

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

<b>Manuscript Number:</b>	GIGA-D-17-00159R1	
<b>Full Title:</b>	Genome assembly of the pink Ipê ( <i>Handroanthus impetiginosus</i> , Bignoniaceae), a highly-valued ecologically keystone Neotropical timber forest tree	
<b>Article Type:</b>	Data Note	
<b>Funding Information:</b>	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
<b>Abstract:</b>	<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 and the most logged species in Brazil, under significant illegal trading pressure. The plant produces large amounts of quinoids, specialized metabolites with documented antitumor 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 covered 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>	
<b>Corresponding Author:</b>	Rosane Garcia Collevatti, PhD Universidade Federal de Goias Goiania, GO BRAZIL	
<b>Corresponding Author Secondary Information:</b>		
<b>Corresponding Author's Institution:</b>	Universidade Federal de Goias	
<b>Corresponding Author's Secondary Institution:</b>		
<b>First Author:</b>	orzenil.silva@embrapa.br B Silva-Junior, PhD	
<b>First Author Secondary Information:</b>		
<b>Order of Authors:</b>	orzenil.silva@embrapa.br B Silva-Junior, PhD	
	Dario Grattapaglia, PhD	

	Evandro Novaes, PhD
	Rosane Garcia Collevatti, PhD
<b>Order of Authors Secondary Information:</b>	
<b>Response to Reviewers:</b>	<p>Goiânia, 27 September 2017</p> <p>Professor Hans Zauner Assistant Editor GigaScience</p> <p>Dear Prof. Zauner,</p> <p>Enclosed you will find the manuscript GIGA-D-17-00159-R1, the revised version of the originally submitted one entitled “Genome assembly of the pink Ipê (<i>Handroanthus impetiginosus</i>, Bignoniaceae), a highly-valued ecologically keystone Neotropical timber forest tree” that we are re-submitting for publication in GigaScience as Data Note. Thank you for the opportunity to revise our manuscript. We have responded in great detail to all comments and suggestions made by the reviewers in the text following this letter. We sincerely appreciated the reviewers’ suggestions that really helped us improving the manuscript and we hope that all the issues raised were addressed at satisfaction.</p> <p>We reiterate our confidence that this work opens exceptional opportunities to empower molecular studies of the several species of the <i>Tabebuia</i> Alliance based on useful genomic resources for genetic and functional analysis in the species. In the revised version we added more data and analyses to show that this read set, genome assembly and corresponding annotations are consistent contributions to the genomics community. Their availability should inspire the development of similar genomic resources for the still largely neglected forest trees of the mega-diverse tropical biomes.</p> <p>To facilitate the visualization of how we have dealt with each single point, we have reiterated each concern or suggestion made by the reviewer in bold and placed our reply (RESPONSE:) immediately after, in the section following this cover letter, the entire document entitled “Cover Letter+Responses_to_GigaScience MS GIGA-D-17-00159-R1”. Hope you and reviewers will consider this as a more acceptable version for publication in GigaScience</p> <p>Thank you for your consideration. We look forward to hearing from you.</p> <p>Sincerely yours, Rosane Collevatti (on behalf of all co-authors) Laboratório de Genética &amp; Biodiversidade, Instituto de Ciências Biológicas, Universidade Federal de Goiás, CP 131, 74001-970, Goiânia, GO, Brasil. E-mail: rosanegc68@hotmail.com. Phone: +55 62 3521-1729</p> <p>Responses to Reviewer #1 The main problem I have is not being able to access the data on NCBI. I can see two biosamples and a bioproject, but I cannot find the SRA records, the genome or the annotations. The biosample record page will usually have a link at the bottom to the SRA, but this is missing (both only link back to the bioproject). Further if you search by the taxonomic name of the tree through Entrez, the number of records for genome is 0 and for SRA is also 0. Possibly the authors are waiting for final release until this is published? But this definitely needs to be taken care of prior to publication. RESPONSE: The sequence data has been deposited in the SRA of NCBI and are available through the following accession identifiers: 1.Study: whole-genome sequencing 1.1Sample name: Himp-UFG1 1.2BioSample: SAMN05195323 1.3SRA: SRS1483442; Run accessions (5): SRR3624821 - SRR3624825 2.Study: mRNA-Seq 2.1Sample name: HIMP-UFG1-RNA-POOL</p>

2.2BioSample: SAMN07346903  
2.3SRA: SRS2349699; Run accession (1): SRR5820886  
3.Genome assembly: this Whole Genome Shotgun project has been deposited at DDBJ/ENA/GenBank under the accession NKXS00000000. The version described in this paper is version NKXS01000000.  
Lines 91-92 "beyond a relatively small numbers of microsatellites with their caveats for more sophisticated genetic analyses." Numbers should be number and also clarification of what the word caveats is referring to.  
RESPONSE: We clarified the sentence and added the number of microsatellites available before our study was developed and the corresponding references. We also added a reference regarding the well known caveats of microsatellites for evolutionary analyses. The sentence now reads: "The species has virtually no genomic tools and resources, beyond a handful of 21 microsatellites [17] with their caveats for more sophisticated genetic analyses in the areas of population structure, admixture, and effective population size and evolution [18]."  
Lines 112-113 - mention of other species being sequenced, but this isn't addressed in the results. Should remove from methods if these results are not addressed in this manuscript  
RESPONSE: Thanks. We removed the sentence. In fact, this pooled sample was used in a second study that addressed SNP discovery, which is the subject of a separate manuscript.  
Line 115 - Min should be Mini  
RESPONSE: Corrected  
Line 123 - I'm not sure if "jump" is the official word or just jargon, maybe consider "fragment lengths of"  
RESPONSE: Actually, the jargon "jump" is used for such mate-pair library. However we agree that "fragment lengths of" reads better. We changed accordingly.  
Line 133 - mention of a "perl script" - this needs to be made available through github or somewhere else  
RESPONSE: The perl script TrimAdaptor.pl was provided to our group by its authors at the High-Throughput Sequencing and Genotyping Center Unit of the University of Illinois Urbana-Champaign. We added this information in the section "AVAILABILITY OF SUPPORTING DATA". The script was uploaded to the GigaScience's data repository following the permission of the original authors.  
Line 139 - same comment about "jumping" - I would consider changing to "mate pair"  
RESPONSE: Changed to mate-pair  
Line 402 - "showed" to something else, perhaps "illustrated"  
RESPONSE: Changed to "depicted".  
Line 432 - adequately doesn't really make sense in this sentence, perhaps remove  
RESPONSE: Agreed. Removed.  
Line 439 - The citation to Fig S6A doesn't make much sense, the text refers to blast results but the figure shows Gene ontology terms (and the figure is cited later in the GO section), maybe a supplemental figure is missing?  
RESPONSE:  
We corrected the figure Fig S6A to better convey the result presented. The figure regarding gene ontology terms analysis is now depicted as Fig S7.  
Line 477 - It's not clear how the search for the genes was done (BLAST?)  
RESPONSE: We clarified this in the manuscript by adding the marked phrases below:  
"Given their medicinal and biological relevance, we have searched the H. impetiginosus annotated genes for the enzymes involved in the biosynthesis of naphthoquinones. By searching for the KEGG identifiers of these enzymes (e.g. K01851) in the InterPro annotation results, we have found all the important known enzymes that lead to the biosynthesis of lapachol (Fig. 6). Unfortunately, however, the last two steps of the lapachol biosynthesis pathway still constitute unidentified enzymes [79]. For comparative purposes, we downloaded the annotation file of five other species from the Phytozome version 11 (Fig. 6). The number of H. impetiginosus genes encoding for the enzymes of each step in the pathway is comparable to the numbers found in other species."

