPing W, Jun Z, HongYan Y, ShuangXiu W, Zhang Z, Jun Y:
34 878 **Systematic analysis of intron size and abundance parameters in diverse lineages**.
35 879 *Science China Life Sciences* 2013, **56**(10):968-974.
36 880 73. Yu J, Yang Z, Kibukawa M, Paddock M, Passey D, Wong G: **Minimal Introns Are Not**
37 881 **"Junk"**. *Genome Research* 2002, **12**(8):1185-1189.
38 882 74. Zhu J, He F, Wang D, Liu K, Huang D, Xiao J, Wu J, Hu S, Yu J: **A Novel Role for Minimal**
39 883 **Introns: Routing mRNAs to the Cytosol**. *PLOS ONE* 2010, **5**(4):e10144.
40 884 75. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO: assessing**
41 885 **genome assembly and annotation completeness with single-copy orthologs**.
42 886 *Bioinformatics* 2015, **31**(19):3210-3212.
43 887 76. Ye J, Fang L, Zheng H, Zhang Y, Chen J, Zhang Z, Wang J, Li S, Li R, Bolund L *et al*: **WEGO:**
44 888 **a web tool for plotting GO annotations**. *Nucleic Acids Research* 2006, **34**(Web Server
45 889 issue):W293-W297.
46 890 77. Cooper GM: **The Cell, 2nd edition. A Molecular Approach**. Sunderland (MA): Sinauer
47 891 Associates; 2000.
48 892 78. Sulpice R, Trenkamp S, Steinfath M, Usadel B, Gibon Y, Witucka-Wall H, Pyl E-T, Tschoep
49 893 H, Steinhauser MC, Guenther M *et al*: **Network Analysis of Enzyme Activities and**
50 894 **Metabolite Levels and Their Relationship to Biomass in a Large Panel of -Arabidopsis**
51 895 **Accessions**. *The Plant Cell* 2010, **22**(8):2872-2893.
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4 896 79. Hamilton JP, Robin Buell C: **Advances in plant genome sequencing**. *The Plant Journal*
5 897 2012, **70**(1):177-190.
6 898 80. Barthelson R, McFarlin AJ, Rounsley SD, Young S: **Plantago: Modeling Whole Genome**
7 899 **Sequencing and Assembly of Plant Genomes**. *PLOS ONE* 2011, **6**(12):e28436.
8 900 81. Hellsten U, Wright KM, Jenkins J, Shu S, Yuan Y, Wessler SR, Schmutz J, Willis JH, Rokhsar
9 901 **DS: Fine-scale variation in meiotic recombination in Mimulus inferred from population**
10 902 **shotgun sequencing**. *Proceedings of the National Academy of Sciences of the United*
11 903 *States of America* 2013, **110**(48):19478-19482.
12 904 82. Cruz F, Julca I, Gómez-Garrido J, Loska D, Marcet-Houben M, Cano E, Galán B, Frias L,
13 905 Ribeca P, Derdak S *et al*: **Genome sequence of the olive tree, Olea europaea**.
14 906 *GigaScience* 2016, **5**(1):29.
15 907 83. Plomion C, Aury JM, Amsellem J, Alaeitabar T, Barbe V, Belser C, Berges H, Bodenes C,
16 908 Boudet N, Boury C *et al*: **Decoding the oak genome: public release of sequence data,**
17 909 **assembly, annotation and publication strategies**. *Molecular ecology resources* 2016,
18 910 **16**(1):254-265.
19 911 84. Park B-S, Lee H-K, Lee S-E, Piao X-L, Takeoka G, Wong R, Ahn Y-J, Kim J-H: **Antibacterial**
20 912 **activity of Tabebuia impetiginosa Martius ex DC (Taheebo) against Helicobacter pylori**.
21 913 *Journal of Ethnopharmacology* 2006, **105**(1-2):255-262.
22 914 85. Gómez Castellanos R, Prieto J, Heinrich M: **Red Lapacho (Tabebuia impetiginosa)—A**
23 915 **global ethnopharmacological commodity?** *Journal of Ethnopharmacology* 2009,
24 916 **121**(1):1-13.
25 917 86. Byeon S, Chung J, Lee Y, Kim B, Kim K, Cho J: **In vitro and in vivo anti-inflammatory**
26 918 **effects of taheebo, a water extract from the inner bark of Tabebuia avellaneda**.
27 919 *Journal of Ethnopharmacology* 2008, **119**(1):145-152.
28 920 87. Koyama J, Morita I, Tagahara K, Hirai K-I: **Cyclopentene dialdehydes from Tabebuia**
29 921 **impetiginosa**. *Phytochemistry* 2000, **53**(8):869-872.
30 922 88. Hussain H, Krohn K, Ahmad VU, Miana GA, Green IR: **Lapachol: An overview**. *Arkivoc*
31 923 2007, **2007**(2):145.
32 924 89. Widhalm J, Rhodes D: **Biosynthesis and molecular actions of specialized 1,4-**
33 925 **naphthoquinone natural products produced by horticultural plants**. *Horticulture*
34 926 *Research* 2016, **3**:16046.
35 927 90. Romagnoli M, Segoloni E, Luna M, Margaritelli A, Gatti M, Santamaria U, Vinciguerra V:
36 928 **Wood colour in Lapacho (Tabebuia serratifolia): chemical composition and industrial**
37 929 **implications**. *Wood Science and Technology* 2013, **47**(4):701-716.
38 930 91. Alkan C, Sajjadian S, Eichler EE: **Limitations of next-generation genome sequence**
39 931 **assembly**. *Nat Meth* 2011, **8**(1):61-65.
40 932 92. Silva-Junior OB, Grattapaglia D, Novaes E, Collevatti RG: **Supporting data for "Genome**
41 933 **assembly of the pink Ipê (Handroanthus impetiginosus, Bignoniaceae), a highly-valued**
42 934 **ecologically keystone Neotropical timber forest tree"**. *GigaScience* database 2017.
43 935 <http://dx.doi.org/10.5524/100379>
44 936
45
46
47
48
49
50
51
52
53 937
54
55
56
57
58
59
60
61
62
63
64
65

938 **Table 1.** *Handroanthus impetiginosus* genome assembly statistics. The final assembly for each
 939 step contains scaffolds of length 1 kbp or longer.