Responses to Reviewer #2  
The article describes the genome assembly of Handroanthus impetiginosus, a neotropical timber tree. Because of heterozygosity in the diploid genome that was sequenced, the final assembly is fragmented (N50=81,316bp, L50=1906). The

assembly fragmentation might be an issue for future analyses, and the authors should be more specific about that. Some parts of the text should be rewritten in order to acknowledge the fact that the assembly obtained is highly fragmented. For instance, the sentence (line 337) "Our genome assembly metrics are similar to recent reports of genome assemblies of other highly heterozygous forest tree genomes", should be discarded for two reasons: first, if other heterozygous genomes were assembled in a highly fragmented way, the authors should not be satisfied with doing "as bad", but should aim at doing better. Second, the metrics obtained for *Quercus robur* were actually better than those obtained for pink Ipê (N50=260kb, L50=1468)

RESPONSE: We have looked at these issues in depth and revised our manuscript accordingly. We agree with the reviewer that we need to be more specific about the limitation of the assembly for further analysis. We disagree, however, with the consideration that the assembly fragmentation might be an issue for future analysis. At most, it might be an issue for comparative genomics analysis and studies of genome evolution, but this is a well-known limitation of unfinished de novo assembly generated primarily with short reads like this one (see Alkan et al. 2011. Limitations of next-generation genome sequence assembly. *Nature Methods*. doi: 10.1038/nmeth.1527). Evolutionary issues and large-scale comparative genomics are beyond the scope of our manuscript. We added a word of caution to end users regarding this particular issue. It was not our intention or claims to provide a fully finished resource for the research community but only a quality level, well-curated, extensively supervised assembly. We have described in the manuscript that our main goal was to provide useful genomic resources for genetic and functional analysis in the species (Lines 341-344; 420-422). Nevertheless, we admit that in parts of the manuscript we have suggested that these resources can empower evolutionary studies (Lines 56; 71; 502). Although we are confident that such studies can be carried out using our assembly, at least at the gene-level of gene-family level, we did not validate it in this present research. We thus rewrote those sentences to tone down any suggestion about using this genome assembly for deeper evolutionary analysis.

Difficulties in obtaining a highly contiguous assembly for this highly repetitive and heterozygous genome was anticipated and discussed in the manuscript. Comparison provided with other similar assemblies was clarified. The N50 metric mentioned by the reviewer is however a misleading one as discussed elsewhere (See Bradnam et al. 2013 Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. DOI: 10.1186/2047-217X-2-10), but we do agree that we should have added more information in the manuscript regarding this comparison. Our brevity in this respect certainly contributed to the misinformed judgment about the quality of our assembly when the reviewer alluded that we "were satisfied with doing this genome assembly as bad". A careful and deeper comparison of our assembly with assemblies for other heterozygous trees clearly shows that our assembly is, in fact, better and particularly so when compared to the one for *Quercus* that the reviewer mentions. Please, find as additional supplementary material to the manuscript a new figure and a file that brings a detailed comparison of our assembly metrics with other three recently published assemblies for highly heterozygous trees (Figure S8 and File S2).

Certainly, there are still inaccuracies at the base and assembly level in this unfinished assembly but we tried hard not to deliver it to end users without an appropriate documentation, which make this initial read set and sequence a primary and very reliable object for further improvement. To clarify the main features of our assembly that highlight its quality and to better respond to the reviewer's suggestion, we also improved the Re-use Potential section.

Abstract, line 47: the terms "redundancy in the consensus determination" are not clear. Figure 2B shows that most scaffolds correspond to a consensus between the two haplotypes. What does "redundancy" mean? Does it mean that for some parts of the genome the two haplotypes were assembled separately? The sentence should be clarified.

RESPONSE: Thanks for the remark. To complement the depth of read coverage analyses, we added a new analysis – contig termini analysis – carried out to identify the most probable causes of breaks in the assembly. Contig termini are the positions of the terminal nucleotides of each contig from the genome assembly created by cutting at each gap (of at least one base pair, i.e. one or more Ns). This analysis was developed using a protocol described elsewhere (Tørresen et al. 2017. *BMC Genomics* 2017, 18(1):95). As noted by the reviewer, figure 2B shows that most scaffolds correspond to a consensus between the two haplotypes. The contig termini analysis reinforces this

view showing that virtually none of the breaks are caused by allelic variants that could arise from splitting contigs of divergent haplotypes within a locus. Thus, there is no redundancy in the assembly in the sense that the two haplotypes might have been assembled separately. We removed this sentence in the manuscript.

Line 156: Figure S1 is called but it does not correspond to the pipeline (it should be figure S3). Along the manuscript numerous calls to figures do not correspond to the actual list of figures: the numbering and calling of figures should be checked carefully.

RESPONSE: All the figure and table numbers both in the main text and supplementary material were checked carefully and corrected where necessary.

Line 262 (316, table 1): The metrics (N50, L50...) are most of the time given for the whole assembly and for sequences longer than 20kb. What is the size of the smaller scaffolds in the assembly (was a threshold set)? It would be interesting to provide metrics (in table 1 also) for scaffolds >1kb or >2kb: how much of the genome is included in such scaffolds? Filtering out short scaffolds from the assembly should be envisaged, since no genes can be annotated in short scaffolds. It would have the advantage of providing better assembly metrics.

RESPONSE: We agree that filtering out short scaffolds is adequate. In fact, in our study, we have considered only scaffolds of length 1 kbp or longer in each step of the assembly process. We explicitly added this information to the manuscript where necessary. We added this phrase in the title of Table 1: "All assembly step only contain scaffolds of length 1 kbp or longer". We added the phrase "of length 1 kbp or longer" in parts of the main manuscript in which we were describing the assembly metrics. We also added a supplementary table providing metrics of how much of the assembly is included in scaffolds of length 2 kbp or longer. To illustrate the quality of our assembly we also included File S1 with a detailed comparison of our assembly with other recent assemblies for diploid, highly heterozygous tree genomes.

Line 285: the N50 is given for scaffolds >1kb, and not for all scaffolds: is it possible to provide the same metrics for all assemblies in order to be able to compare them?

RESPONSE: In our study, we discarded short scaffolds of length smaller than 1 kbp, and metrics were provided considering all scaffolds in the assembly. We rewrote the phrase as follows: "The final genome assembly after REAPR breaks had 19,319 sequences of length 1 kbp or longer, with 576,929,188 bp. N50 size of scaffolds dropped from 97 kbp (L50 = 1,792) to 71 kbp (L50 = 2,379).".

Lines 281-288: What proportion of the sequence-coverage differences called by REAPR correspond to boundaries between regions where alleles were assembled separately vs collapsed? If most of those errors are due to heterozygosity, would it make sense not to use the coverage information to break the scaffolds? Is there an option in REAPR in order to avoid the breaks? The procedure to filter out redundant copies of unmerged haplotypes might then require to split scaffolds, but it might result in less breaks in the assembly. Can the authors discuss on that?

RESPONSE: Thanks for pointing out these issues. We looked at these points in depth and our considerations follow below to provide some clarification. We did not say that "most errors are due to heterozygosity". In lines 274-279, we wrote that the most frequent error reported in the Reapr analysis (~92%) was caused by reads that mapped ambiguously to the assembly. It can be due to interspersed repeats with high-copy families, short-tandem repeats and low-copy repeats such as segmental duplications; all of them are known to be highly represented in eukaryotic genomes. These highly similar sequences reduce the depth of coverage reported for each base but do not imply in breaks of scaffolds unless the calculated fragment coverage distribution error (FCD) is higher than a theoretical value threshold. So, it is not the coverage information at each base that leads to a break but the FCD error at each base related to a background noise (the FCD error threshold). As correctly observed by the reviewer, it does not make sense to use the coverage information to break the scaffolds. Due to the variance in depth expected in the shotgun process and the presence of repeats, it could result in the introduction of large number of false breaks. Our only intention in removing putative pairs of sequences that looked like haplotypes of the same locus possibly assembled separately, was to provide a reliable resource for further variant calling analysis. It was not this procedure that required to split scaffolds. The procedure to identify regions in the assembly that have fragmented distribution around bases causing an FCD error was carried out to evaluate the overall and local accuracy of the assembly process and to prevent misassemblies. It should be a mandatory process for every de novo assembly due to the difficult decisions taken by different programs to resolve repeats and fuse contigs containing allelic variants in a



unique consensus sequence. Of course, we do not advocate the use of FCD but any similar measure of accuracy of the process should be used. For instance, FRCurve measure implemented in Amosvalidate, which is included as part of the AMOS assembly package is another interesting alternative (Vezi et al. Feature-by-Feature – Evaluating De Novo Sequence Assembly PLoSOne (DOI: 10.1371/journal.pone.0031002). We have carried out this Reapr analysis just to eliminate fragrant misassemblies that could hinder further use of this genomic resource.

We added a new analysis in the manuscript (contig termini analysis) to detect the proportion of features with overlap to contigs end in the assembly. Contig termini are the positions of the terminal nucleotides of each contig from the genome assembly created by cutting at each gap (of at least one base pair, i.e. one or more Ns). This analysis was carried out using the protocol described elsewhere (Tørresen et al. 2017. BMC Genomics 2017, 18(1):95 doi: 10.1186/s12864-016-3448-x).

In the final assembly, the proportion of contig ends that overlap to breaks suggested by the FCD analysis using Reapr is very low (<3%), i.e. there are very few detected misassemblies. Moreover, virtually none of contig ends overlap to allelic variants annotated with FreeBayes using read data mapped to the final genome assembly. Contig termini overlap most prominently (~50%) with regions that do not encompass any structural annotation in the assembly or regions that have no depth of coverage (~15%) based on mapped reads. It suggests that contig ends match large repeats not resolved given the short-read sequence data. Another possibility is that interruptions in continuity and contiguity might be due to young euchromatic segmental duplication with higher sequence similarity to the consensus sequence. This hypothesis is compelling given that the size of remaining gaps within scaffolds ( $2,384 \pm 3,167$ ) bp is not longer than the longest fragments used in the mate-pair libraries (>10 kbp). Moreover, the completeness of the assembly related to the empirically determinate genome size (90%) does not suggest that breaks of contiguity contain much longer sequences.

Lines 328, 332: Figure 1C and D are called but the described figures correspond to Figure 2A and B.

RESPONSE: Corrected

Line 348: know -> known

RESPONSE: Corrected

Line 352: Mostly -> Most

RESPONSE: Corrected

Line 356: "expand a wide range of sequence sizes" : the sentence should be corrected

RESPONSE: Corrected. It now reads: "... cover a wide range of sequence sizes from 42 bp..."

Line 361: "unknown non-classified" sequences: are not the two terms redundant?

RESPONSE: Removed the word "unknown"

Line 366: Figure 2 is called, it should be Figure 3

RESPONSE: Corrected

Line 390: 31,668 genes were annotated. This number is relatively high for a plant genome. It would be interesting to explore paralogous gene clusters: are there duplicated genes in the genome? Are these genes more likely to have arisen from WGDs or tandem repeats? If such analyses are possible despite the fragmented nature of the assembly, they would be of great interest.

RESPONSE: Actually, the number of genes annotated in the *H. impetiginosus* (~32k) can be considered quite average for a plant genome. For example, *S. lycopersicum* has ~34k protein coding genes, *Populus* 42K, *Eucalyptus* 36K and soybean 56K. Among other plants from the order Lamiales, monkey flower (*M. guttatus*) has 28k (Hellsten et al. 2013 PNAS 110:19478-19482) protein coding genes, sesame (*S. indicum*) has 23k (Zhang et al. 2013 Genome Biology 14:401) and olive tree (*O. europaea*), which is believed to have undergone whole-genome duplication (WGD) events, has 56k genes (Cruz et al. 2016 GigaScience 5:29). Nevertheless, we included in our manuscript an additional analysis to identify low-confidence set of genes based on signatures derived from loss of information due to the fragmentation of this unfinished genome assembly for *H. impetiginosus*. The results were summarized in a new Supplementary File S1 provided along with the main manuscript. If these low-confidence genes are excluded, the number of genes with reliable annotation in the genome assembly is 28,603, which is actually 90% of the total estimate of 31,668. This high-confidence subset contains approximately the same number of genes reported for

other members of the Lamiales. To account for gene families in the assembly, we explored the catalog of genes in *H. impetiginosus* by predicting orthologous groups (orthogroup) with OrthoFinder (Emms & Kelly 2015 Genome Biology 16:157). For comparative genomics, the orthology analyses was performed also with proteins from the basal plant *A. trichopoda*, the Lamiales *M. guttatus* and *O. europaea*., as well as the forest tree *P. trichocarpa*. The number of genes from each species was analyzed in each orthogroup. The percentage of orthogroups with only one gene from each species was higher in *A. trichopoda*, *H. impetiginosus* and *M. guttatus*, compared to the species (*P. trichocarpa* and *O. europaea*) that are known to have experienced WGD events. *P. trichocarpa* and *O. europaea*, on the other hand, present higher percentages of orthogroups with two or more genes (Figure 1). Given the fragmentation of the genome assembly and the lack of scaffolds mapping onto chromosomes, it is hard to analyze WGD events in the *H. impetiginosus* genome at relevant scale. Nevertheless, these orthology results provide initial evidence that this genome may have not undergone WGD events.

Figure 1. Distribution of the number of genes in orthogroups of the five species analyzed. Results indicate that *H. impetiginosus* has fewer genes in the orthogroups compared to the species *O. europaea* and *P. trichocarpa* that are known to have experienced recent WGD events.

We also investigated the possibility of tandem duplications in the *H. impetiginosus* genome by analyzing the genome location of genes from the same orthology groups. On average orthogroups with two or more genes presented 87.1% of these elements in different scaffolds. Two genes from the same scaffold were found in the same orthogroup in 10.1% of the cases, on average. The remaining 2.8% are instances where three or more genes from the same orthogroup were in the same scaffold. These results indicate that tandem duplication of genes may have not been frequent throughout the evolution of the *H. impetiginosus* genome. Even considering only scaffold co-localization, regardless of distance, our estimates are far from the high frequency of tandemly duplicated genes (~34%) observed in the *Eucalyptus* genome (Myburg et al. 2014) for example.

Line 439: Figure S6 -> Figure 4A ?

RESPONSE: Corrected

Line 446: Figure 3B -> Figure 4B?

RESPONSE: Corrected

Lines 451-452: BUSCO results were benchmarked using poplar. Poplar is known to have undergone a Whole Genome Duplication; Was the duplication status of *Handroanthus* investigated? If there was no WGD, the duplicate level is probably not comparable to that of *Populus*.

RESPONSE: Duplication status of *H. impetiginosus* was not specifically investigated. This was beyond the scope of this initial work. Following the reviewer's suggestion we have however benchmarked BUSCO results using other lamids, *Erythranthe guttata* and *Olea europaea* as described above.

Lines 454-464: GO terms were compared to those of poplar. Why was such a distant species chosen for the comparison? What about a comparison with other asterids, or even lamids (*E. guttata*, or *Olea*)?

RESPONSE: Thanks. Following the reviewer's suggestion we have performed GO terms comparison using other lamids, *Erythranthe guttata* and *Olea europaea*.

Lines 482-487: The authors report that some steps of the quinoid metabolism are encoded by more genes in pink Ipê than in other species. It would be interesting to elucidate how these genes were amplified in the *Hydroanthus* genome: is it possible to build a phylogeny of these genes ? Are some of them located on the same scaffold? (are they possibly deriving from tandem duplications?..)

RESPONSE: Analyses with OrthoFinder confirmed that many genes of each step of the quinoid pathway are indeed from the same orthologue group. We investigated whether genes from these quinoid orthogroups have arisen from tandem duplications. By analyzing if genes groups were in the same scaffold, we found little evidence of gene family expansion by tandem duplication. Of the eight quinoid orthogroups identified, containing 23 genes, in only one there were two genes localized in the same scaffold. With respect to a phylogenetic analyses of these genes, we believe that this is beyond the scope of this study. The point we want to make is that we indeed found in the genome assembly these genes from this especially important secondary pathway.

**Additional Information:**

Question	Response
<p>Are you submitting this manuscript to a special series or article collection?</p>	<p>No</p>
<p><b>Experimental design and statistics</b></p> <p>Full details of the experimental design and statistical methods used should be given in the Methods section, as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>. 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>	<p>Yes</p>
<p><b>Resources</b></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 <a href="#">Research Resource Identifiers</a> (RRIDs) for antibodies, model organisms and tools, where possible.</p> <p>Have you included the information requested as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>
<p><b>Availability of data and materials</b></p> <p>All datasets and code on which the conclusions of the paper rely must be either included in your submission or deposited in <a href="#">publicly available repositories</a> (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.</p> <p>Have you have met the above requirement as detailed in our <a href="#">Minimum Standards Reporting Checklist</a>?</p>	<p>Yes</p>



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<sup>1,2</sup>, Dario Grattapaglia<sup>1,2</sup>, Evandro Novaes<sup>3</sup>, Rosane G. Collevatti<sup>4\*</sup>  
11  
12 5

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

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

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

20 10 <sup>4</sup>*Laboratório de Genética & Biodiversidade, Instituto de Ciências Biológicas, Universidade Federal*  
21 *de Goiás, CP 131. Goiânia, GO. 74001-970. Brazil.*  
22 11

23 12  
24  
25 13 **\*Corresponding author:** Rosane Garcia Collevatti, Instituto de Ciências Biológicas, Universidade  
26  
27 14 Federal de Goiás, CP 131, 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 33