Scaffold sequences	Allpaths-LG	Allpaths-LG/ Sspace/GapClose	Allpaths-LG/Sspace/ GapClose/Reapr
Number	57,815	16,090	13,206
Total size, without gaps (bp)	469,049,393	565,959,143	476,867,120
Total size, with gaps (bp)	614,626,609	586,542,612	503,314,177
Number > 10 Kbp	10,029	8,602	8,348
Number > 20 Kbp	6,920	6,791	6,647
Number > 100 Kbp	1,100	1,709	1,304
Number > 1 Mbp	2	0	0
Longest sequence (bp)	1,844,569	979,053	558,523
Average size (bp)	10,631	36,454	38,112
N50 length (bp)	57,726	97,266	80,946
L50 count	2,595	1,792	1,906
GC %	33.63	33.57	33.62

940

941

942 **Table 2.** *Handroanthus impetiginosus* gene prediction statistics with respect to the number,
 943 length and base composition of genes, transcripts, exons and introns.

944

	Genes	Transcripts	Exons	Introns
Number	31,688	35,479	154,209	122,521
Average number/gene	-	1.12	4.87	3.87
Average length	3,129	3,342	285	445
N50 length	4,421	4,643	477	839
%GC	38.38	38.22	42.60	32.83
%N	0.43	0.43	0.00	0.29

945

946

1
2
3
4 947 **Table 3.** The distribution of the minimal introns (53–125 bp) and the minimal-intron-containing
5
6 948 genes – as the number of genes with at least one minimal intron – from selected plant species in
7
8 949 comparison to the *H. impetiginosus* genome assembly.
9

10 950

Species	Genome size (Mbp)	Number of intron (bp)	Mean intron length (bp)	Minimal intron (%)	Gene (%)
<i>A. thaliana</i> (Rosids)	120	118,037	164	72.29	57.08
<i>E. guttata</i> (Asterids)	312	117,507	290	47.75	57.63
<i>P. trichocarpa</i> (Rosids)	423	166,809	380	36.96	53.41
<i>E. grandis</i> (Rosids)	691	137,329	425	33.49	48.38
<i>S. indicum</i> (Asterids)	354	101,313	439	38.14	49.76
<i>H. impetiginosus</i> (Asterids)	557	122,521	445	34.36	49.78
<i>S. lycopersicum</i> (Asterids)	900	125,750	543	36.09	47.78

11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28 951
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 952 **Figure Legends**

5
6 953

7
8 954 **Figure 1.** The *Handroanthus impetiginosus* (Mart. ex DC.) Mattos (syn. *Tabebuia*
9 955 *impetiginosa*, Bignoniaceae), tree UFG-1 whose genome was sequenced.

10
11 956

12
13 957 **Figure 2. Depth of coverage analysis.** (A) Histograms of k-mer frequencies in the filtered read
14 958 data for k = 25 (red) and GenomeScope modeling equation on *H. impetiginosus* (blue). The x-axis
15 959 shows the number of times a k-mer occurred (coverage). The vertical dashed dark blue lines
16 960 correspond to the mean coverage values for unique heterozygous k-mers (left peak) and unique
17 961 homozygous k-mers (right peak). (B) Density plot of read depth based on mapping all short
18 962 fragment reads back to the assembled scaffolds (red). Left peak (at depth = 34x) corresponds to
19 963 regions where the assembler created two distinct scaffolds from divergent putative haplotypes.
20 964 The right peak (at depth = 67x) contains scaffolds from regions where the genome is less
21 965 variable, allowing the assembler to construct a single contig combining homologue sequences.
22 966 Histograms of Poisson modeling for read depth in the assembly (green, lambda = 34; blue,
23 967 lambda = 67) are shown.

24 968

25 969 **Figure 3. Depth of coverage analysis for the haplotype-reduced assembly.** (A) Density plot of
26 970 read depth based on mapping all short fragment reads back to the haplotype-reduced
27 971 assembled sequences after identification and removal of redundant sequences due the
28 972 structural heterozygosity in the genome. (B) Density plot for average sequencing coverage per-
29 973 scaffold on the final assembly. The observed number of scaffolds in the final haplotype-reduced
30 974 assembly and the respective read coverage (blue line) is shown in comparison to a Poisson
31 975 process approximation (red line) with lambda = 63x, the observed average sequencing coverage
32 976 in the useful read data.

33 977

34 978 **Figure 4. Repeat content of the *H. impetiginosus* genome assembly.** (A) The density of
35 979 interspersed and tandem repeat as percent of the assembly. The size of the circles represents
36 980 the number of copies in the assembly for each family of repeats; (B) Distribution of sizes of the
37 981 consensus sequences for repeat families identified using *de novo* and homology methods for
38 982 repeat characterization.

39 983
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 984 **Figure 5. Transcriptome quality assessment** (A) similarity search of *H. impetiginosus* putative
5 peptides against source database of plant protein sequences using BLASTP algorithm (e-value
6 985 1e-6). Transcript count means the number of peptides of *H. impetiginosus* with best hit against
7 986 the source database using bit-score and grouping results by taxon name. Transcript score
8 987 corresponds to the average bit-score overall hits for each group using the best hit. We ordered
9 988 taxon groups by their average bit-score overall hits and used Welch's t-test to compare the
10 989 distributions of bit-score hits between two adjacent groups with p-values <0.01 (ns = non-
11 990 significant; *** significant); (B) Completeness of the expected gene space of the genome
12 991 assembly, estimated with BUSCO. The estimates were compared with genome annotations for
13 992 other lamids, *Erythranthe guttata* and *Olea europaea*.
14 993
15 994

16 995 **Figure 6.** Contig termini analysis to investigate the possible genomic features associated with
17 996 gaps in the genome assembly. Contigs were created from the genome assembly with the "cutN -
18 997 n 1" command from seqtk program, which cut at each gap (of at least one basepair, i.e. one or
19 998 more Ns). The figure shows the percentage of contig termini (the position of the terminal
20 999 nucleotides of each contig) intersecting with different annotations of the genome.
21
22
23

1000

24 1001 **Figure 7. Genes of the biosynthetic pathway of specialized quinoids.** O-succinylbenzoate (OSB)
25 1002 pathway depicting the number of *H. impetiginosus* (Himp) annotated genes for the known
26 1003 enzymes that lead to the biosynthesis of the naphthoquinones, including lapachol. For
27 1004 comparison, it also shows the numbers of genes for the closely related *Mimulus guttatus*
28 1005 (Mgut), *Solanum lycopersicum* (Slyc), for the model *Arabidopsis thaliana* (Ath), and for the tree
29 1006 species *Eucalyptus grandis* (Egr) and *Populus trichocarpa* (Potri). The pathway was modified
30 1007 from [89].
31 1008

1009 **Supplementary material**

1010

1011 **Table S1.** Summary of the sequence data generated for the genome assembly of *Handroanthus*
1012 *impetiginosus* based on the ALLPATHS-LG algorithm.
1013

1014 **Figure S1.** Flow cytometry results of the sequenced tree UFG-1 of *H. impetiginosus*. Flow
1015 cytometry estimate of the nuclear DNA content was carried out using young leaf tissue on a BD

61
62
63
64
65

1
2
3
4 1016 Accuri™ C6 Plus personal flow cytometer. *Pisum sativum* (genome size 9.09 pg/2C or ~4380
5
6 1017 Mb/1C) was used as standard for comparison (M2). The estimate of nuclear DNA content for *H.*
7
8 1018 *impetiginosus* (M1) averaged over 10 readings was 1.155 pg/2C or 557.3 ± 39 Mb/1C.

9
10 1019

11 1020 **Figure S2.** Overview of the analytical pipeline with the bioinformatics steps and tools employed
12
13 1021 for genome (black arrows) and transcriptome assembly (red arrows), and for gene prediction
14
15 1022 and annotation (blue arrows). Bioinformatics programs are indicated in italic, blue, and the main
16
17 1023 file formats in red. The input sequences are highlighted in yellow boxes and the main products
18
19 1024 in green.

20 1025

21
22 1026 **Figure S3.** Distribution and characterization of simple sequence repeats in *Handroanthus*
23
24 1027 *impetiginosus* genome (A) Histogram of different motifs ranging from 1 to 6 bp (B) Distribution
25
26 1028 of the simple sequence repeats length detected in the genome assembly.

27
28 1029

29 1030 **Figure S4.** Comparison of the gene features parameters, such as number and length, between *H.*
30
31 1031 *impetiginosus* and the other selected dicot plant across distinct lineages of Rosids (*A. thaliana*
32
33 1032 and *P. trichocarpa*) and Asterids (*E. guttata* and *S. lycopersicum*). Frequency histograms are
34
35 1033 shown according to the whole-genome gene content annotation for (A) the complete predicted
36
37 1034 gene structure (B) exons and (C) introns. Dashed vertical lines are the average lengths for the
38
39 1035 gene features.

40 1036

41 1037 **Figure S5.** Histograms for Gene Ontology broader term annotations in the *H. impetiginosus*
42
43 1038 genome assembly. Terms for the Biological Process ontology were summarized with WEGO by
44
45 1039 the second tree level setting. The Pearson Chi-Square test was applied to indicate significant
46
47 1040 relationships between *H. impetiginosus* and the lamid *Erythranthe guttata* regarding the
48
49 1041 number of genes (at $\alpha \geq 5\%$). (A) Terms displaying remarkable relationship between the two
50
51 1042 datasets; (B) terms with a significant difference between the two datasets.

52 1043

53
54 1044 **Figure S6.** Same as Figure S6 but showing comparison between numbers of genes assigned to
55
56 1045 GO broader terms for *H. impetiginosus* and the lamid *Olea europaea*.