1  
2  
3  
4 34 **Abstract**

5  
6 35

7 36 **Background:** *Handroanthus impetiginosus* (Mart. ex DC.) Mattos is a keystone Neotropical  
8 37 hardwood tree widely distributed in seasonally dry tropical forests of South and Mesoamerica.  
9 38 Regarded as the "new mahogany", it is the second most expensive timber and the most logged  
10 39 species in Brazil, under significant illegal trading pressure. The plant produces large amounts of  
11 40 quinoids, specialized metabolites with documented antitumorous and antibiotic effects. The  
12 41 development of genomic resources is needed to better understand and conserve the diversity of  
13 42 the species, to empower forensic identification of the origin of timber and to identify genes for  
14 43 important metabolic compounds.  
15 44

16  
17  
18  
19  
20  
21  
22 44

23 45 **Findings:** The genome assembly covered 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 allowing  
25 47 masking ~31% of the assembly. Depth of coverage indicated that consensus determination  
26 48 adequately removed haplotypes assembled separately due to the extensive heterozygosity of the  
27 49 species. Automatic gene prediction provided 31,688 structures and 35,479 mRNA transcripts,  
28 50 while external evidence supported a well-curated set of 28,603 high-confidence models (90% of  
29 51 total). Finally, we used the genomic sequence and the comprehensive gene content annotation  
30 52 to identify genes related to the production of specialized metabolites.  
31 53

32  
33  
34  
35  
36  
37  
38 53

39 54 **Conclusions:** This genome assembly is the first well-curated resource for a Neotropical forest tree  
40 55 and the first one for a member of the *Bignoniaceae* family, opening exceptional opportunities to  
41 56 empower molecular, phytochemical and breeding studies. This work should inspire the  
42 57 development of similar genomic resources for the largely neglected forest trees of the mega-  
43 58 diverse tropical biomes.  
44 59

45  
46  
47  
48  
49 59

50  
51 60

52 61 **Keywords:** heterozygous genome, RNA-seq, transposable elements, quinoids, *Bignoniaceae*  
53 62

54 62

55 63

56 64

57 65

58

59

60

61

62

63

64

65

66

67 **DATA DESCRIPTION**

68

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

91

92 Species of *Handroanthus* and *Tabebuia* have virtually no genomic tools and resources, beyond a  
93 handful of 21 microsatellites [17] with their caveats for more sophisticated genetic analyses in the  
94 areas of population structure, admixture, and effective population size and evolution [18]. Whole-  
95 genome sequencing has now become accessible to a point that efforts to develop improved  
96 genomic resources for such species are possible and warranted. We built a preliminary assembly  
97 of the nuclear genome of a single individual of *Handroanthus impetiginosus* based on short-reads

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 98 and longer mate-pair DNA sequence data to provide the necessary framework for the  
5  
6 99 development of genomic resources to support multiple genomic and genetic analyses of this  
7  
8 100 keystone Neotropical hardwood tree regarded as the "new mahogany". It is the second most  
9  
10 101 expensive timber and the most logged species in Brazil [19], exported largely to North America  
11  
12 102 for residential decking and currently under significant illegal trading pressure. Additionally, the  
13  
14 103 tree produces large amounts of natural products such as those of quinoid systems (1,4-  
15  
16 104 anthraquinones, 1,4-naphthoquinones, and 1,2-furanonaphthoquinones), specialized  
17  
18 105 metabolites with promising antitumor, anti-inflammatory and antibiotic effects [20, 21]. The  
19  
20 106 high pressure of logging and illegal trading on this species with a notable ecological keystone  
21  
22 107 status urges conservation efforts of existing populations.  
23

## 24 109 **METHODS**

25 110  
26  
27 111 **Sample collection and sequencing.** DNA of a single adult tree of *H. impetiginosus* (UFG-1) (Figure  
28  
29 112 S1) was extracted using Qiagen DNeasy Plant Mini kit (Qiagen, DK). Flow cytometry was used to  
30  
31 113 check the genome size of tree UFG-1 indicating a genome size of (557 ±39) Mb /1C (Figure S2)  
32  
33 114 consistent with published estimates [22]. Total RNA from shoots of five seedlings and from the  
34  
35 115 differentiating xylem of the adult tree (UFG-1) was extracted using Qiagen RNeasy Plant Mini kit  
36  
37 116 (Qiagen, DK) and pooled for RNA sequencing. DNA and RNA sequencing was performed at the  
38  
39 117 High-Throughput Sequencing and Genotyping Center of the University of Illinois Urbana-  
40  
41 118 Champaign, USA. The following libraries were generated for sequencing: (1) two shotgun genomic  
42  
43 119 libraries of short fragments (300bp and 600bp) from tree UFG-1 (2) one shotgun library from  
44  
45 120 combined pools of five RNA samples tagged with a single index sequence. Paired-end sequencing,  
46  
47 121 2x150 nt, was performed in two lanes of an Illumina HiSeq 2500 instrument (Illumina, CA, USA).  
48  
49 122 Three additional mate-pair libraries (fragment lengths of 4kb to 5.5kb, 8kb to 10kb and 15kb to  
50  
51 123 20kb) for UFG-1 were also sequenced in two lanes of an Illumina HiSeq 2000 instrument (2x101  
52  
53 124 bp). This long-range sequence resource was used to generate the final genome assembly for  
54  
55 125 annotation. A complete overview of the genome assembly and annotation pipeline is provided  
56  
57 126 (Figure S3).  
58

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

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

20 138  
21 139 ***Transposable elements and repetitive DNA.*** Repetitive elements were detected and annotated  
22  
23 140 on the genome assembly with the RepeatModeler *de novo* repeat family identification and  
24  
25 141 modeling package [28]. Using RECON, RepeatScout and Tandem Repeat Finder, repetitive  
26  
27 142 sequences were detected in the scaffolds longer than 10 kb using a combination of similarity-  
28  
29 143 based and *de novo* approaches. The TE sequences were evaluated using modeling capabilities of  
30  
31 144 the RepeatModeler program, with default settings, to compare the TE library against the entire  
32  
33 145 assembled sequences and to refine and classify consensus models of putative interspersed  
34  
35 146 repeats. A complementary analysis intended to augment the number of TE sequences classified  
36  
37 147 according to current criteria [29] was performed using the PASTEC program [30]. RepeatMasker  
38  
39 148 Open-4.0 [31] was used with the sequences from the *de novo* repetitive element library to  
40  
41 149 annotate the interspersed repeats and to detect simple sequence repeats (SSRs) on the genome  
42  
43 150 assembly.

44 151  
45 152 ***Protein-coding genes annotation.*** Protein-coding genes annotation was performed with a  
46  
47 153 pipeline that combines RNA-seq assembled transcript and protein alignments to the reference  
48  
49 154 with *de novo* predictions methods (Figure S3). RNA-Seq reads were screened for the presence of  
50  
51 155 adapters, which were removed using *Fastq-mcf* [23]. *Trimmomatic* [32] was used to (1) remove  
52  
53 156 low quality, no base called segments (N's) from sequencing reads; (2) scan the read with a 4-base  
54  
55 157 sliding window, cutting when the average quality per base dropped below 15; and (3) remove  
56  
57 158 reads shorter than 32 bp after trimming. Trimmed reads mapped to mitochondrial, chloroplast  
58  
59 159 and ribosomal sequences from plants with *Bowtie1* [24] (options  $-v\ 3\ -a\ -m\ 1$ ) were also  
60  
61 160 removed. Transcript *de novo* assemblies were performed using *SOAP-Transdenovo* [33] and  
62  
63 161 *Trinity de-novo* [34] from the processed reads. The assemblies were concatenated and used as  
64  
65



1  
2  
3  
4 162 input to *EvidentialGene* [35], a comprehensive transcriptome pipeline to identify likely complete  
5  
6 163 coding regions and their proteins in the final, combined, transcriptome assembly. Gene modeling  
7  
8 164 was carried out using standard procedures and tools described, for instance, in [36]. In summary,  
9  
10 165 a genome-guided transcriptome assembly of *H. impetiginosus* was performed with the JGI  
11  
12 166 PERTRAN RNA-seq Read Assembler pipeline [37] using both the RNA-Seq trimmed reads and  
13  
14 167 sequences from the *de novo* transcript assembly. Loci were identified by the assembled transcript  
15  
16 168 alignments using BLASTX [38] and EXONERATE [39] alignments of peptide sequences to the  
17  
18 169 repeat-soft-masked genome using RepeatMasker [40], based on a transposon database  
19  
20 170 developed as part of this genome assembly annotation. Known peptide sequences included  
21  
22 171 manually curated data sets for plant species available from UniProtKB/Swiss-Prot [41] and  
23  
24 172 sequences available from Phytosome [1] version 11 for *Arabidopsis thaliana*, *Oryza sativa*,  
25  
26 173 *Erythranthe guttata*, *Solanum lycopersicum*, *Solanum tuberosum*, *Populus trichocarpa* and *Vitis*  
27  
28 174 *vinifera*. Gene structure were predicted by homology-based predictors, FGENESH++,  
29  
30 175 FGENESH\_EST [42, 43] and GenomeScan [44]. Gene predictions were improved by PASA [45],  
31  
32 176 including adding UTRs, correcting splicing and adding alternative transcripts. PASA-improved gene  
33  
34 177 model peptides were subjected to peptide homology analysis with the above-mentioned  
35  
36 178 proteomes to obtain Cscore values and peptide coverage. Cscore is the ratio of the peptide  
37  
38 179 BLASTP score to the mutual best hit BLASTP score, and peptide coverage is the highest percentage  
39  
40 180 of peptide aligned to the best homolog. A transcript was selected if its Cscore value was greater  
41  
42 181 than or equal to 0.5 and its peptide coverage was greater than or equal to 0.5 or if it had transcript  
43  
44 182 coverage but the proportion of its coding sequence overlapping repeats was less than 20%. For  
45  
46 183 gene models where greater than 20% of the coding sequence overlapped with repeats, the Cscore  
47  
48 184 value was required to be at least 0.9 and homology coverage was required to be at least 70%  
49  
50 185 to be selected. Selected gene models were then subjected to classification analysis using  
51  
52 186 *InterProScan 5* [46] for PFAM domains, PANTHER, Enzyme Commission Number (EC) and KEGG  
53  
54 187 categories. Gene ontology annotation was obtained, where possible, from Interpro2GO and  
55  
56 188 EC2GO mappings.

## 57 189

### 58 190 **DATA VALIDATION AND QUALITY CONTROL**

59 191

60 192 ***Global properties of the H. impetiginosus tree genome from the unassembled reads.*** Sequencing  
61  
62 193 of the *H. impetiginosus* tree genome generated c. 599 million reads, comprising 73 Gbp of  
63  
64  
65

1  
2  
3  
4 194 sequence data. This represents nearly 132× the expected sequence coverage. After removal of  
5  
6 195 adaptors, followed by standard error correction and trimming with ALLPATHS-LG, with default  
7  
8 196 options, c. 46 Gbp of data was found useful for the assembly process, yielding sequencing  
9  
10 197 coverage of 82x (63x from the fragments libraries and 19x from the mate pair libraries). The  
11  
12 198 estimated physical coverage was 400x based on the observed fragment size distributions (Table  
13  
14 199 S1). ALLPATHS-LG k-mer spectrum frequency analysis (at K=25) on useful reads, error corrected  
15  
16 200 reads, estimated a haploid genome size of 540,968,531 bp, a repeat fraction of 38.0%, and a SNP  
17  
18 201 rate of 1/88 bp (1.14%). An alternative analysis of the k-mer frequencies using GenomeScope [26]  
19  
20 202 produced a haploid genome size estimate of 503,748,072 bp, repetitive content of 36.6% and SNP  
21  
22 203 rate of 1/60 bp (1.65%). Both estimates (Figure 1A) are consistent with the flow cytometry  
23  
24 204 estimates and in line with the expectations regarding the heterozygous content of the *H.*  
25  
26 205 *impetiginosus* genome, a predominantly outcrossed tree [47]. Sequencing errors caused an  
27  
28 206 extreme peak at  $k = 1$  in the k-mer frequency distribution. Both k-mer histograms display two  
29  
30 207 distinct peaks comprising the largest area of each histogram at depths 27 and 55. The bimodal  
31  
32 208 distributions characterize the expected behavior for k-mer frequencies of a heterozygous diploid  
33  
34 209 genome as seen, for example, in the recently reported Oak genome [48]. In the right homozygous  
35  
36 210 peak (at K=55), k-mers are shared between the two homologous chromosomes. The left or  
37  
38 211 heterozygous peak, with half the k-mer depth of the homozygous peak, contains k-mers that are  
39  
40 212 unique to each haplotype due to heterozygosity. The difference in height between these peaks  
41  
42 213 (heterozygous/homozygous ratio) is a measure of the heterozygosity within the genome, which  
43  
44 214 is 1.65% according to the GenomeScope modeling equation.

215

216 *[Insert Figure 1 here]*

217

218 **Genome assembly.** State-of-the-art haploid genome assembler pipelines from short-reads  
219 ALLPATHS-LG [27] and SOAPdenovo2 [49] were considered for an initial evaluation on the dataset  
220 of reads. Two relatively new algorithms specifically developed for de novo assembly of  
221 heterozygous genomes, MaSuRCA [50] and PLATANUS [51], were also attempted as alternatives  
222 to the other two assemblers designed for genomes of low heterozygosity. Reads were first  
223 preprocessed and error corrected using the algorithms provided by each assembler. PLATANUS  
224 was set to run but after 10 weeks it did not produce any result in an Intel(R) Xeon(R) server with  
225 64 X7560 2.27GHz CPUs, 256 GB RAM, except for the k-mer count table on the input trimmed

1  
2  
3  
4 226 reads. After 9 week-long runtimes in an Intel(R) Xeon(R) server with 64 X7560 2.27GHz CPUs, 512  
5  
6 227 GB RAM, MaSuRCA successfully completed the generation of the super-reads from the trimmed  
7  
8 228 reads but the process was aborted on the overlap-correction process in the Celera Assembler due  
9  
10 229 to excessive CPU usage. SOAPdenovo2 ran very fast (3 days) but produced an assembly with total  
11  
12 230 scaffold size of 860 Mbp. Analysis with SOAPdenovo2 was run with different k-mer sizes, from 31  
13  
14 231 to 71, step of 10, but none of them produced a reasonable assembly size in view of the expected  
15  
16 232 size estimated by flow cytometry and the k-mer frequency. ALLPATHS-LG was therefore used to  
17  
18 233 assemble the genome with default options. The short reads from fragmented libraries were error-  
19  
20 234 corrected using default settings (K-mer size of 24, ploidy of 2), fragment-filled and assembled into  
21  
22 235 initial unipaths (k-mer size of 96, ploidy of 2). Jumping reads from the mate-pair libraries were  
23  
24 236 then aligned to the unipaths and all alignments were processed in a seed-extension strategy with  
25  
26 237 junction point recognition within the read aimed to remove invalid and duplicate fragments to  
27  
28 238 perform error correction and initial scaffolding. This initial process produced an assembly graph  
29  
30 239 that was turned into scaffolds by analyzing branch points in the graph topology. This late process  
31  
32 240 converted single-base mismatches into ambiguous base codes at branch. It also flattened some  
33  
34 241 other structural features of the assembly including short indels. The contig assembly comprised  
35  
36 242 109,064 sequences of length 500 bp or longer with total length of 466,314,780 bp. Genome  
37  
38 243 assembly after scaffolding comprised 57,815 scaffolds of length 1 kbp or longer with total length  
39  
40 244 of 610,091,865 bp and N50 of 57 Kbp. The fraction of bases captured in gaps was 23.9% and the  
41  
42 245 rate of ambiguous bases for all bases captured in the assembly was 0.24%. This assembly was only  
43  
44 246 slightly larger in size (<10%) than the empirically determined genome size using flow cytometry.  
45  
46 247

47  
48 248 **Alternative scaffold and gap-filling.** Although the ALLPATHS-LG performance was good in  
49  
50 249 recovering the expected genome size in the assembled contigs there was a high fraction of the  
51  
52 250 bases captured in gaps in the scaffolds ( $\sim \frac{1}{4}$  of the total genome assembly). *De novo* assembly  
53  
54 251 algorithms applied to moderate-to-high levels of heterozygosity cannot match the performance  
55  
56 252 achieved in assemblies of homozygous genomes, especially at the contig assembly level [52]. We  
57  
58 253 thus used the assembled contigs to perform an alternative scaffolding step with SSPACE [53] using  
59  
60 254 the error-corrected short fragment reads and the jumping reads. In this approach, genome  
61  
62 255 assembly comprised 16,090 scaffolds of length 1 kbp or longer with total length of 577,446,088  
63  
64 256 bp and N50 of 95 Kbp, respectively. The fraction of bases captured in gaps dropped from 23.9%  
65  
66 257 to 18.9% in contrast to ALLPATHS-LG scaffolding, totaling 109,533,288 bp. The rate of ambiguous

1  
2  
3  
4 258 bases for all bases captured in the assembly dropped from 0.24% to 0.13%. All preprocessed reads  
5  
6 259 were reused in an attempt to close the intra-scaffold gaps using the *GapCloser* [54] algorithm.  
7  
8 260 Genome assembly after gap-filling was 586,206,884 bp in 15,671 scaffolds of length 1 kbp or  
9  
10 261 longer and only 20,583,469 bp (3.51% of the genome assembly) remained in 24,907 gaps. N50 of  
11  
12 262 scaffolds of length 1 kbp or longer, with gaps, was 97,344 Kb (L50 = 1,792). Sequences longer than  
13  
14 263 20 kb were assembled in only 6,791 scaffolds totaling 538,102,146 bp, ~97% of the genome size  
15  
16 264 estimated from flow cytometry (557 Mb).