57 1046
58
59
60
61
62
63
64
65

1
2
3
4 1047 **Figure S7.** Sequence length distribution from the assemblies of *H. impetiginosus* and other two
5
6 1048 highly heterozygous trees of the genus *Quercus*. Figure shows density plots for the size of
7
8 1049 scaffolds with 2 kbp or longer in the three assemblies. Contigs metrics were computed by
9
10 1050 cutting at each gap (of at least 25 base pair, i.e. 25 or more Ns). Scaffolds and contigs length
11
12 1051 were plotted using the common logarithm to respond to skewness towards large values.
13

14 1052

15 1053 **File S1.** Evidences adopted to support protein-coding loci identification and assignment in the
16
17 1054 *H. impetiginosus* genome assembly. Two qualifiers – high-confidence and low-confidence – were
18
19 1055 added to the locus based on the reported evidences.
20

21 1056

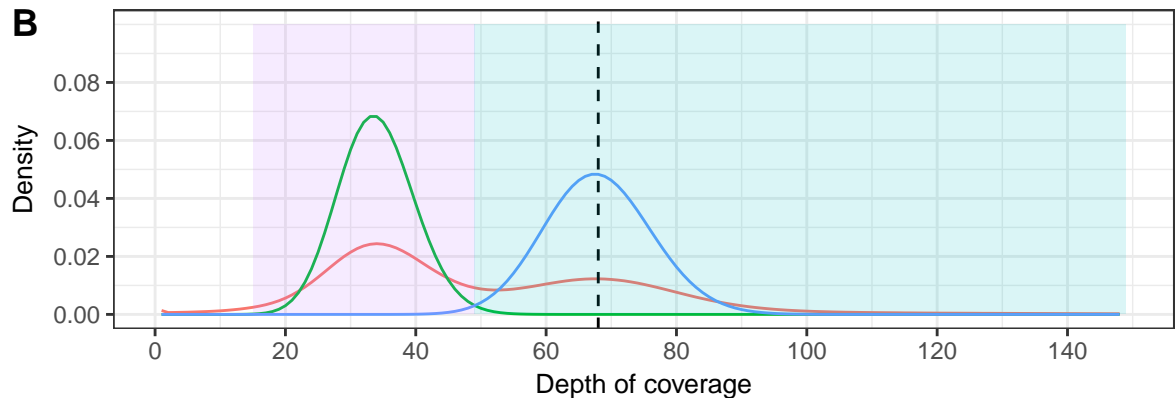
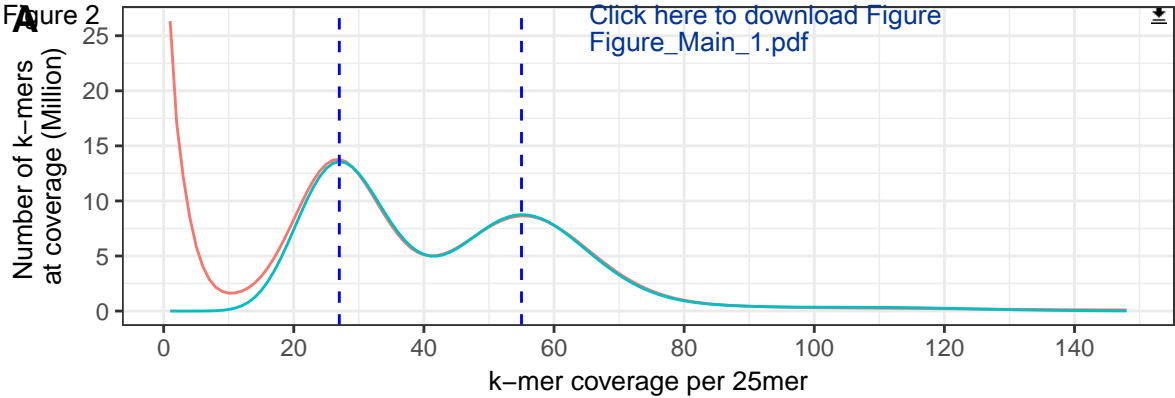
22 1057 **File S2.** Genome assembly metrics from the assemblies of *H. impetiginosus* and other two highly
23
24 1058 heterozygous trees of the genus *Quercus*. Comparison between metrics based on the
25
26 1059 `assemblathon_stats` script part of the `assemblathon2-analysis` package
27
28 1060 (<https://github.com/ucdavis-bioinformatics/assemblathon2-analysis>). Metrics were computed
29
30 1061 for scaffolds with 2 kbp or longer in length. Genomic sequences in scaffolds for *Quercus lobata*
31
32 1062 was obtained from <https://valleyoak.ucla.edu/genomicresources/> (accessed on 9/20/2017). For
33
34 1063 *Quercus rubra*, genomic sequences in scaffolds were downloaded from the ENA (European
35
36 1064 Nucleotide Archive) repository, accessions LN776247-LN794156.
37

38 1065

39 1066

40 1067
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65





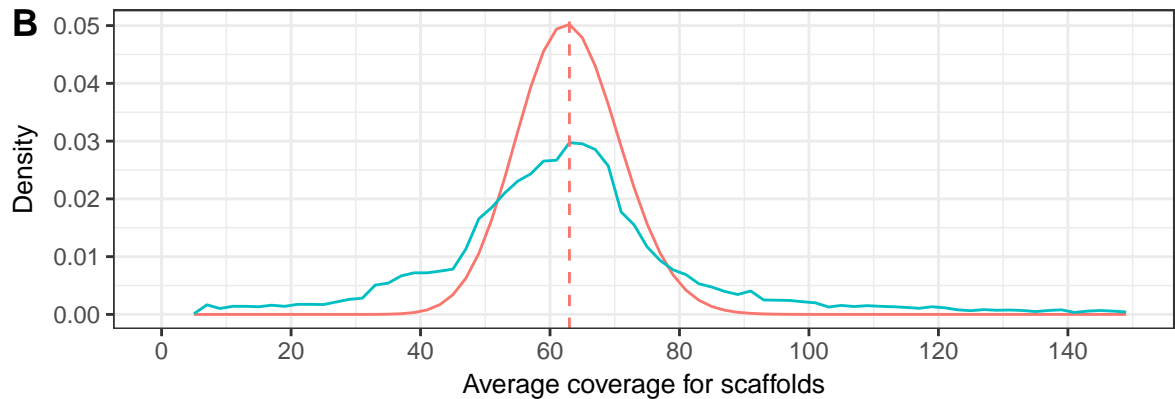