17 265

18 266 ***Evaluation of accuracy of the genome assembly.*** A subset of fragments and jumping read pairs  
19  
20 267 (~15x sequencing coverage each) were used to uncover inaccuracies in the genome assembly.  
21  
22 268 Scaffolds with identified errors were broken or flagged for inspection. REAPR [55] was used to test  
23  
24 269 each base of the genome assembly looking for small local errors (such as a single base  
25  
26 270 substitutions, and short insertions or deletions) and structural errors (such as scaffolding errors)  
27  
28 271 located by means of changes to the expected distribution of inferred sequencing fragments from  
29  
30 272 the mapped reads using SMALT v0.7.6 [56]. REAPR reported that only 343,588,027 (~60%) bases  
31  
32 273 in the assembly should be free of errors, with 5,476 reported (1,658 within contigs, 3,818 over  
33  
34 274 gaps) in the remaining 242,618,857 bp. The most frequent (~92%) type of inaccuracy reported  
35  
36 275 was *Perfect\_cov* and *Link*. *Perfect\_cov* means low coverage of perfect uniquely mapping reads  
37  
38 276 while *Link* describes situations in which reads map elsewhere in the assembly. The recognition of  
39  
40 277 this inaccuracy at the base pair level should thus reflect the repetitive nature of the genome as  
41  
42 278 inferred from the k-mer frequency spectra analysis (~36-38% of repeats). Besides the base pair  
43  
44 279 inaccurate calls due to repeats, other structural problems in the assembly were identified based  
45  
46 280 on sequence-coverage differences from the expected fragment size distribution and the program  
47  
48 281 used this information to break these. Given the high heterozygosity and divergence between  
49  
50 282 haplotypes on this diploid genome sequence, homologous sequences can assemble separately or  
51  
52 283 merge. Moreover, unresolved repeat structures in the assembly might also contribute heavily to  
53  
54 284 this issue. Structural errors in REAPR were likely called at the boundaries of these regions. The  
55  
56 285 final genome assembly after REAPR breaks had 19,319 sequences of length 1 kbp or longer, with  
57  
58 286 576,829,188 bp. N50 size of scaffolds dropped from 97,344 Kb (L50 = 1,792) to 71,491 bp (L50 =  
59  
60 287 2,379). The number of remaining gaps in the assembly was 21,417 totaling 30,066,113 bp (5.05%).

61 288

62 289

1  
2  
3  
4 290 Paired-end reads from the short fragment libraries were aligned back independently to this  
5  
6 291 genome assembly using SMALT (map -r 0 -x -y 0.5; default alignment penalty scores). Per-scaffold  
7  
8 292 depth of coverage was computed, regardless of mapping quality, using GATK DepthofCoverage.  
9  
10 293 The mean read depth across the scaffolds resulted in 66.45x. The mean read length of the mapped  
11  
12 294 reads was 139.8 bp and the corresponding k-mer coverage for size of 25 was 55.04x which  
13  
14 295 matches with the homozygous peak computed from the k-mer frequency distribution from the  
15  
16 296 unassembled reads. The read depth frequencies are shown in Figure 1B. The  
17  
18 297 heterozygous/homozygous peak height (> 1) in the distribution suggests that the assembly  
19  
20 298 contains redundant copies of unmerged haplotypes due to the structural heterozygosity of the  
21  
22 299 diploid genome of the species. To specifically deal with the heterozygosity we introduced a step  
23  
24 300 to, leniently, recognize and remove alternative heterozygous sequences. Sequences of scaffolds  
25  
26 301 were aligned one versus all using BLAT [57] and results were concatenated in a single file of  
27  
28 302 alignments and sorted. Similar sequences were identified on the base of pairwise similarity using  
29  
30 303 filterPSL utility from AUGUSTUS [58] with default parameters, and retaining all best matches to  
31  
32 304 each single sequence queried against all others that satisfy minimal percentage of identity  
33  
34 305 (minId=92%) and minimal percentage of coverage of the query read (minCover=80%). We  
35  
36 306 considered as heterozygous redundant those scaffolds that showed pairwise similarity to exactly  
37  
38 307 another sequence and their depth of coverage fell in a Poisson distribution with parameters given  
39  
40 308 by the heterozygous peak of the read depth distribution over all scaffolds ( $\lambda = 34$ ; Figure  
41  
42 309 1B). The final step was to keep only one copy – the largest one – of the heterozygous scaffolds  
43  
44 310 among pairs with high similarity.

45  
46 311  
47 312 **A preliminary assembly of the *H. impetiginosus* genome.** At the end of the accuracy evaluation  
48  
49 313 processes, the genome assembly had a total size of 503,308,897 bp, with gaps, in 13,206 scaffolds.  
50  
51 314 The N50 of scaffolds of 1 kbp or longer was 80,946 bp (L50 = 1,906), the average size of the  
52  
53 315 sequences was 38,118 bp. Using 20 kbp as an approximate value of longest plant gene length [59,  
54  
55 316 60], the percentage of scaffolds that equaled or surpassed this value in relation to the empirically  
56  
57 317 determined genome size is 83%, which corresponds to over 92% of the assembly total size. Contigs  
58  
59 318 generated by cutting scaffolds at each gap (of at least 25 base pair, i.e. 25 or more Ns) produced  
60  
61 319 N50 of 40,064 bp (L50 = 3,551) with average sequence size of 19,765 bp. The remaining gaps  
62  
63 320 comprised 26,447,057 bp (5.25% of the genome assembly) in 11,094 segments, with size of 2,384  
64  
65 321  $\pm 3,167$  bp. The total assembly size represents over 90% of the flow cytometry genome estimate



1  
2  
3  
4 322 (557 Mb) and should provide a good start to build a further improved reference genome assembly  
5  
6 323 of the species using long-range scaffolding techniques such as whole genome maps using either  
7  
8 324 imaging methods [61] or contact maps of chromosomes based on chromatin interactions [62].  
9  
10 325 Table 1 summarizes the main statistics of the *Handroanthus impetiginosus* genome assembly with  
11  
12 326 respect to the decisions made in the assembly process.

13 327

14  
15 328 A reassessment of the assembly accuracy was carried out using REAPR on the final genome  
16  
17 329 assembly. A total of 121 errors within a contig were still recognized, a much smaller number than  
18  
19 330 previously annotated (1,658 errors). Figure 2A shows the frequency distribution for the read  
20  
21 331 depth computed from the paired-end read alignment to the scaffolds sequences. It indicates the  
22  
23 332 expected effect on the distribution in comparison to the previous more redundant assembly. The  
24  
25 333 height of the heterozygous peak was successfully lowered by removing unmerged copies of the  
26  
27 334 same heterozygous loci. Figure 2B shows the relation between the observed number of scaffolds  
28  
29 335 in the final assembly and their read coverage in comparison to a Poisson approximation with  
30  
31 336  $\lambda = 63$  which was the observed average sequencing coverage for reads set from short  
32  
33 337 fragment libraries. Loss of information due to repeat sequences is clearly a limitation of this *H.*  
34  
35 338 *impetiginosus* assembly. Given the high rate of non classified consensus sequences we can infer  
36  
37 339 that most families/subfamilies of repeats might be underrepresented.