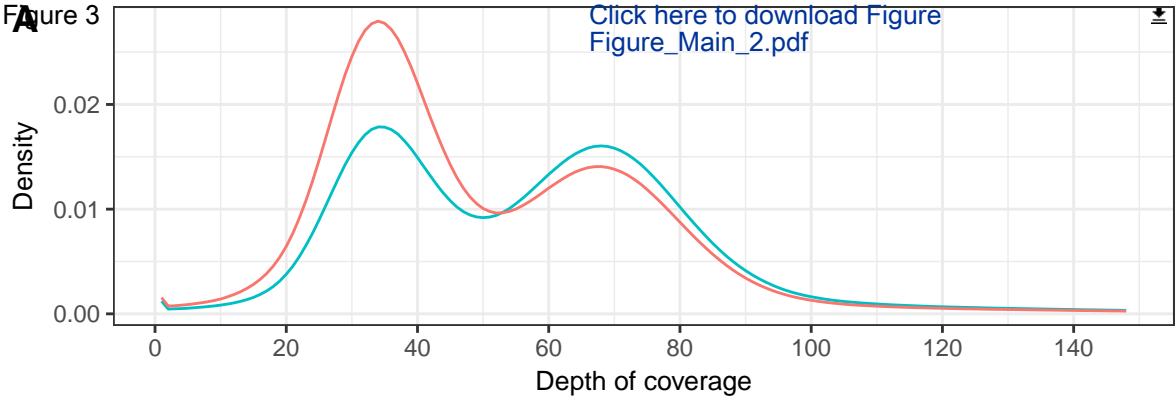
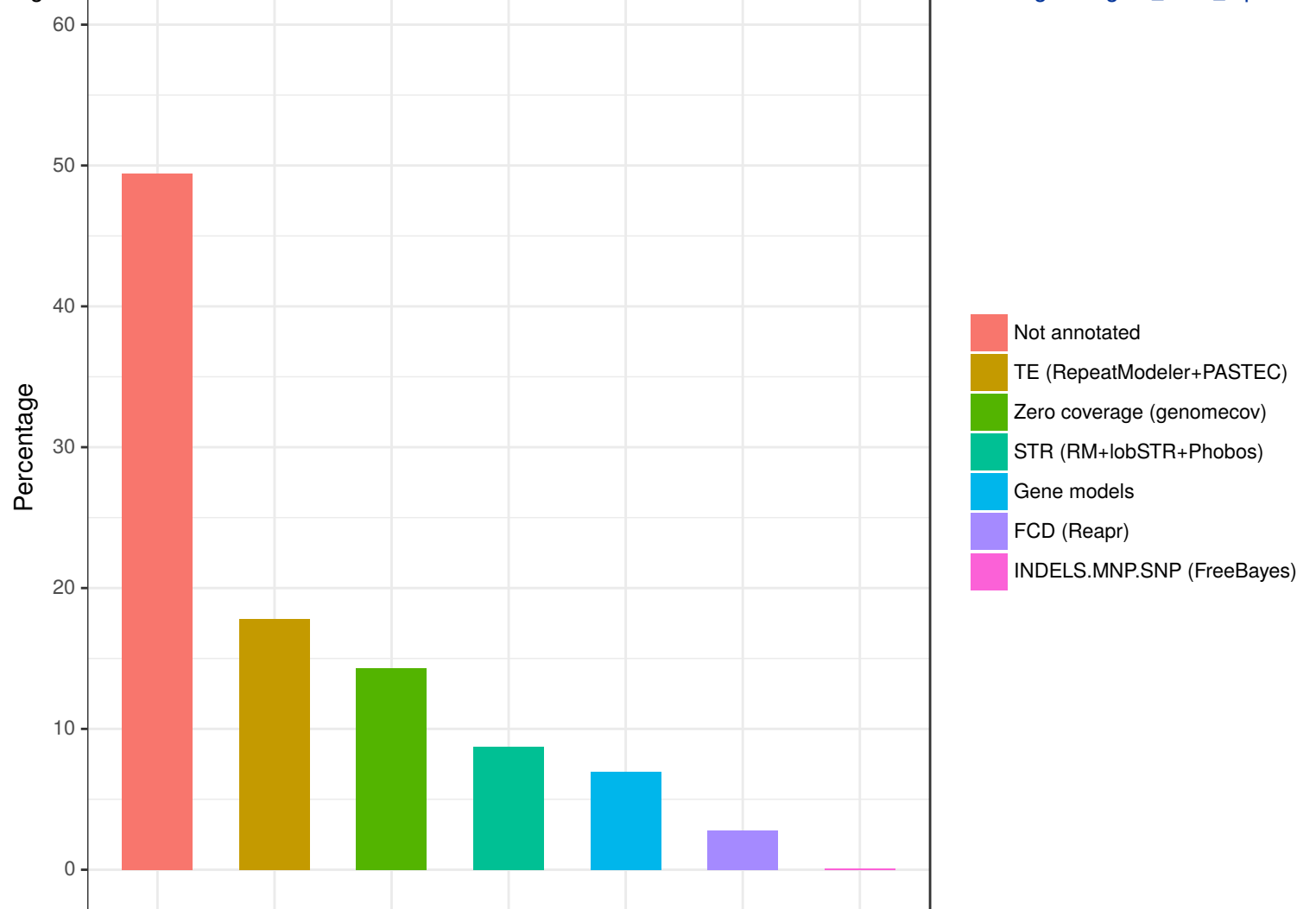
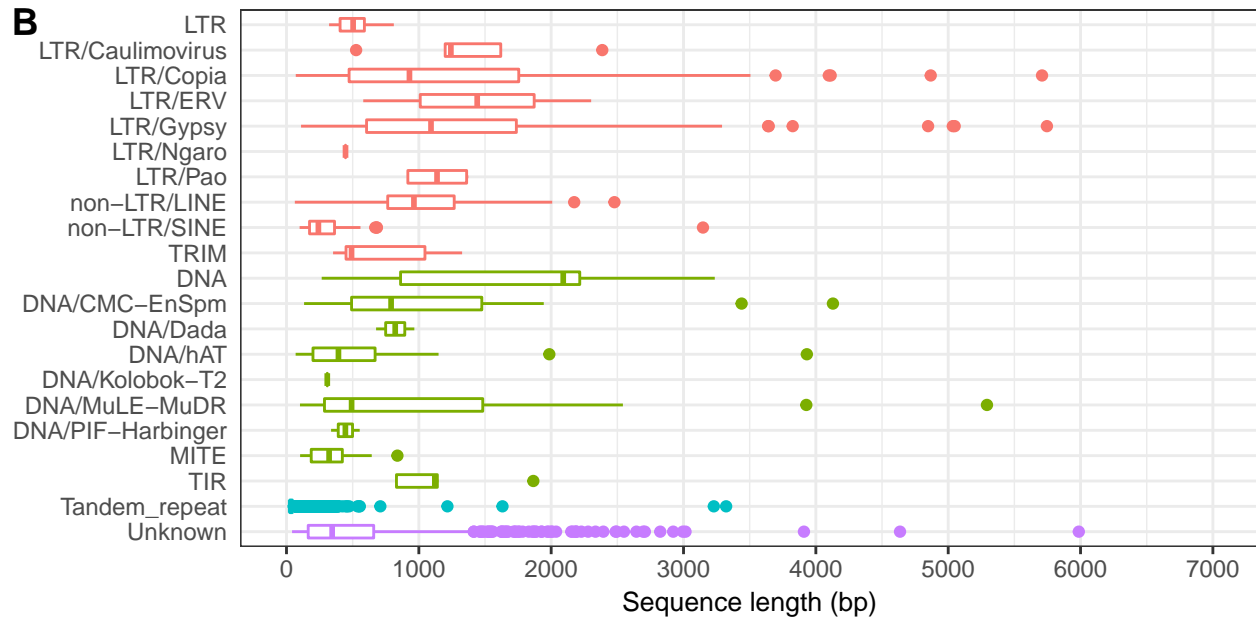
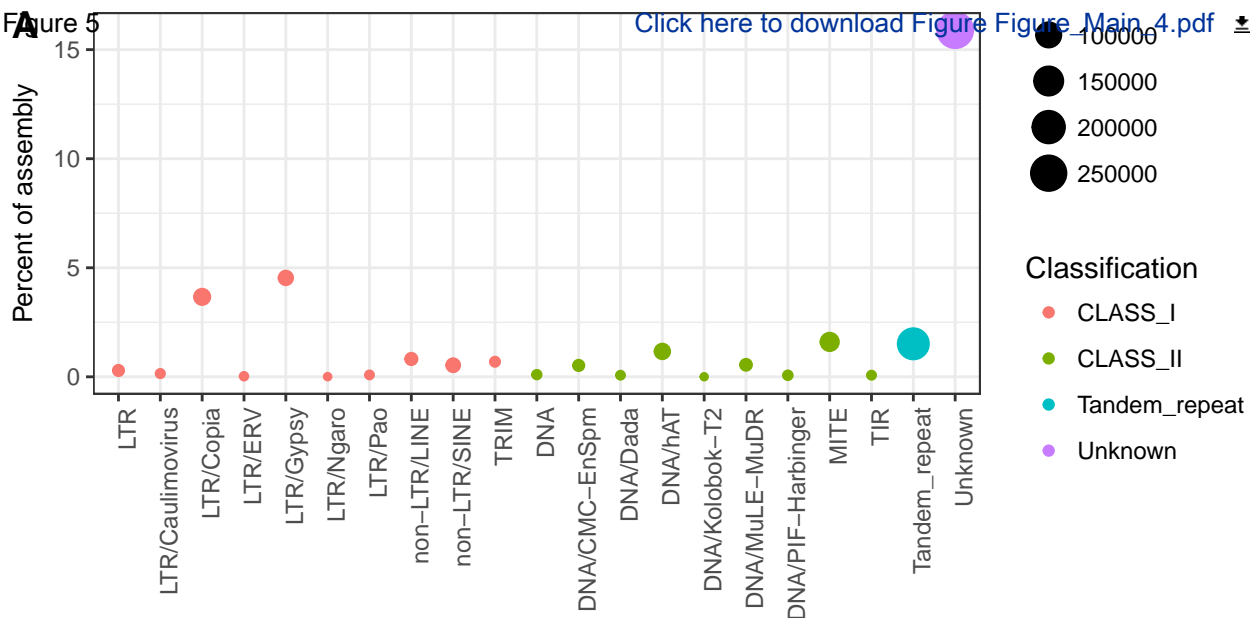
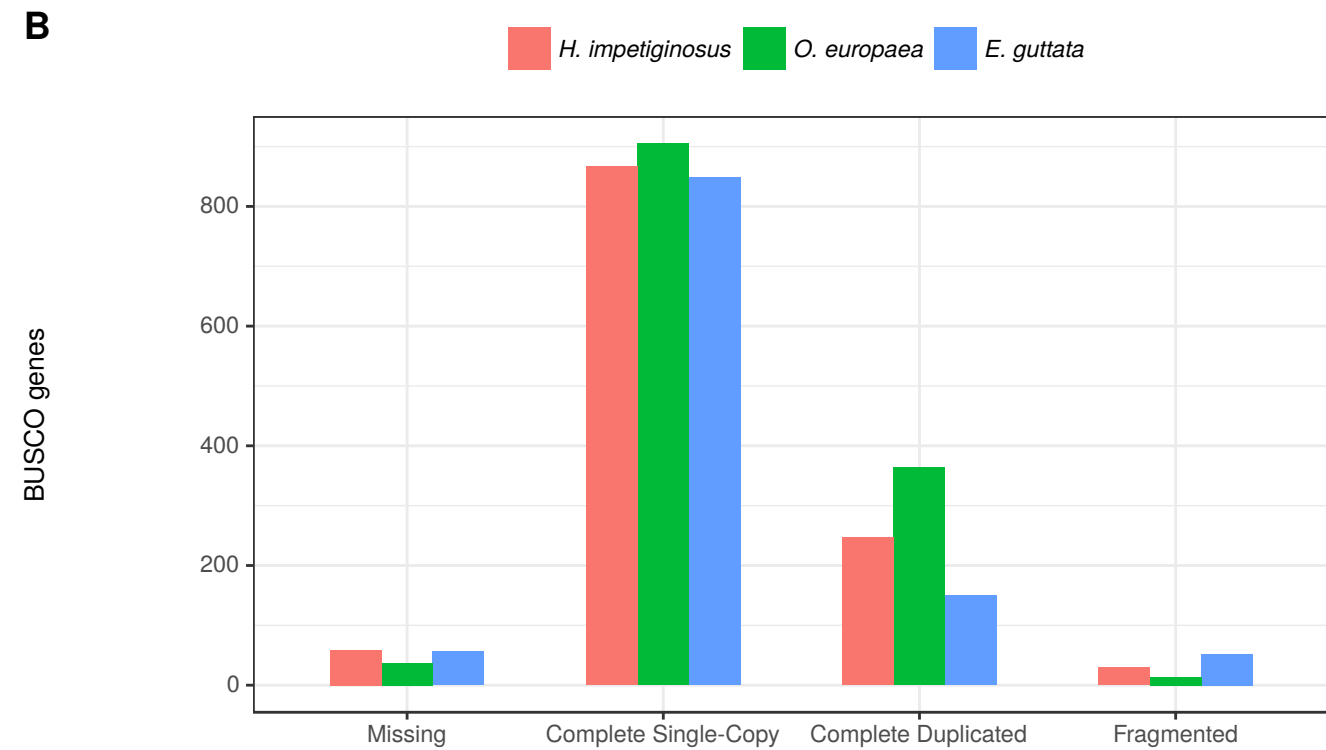
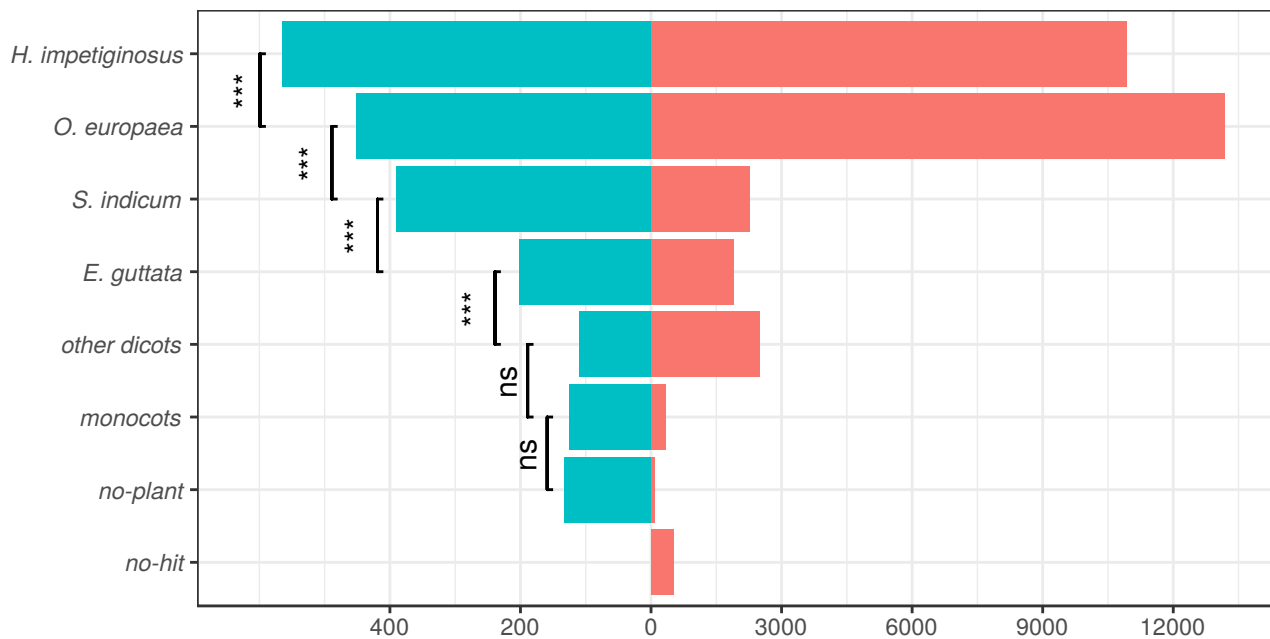