38 340

39 341 *[Insert Figure 2 here]*

40 342

41  
42 343 To complement the depth of read coverage analyses, we performed additional analyses to  
43  
44 344 identify the most probable causes of breaks in the assembly. We inspected contig termini defining  
45  
46 345 the positions of the terminal nucleotides of each contig from the genome assembly created by  
47  
48 346 cutting at each gap (of at least one base pair, i.e. one or more Ns). This analysis was developed  
49  
50 347 using a protocol described elsewhere [63] and results are summarized in Figure 3. Contig termini  
51  
52 348 overlap most prominently (~50%) with regions that do not encompass any annotated feature or  
53  
54 349 regions that have no depth of coverage (~15%) based on mapped reads to the assembly. It  
55  
56 350 suggests that contigs end in large repeats not yet resolved given the inherent limitations of short-  
57  
58 351 read sequence data. Another possibility is that these regions can contain low-copy young  
59  
60 352 euchromatic segmental duplication with higher sequence similarity to the consensus sequence.  
61  
62 353 Annotated interspersed repeats (~18%) and short tandem repeats (~9%) were the most  
63  
64  
65

1  
2  
3  
4 354 prominently annotated features with overlap to contigs ends. Less than 8% (2,473 of 31,668) of  
5  
6 355 annotated gene models were found to overlap contigs ends, indicating that very few are likely to  
7  
8 356 be interrupted in this unfinished assembly. It is a trend that was confirmed using BUSCO analysis  
9  
10 357 which reported only 3% of fragmented genes. Based on variant identification analysis with  
11  
12 358 FreeBayes using read data mapped to the genome assembly, we found virtually no allelic variants  
13  
14 359 located at contigs end, suggesting that interruption of continuity and contiguity in the assembly  
15  
16 360 is not related to differences between haplotypes.

17 361

18 362 *[Insert Figure 3 here]*

19  
20 363

21  
22 364 **Repetitive DNA.** A total of 1,608 consensus sequences (average length = 773 bp and totaling  
23  
24 365 1,281,536 bp) representing interspersed repeats in the genome assembly were found. Search for  
25  
26 366 domains in these sequences with similarity to known large families of genes that could confound  
27  
28 367 the identification of true repeats indicated 85 false positives in the consensus library of repeats.  
29  
30 368 Further 50 sequences were annotated with predicted protein domains frequently associated with  
31  
32 369 protein coding genes. These 135 sequences were wiped out from the consensus library. Most of  
33  
34 370 the remaining 1,473 sequences (71.1%) could not find classification in the hierarchical well-known  
35  
36 371 classes of Transposable Elements [64] but 16.6% could be classified as Class I (retrotransposons)  
37  
38 372 including three orders: LTR (12.8%), LINE (1.6%) and SINE (2.2%); 8.4% are Class II (DNA  
39  
40 373 transposons). Other categories comprised non-autonomous TEs: TRIM (0.4%) and MITE (3.5%).  
41  
42 374 Unknown non-classified sequences in the consensus library cover a wide range of sequence sizes  
43  
44 375 from 42 bp up to 5,987 bp (average = 345 bp, median = 503 bp). The 1,473 sequences representing  
45  
46 376 interspersed repeats in the consensus repeat library were used to mask the genome with  
47  
48 377 RepeatMasker. The masked fraction of the genome assembly comprised 155,348,349 bp, i.e.  
49  
50 378 30.9% of the total assembled genome of 503 Mbp. Remarkably, if we add to these ~155 Mbp the  
51  
52 379 54 Mbp of non-captured base pairs in the assembly when considering the empirically determined  
53  
54 380 genome size (= 557–503), the repetitive fraction of the genome approximates 37.5% (209 Mbp  
55  
56 381 out 557 Mbp). This is within the expected range (36.6% - 38.0%) for the repetitive fraction of the  
57  
58 382 genome estimated from the reads set using k-mer profiling approaches.

59 383

60 384 More than 50% of the masked bases in the assembly, or 80 Mbp, came from non-classified  
61  
62 385 sequences in the consensus library. In the well-known repeats, retrotransposons are the most

1  
2  
3  
4 386 abundant class in the assembly comprising 50 Mbp (~1/3 of the masked bases) with prominence  
5  
6 387 of LTR/Gypsy (~23 Mb) and LTR/Copy (18 Mb) families of repeats. DNA transposons and non-  
7  
8 388 autonomous orders of transposons masked 12 Mbp and 11 Mbp (~1/6 of the masked bases),  
9  
10 389 respectively, highlighting the prominence of DNA/hAT families of class II and MITE (Figure 4).  
11  
12 390 Simple sequence repeats (SSRs) detection using RepeatMasker identified a total of 182,115  
13  
14 391 microsatellites with a density of 2.76 kb per SSR in the genome assembly. This density  
15  
16 392 corroborates the general finding that the overall frequency of microsatellites is inversely related  
17  
18 393 to genome size in plant genomes [65]. This SSRs density in *H. impetiginosus* (genome size of 557  
19  
20 394 Mbp/SSR density of 362 per Mbp) is higher than in larger plant genomes such as those of maize  
21  
22 395 (1,115 Mbp/163 SSRs per Mbp), *S. bicolor* (738 Mbp/175 per Mbp), *G. raimondii* (761 Mbp/74.8  
23  
24 396 per Mbp) [66] but lower than densities in smaller genomes such as those of *A. thaliana* (120 Mbp/  
25  
26 397 418 per Mbp), *Medicago truncatula* (307 Mbp/ 495 per Mbp) and *C. sativus* (367 Mbp/ 552 per  
27  
28 398 Mbp) [67]. Different SSR motifs ranging from 1 to 6 bp showed that the di-nucleotide repeats  
29  
30 399 were the most abundant repeats followed by the mono- (Figure S4A). The frequency of SSR  
31  
32 400 decreased with increase in motif length (Supplementary Figure S4B), which is a trend usually  
33  
34 401 observed both in monocots and dicots [67].  
35  
36 402

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  
403 *[Insert Figure 4 here]*

404  
405 **Transcriptome assembly and gene content annotation and analysis.** A single run of Illumina  
406  
407 HiSeq 2500 sequencing, from a pool of RNA samples, generated nearly 148 million of paired end  
408  
409 reads. After adapter removal, trimming and coverage normalization, 55.2 million high-quality  
410  
411 reads (38%) were used to assemble the transcriptome using *de novo* (Trinity and SOAP-Trans-  
412  
413 denovo transcripts combined with the EvidentialGene pipeline) and genome guided methods  
414  
415 (PERTRAN). The PASA pipeline was used to integrate transcripts alignments to the genome  
416  
417 assembly from these set of sequences, generating 54,320 EST assemblies representing putative  
protein-coding loci in the genome assembly. Loci were identified by the assembled transcript  
alignments using BLASTX [36] and EXONERATE [37] alignments of plant peptides to the repeat-  
soft-masked genome using RepeatMasker. After gene model prediction and refinements, a total  
of 36,262 gene models were found in the genome assembly and 31,668 of them were retained  
after quality assessment based on Cscore, protein coverage, and overlap to repeats as described  
in Methods. The number of predicted mRNA transcripts was 35,479.

418

419 Structural features of the gene content are shown in Table 2 and Table 3. The average number of  
420 exons per gene was ~5 and its average length was 285 bp. The average number of introns per  
421 gene was ~4 and its average length was 445 bp. The GC content is significantly different between  
422 exons and introns (t-test p-value < 0.0001). Coding sequences have ~43% of GC, while introns  
423 have less with ~33% (Table 2). GC content tends to be higher in coding (exonic) than in non-coding  
424 regions [68], which may be related to gene architecture and alternative splicing [69-71]. A  
425 comparison of the gene features parameters, such as number and length (Figure S5A), was carried  
426 out between *H. impetiginosus* and *Erythranthe guttata*, another plant in the order Lamiales  
427 (Asterids), the model plant *A. thaliana* and the model tree *P. trichocarpa* (Rosids). As depicted in  
428 the frequency histograms, the exons parameters are stable among these species (Figure S5B). For  
429 the introns (Supplementary Figure S5C), frequency histograms have a sharp peak around 90 bp  
430 and a larger peak that is much lower in density. There is a small intron-size variability from species  
431 to species in the distributions, especially for larger introns, which rarely go beyond than 10,000  
432 bp. The intron length distributions in these four species is similar to those observed in lineages  
433 that are late in the evolutionary time scale, such as plants and vertebrates [72]. The sharp peak in  
434 the distributions at their “minimal intron” size is supposed to affect function by enhancing the  
435 rate at which mRNA is exported from the cell nucleus [73, 74]. In the model plant *A. thaliana*, a  
436 minimal intron group was previously defined [73] as anything that lies within three standard  
437 deviations of the optimum peak at 89±12 bp (53 bp – 125 bp). According to this definition, Table  
438 3 summarizes the distribution of the minimal intron among genes of *H. impetiginosus* and other  
439 selected plant species in the Asterids and Rosids lineages. We have calculated the percentages of  
440 minimal introns out of the total introns and the fraction of minimal-intron-containing genes with  
441 at least one minimal intron. Computed values were similar between *H. impetiginosus* and those  
442 of selected species with higher number of large introns (smaller minimal intron peak) but were  
443 more distinctive with those species such as *A. thaliana* and *E. guttata* in which the number of  
444 large introns was lower (larger minimal intron peak). This is thought as a general trend and was  
445 also observed in previous work [73]. These comparative analyses about the structural properties  
446 of the predicted genes indicate that the genome assembly of *H. impetiginosus* contains highly  
447 accurate gene structures.