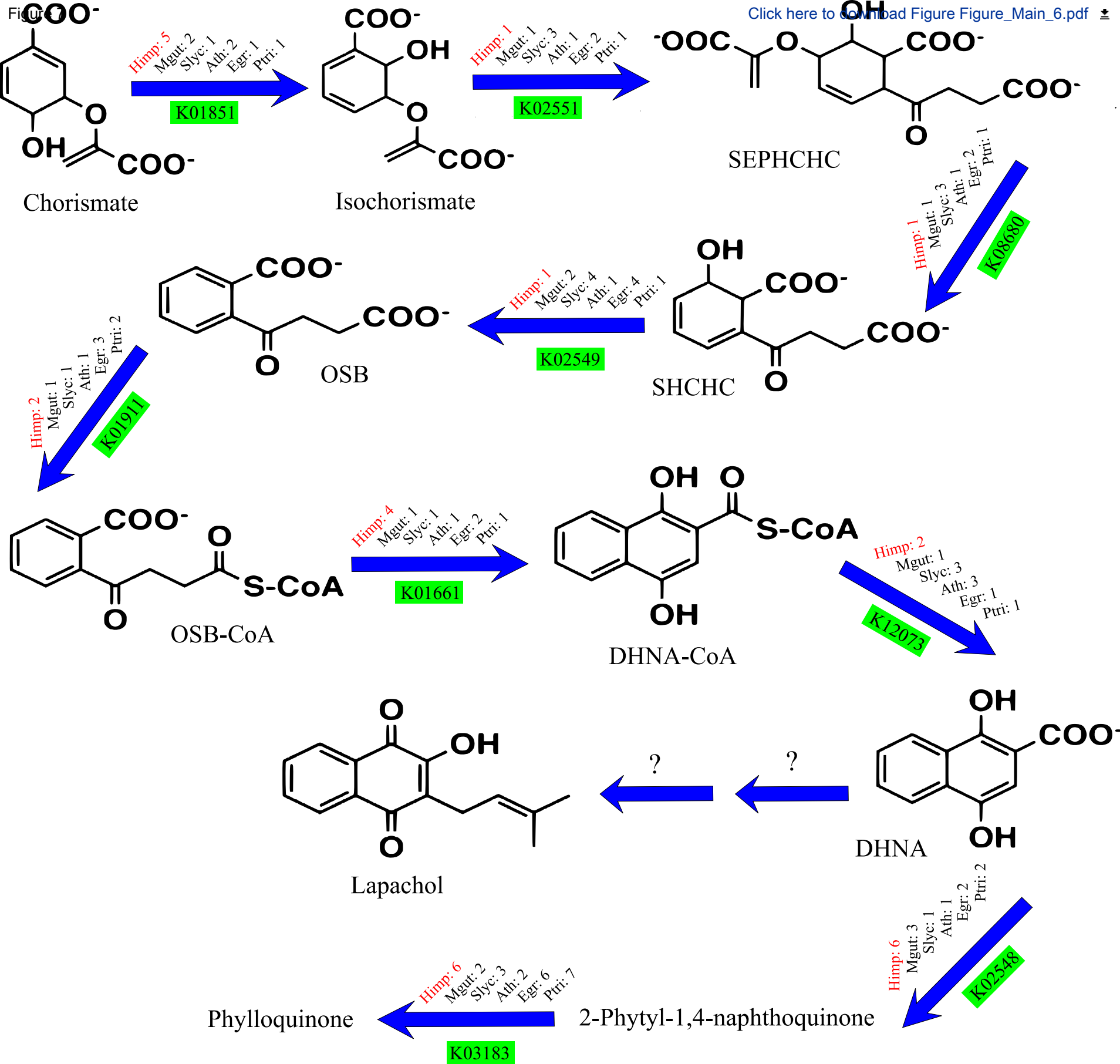
Figure 4

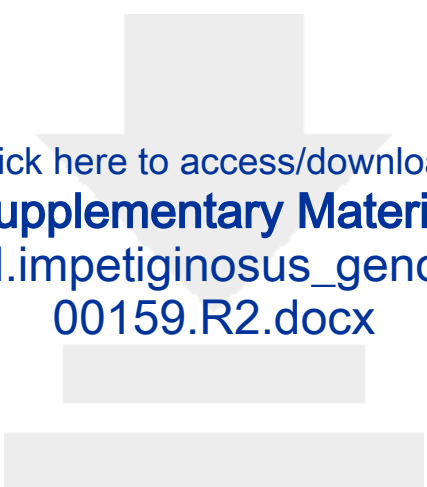
[Click here to download Figure Figure_Main_3.pdf](#)











Click here to access/download

Supplementary Material

Supp_Material_H.impetiginosus_genome_GIGA-D-17-
00159.R2.docx



Click here to access/download
Supplementary Material
Supplementary_FileS1.xlsx





Click here to access/download
Supplementary Material
Supplementary_FileS2.xlsx