448

449 *[Insert Figure 5 here]*

450

451 To further validate the gene content annotation, we used the transcript assemblies and selected  
452 plant proteomes to inspect if these sequences could align in its entirety to the genomic sequence.  
453 Out of the 31,668 primary mRNA transcripts (considering only the longest one when isoforms  
454 were predicted) in the genome, 11,488 have 100% of their CDS covered by EST assemblies. The  
455 remaining 20,054 transcripts have either a minimum of 80% of their CDS covered by EST  
456 assemblies or a cscore  $\geq 0.5$ . From these latter, the encoded putative peptides have excellent  
457 sequence similarity support from BLASTP comparisons with dicot species *Erythranthe guttata*  
458 (5,224 genes), *Sesamum indicum* (4,625 genes), potato or tomato (2,777 genes), soybean (1,484  
459 genes) and the poplar tree (1,424 genes) reflecting the taxonomic relationship between *H.*  
460 *impatiginosus* and these other related dicots. Gene models support was also found from more  
461 distantly related dicots (1,826 genes) and monocots (1,042 genes). Altogether, 31,048 gene  
462 models (98%) show well-supported similarity hits to other known plant protein sequences.  
463 Additional 517 predicted protein sequences did not produce hits and 103 sequences produced  
464 ambiguous hits from non-target species or represent possible contaminants in the assembly such  
465 as endophytic fungi (ascomycetes, 42 sequences; basidiomycetes, 17 sequences). Figure 5A  
466 summarizes the main finding regarding the similarity analyses with known proteins.

467

468 BUSCO [75] single-copy genes plant profiles were used to estimate completeness of the expected  
469 gene space as well as the duplicate fraction of the genome assembly. Out of the 956 profiles  
470 searched on the assembly, 59 (6.1%) were reported missing and 30 (3.1%) returned fragmented.  
471 From the profiles with complete match to the assembly, 867 (90.7%) were reported as single-copy  
472 and 247 (25.8%) were found completely duplicated. We benchmarked our results by searching  
473 the BUSCO profiles on the genomes of other lamids, *Erythranthe guttata* and *Olea europaea*. In  
474 *E. guttata* the analysis reported completeness level of 88% (848 single-copy profiles with  
475 complete match) while fragmented genes were 52 (5.4%). In *O. europaea*, the completeness level  
476 was 94% (905 complete single-copy profiles) and fragmented genes were only 14 (1.4%).  
477 Summary of BUSCO analysis is presented in Figure 5B.

478

479 Databases for gene ontology (GO) annotation are rich resources to describe functional properties  
480 of experimentally derived gene sets. To explore relationships between the GO terms in the *H.*  
481 *impatiginosus* and related, well-curated, genomes we used WEGO [76] to perform a genome-wide



1  
2  
3  
4 482 comparative analyses among broad functional GO terms with other lamids. The P-value of  
5  
6 483 Pearson Chi-Square test was considered to indicate significant relationships between the  
7  
8 484 proportions of genes of each GO term in these two datasets and to suggest patterns of enrichment  
9  
10 485 (Figures S6 and S7). These analyses revealed several GO terms in which the proportion of genes  
11  
12 486 in the two compared species were related. For the terms in which the comparison did not indicate  
13  
14 487 a significant relationship of gene proportions between the two datasets, the compared GO terms  
15  
16 488 suggested enrichments in *H. impetiginosus* for GO terms involved in metabolic processes and  
17  
18 489 catalytic activity in comparison to *E. guttata* and *O. europaea*.

19 490

20 491 [Insert Figure 6 here]

21  
22 492

23  
24 493 The central role of enzymes as biological catalysts is a well-studied issue related to the chemistry  
25  
26 494 of cells [77]. An important feature of most enzymes is that their activities can be regulated to  
27  
28 495 function properly to comply with physiological needs of the organism. We observed that GO term  
29  
30 496 for enzyme regulatory activity encompass a higher proportion of genes in *H. impetiginosus* than  
31  
32 497 in the two other lamids, albeit the difference did not reach significance in *E. guttata*. Research in  
33  
34 498 *Arabidopsis*, an herbaceous plant, has found little connectivity between metabolites and enzyme  
35  
36 499 activity [78]. In comparison to *Arabidopsis* broader GO terms, *H. impetiginosus* showed, as  
37  
38 500 discussed above, enrichment for the proportion genes assigned to metabolic process (49.1% >  
39  
40 501 47.4%; p-value 0.002) and catalytic activity (46.2% > 42.9%; p-value = 0). The proportion of genes  
41  
42 502 for enzyme regulatory activity was also higher in *H. impetiginosus* than *A. thaliana*, though not  
43  
44 503 statistically significant (p-value = 0.083). Investigations into whether and how metabolic process  
45  
46 504 and enzyme activities relate and how it could influence the known richness of metabolites for  
47  
48 505 forest trees of the mega diverse tropical biomes, particularly in the genus *Tabebuia* and  
49  
50 506 *Handroanthus*, shall be an interesting issue for future molecular and chemistry studies.

51  
52 507

53  
54 508 **Benchmarking the genome assembly of *H. impetiginosus*.** Based on current standards for plant  
55  
56 509 genome sequence assembly [60, 79, 80] we have provided a quality assembly of high future utility.  
57  
58 510 To support functional analyses we classified the gene models into high-confidence and low-  
59  
60 511 confidence groups. Out of the 31,688 protein-coding loci annotated in the genome assembly,  
61  
62 512 28,603 (90%) produced high-confidence gene models (Supplementary File S1). This subset  
63  
64 513 contains approximately the same number of genes reported in less fragmented genome  
65

1  
2  
3  
4 514 assemblies for other lamids. *E. guttata* (2n=28) reports 28,140 protein-coding genes [81]; *O.*  
5 515 *europaea* (2n = 46) has 56,349 protein-coding genes [82] but its genome has likely undergone a  
6 516 whole genome duplication event. Most of *Tabebuia* and *Handroanthus* species studied so far have  
7  
8 517 2n = 40 [22]. The fraction of gene duplicates in the BUSCO analysis (see Figure 5B) was intended  
9  
10 518 to estimate the level of redundancy in the genome assembly. We benchmarked our results by  
11  
12 519 searching the completed duplicated BUSCO profiles on the genomes of *E. guttata* and *O.*  
13 520 *europaea*. In the first, it was found to be 15% (150 out of 956), while in the latter the duplicated  
14  
15 521 profiles were 38% (364 out of 956). In these three lamids, we can infer that the frequency of small-  
16  
17 522 and large-scale duplications, such as (paleo)polyploidy, can explain the differences in the number  
18  
19 523 of annotated genes and levels of gene duplication (*E. guttata* <= *H. impetiginosus* << *O. europaea*).  
20  
21 524 It suggests that the *H. impetiginosus* genome has not undergone a recent whole-genome  
22  
23 525 duplication event, although a deeper analysis of this question was beyond the scope of this study.  
24  
25 526

26  
27 527 Our genome assembly metrics were benchmarked against comparable genome assemblies of  
28  
29 528 other highly heterozygous forest tree genomes (File S2 and Figure S8). The *H. impetiginosus*  
30  
31 529 assembly has 503 Mbp in 13,206 scaffolds  $\geq 2$  kbp, representing over 90% of the flow cytometry  
32  
33 530 estimated size (557 Mb). For *Quercus robur*, the assembly had 17,910 scaffolds  $\geq 2$  kbp with  
34  
35 531 scaffolds N50 of 260 kbp, but corresponding to 1.34 Gbp, i.e. 81% larger than the expected 740  
36  
37 532 Mbp genome, which is clearly undesirable [83]. For *Quercus lobata* with a genome size of 730  
38  
39 533 Mbp two assemblies were provided: a haplotype-reduced assembly, with 40,158 contigs totaling  
40  
41 534 760 Mb, N50 of 95 kbp and a more complete version for gene models, containing 94,394 scaffolds  
42  
43 535  $\geq 2$  kbp, totaling 1.15 Gbp, with an N50 of 278 kbp [48]. Despite our lower NG50/N50 scaffold  
44  
45 536 length <100 kbp, the *H. impetiginosus* assembly has a large (60%) percentage of scaffolds  $\geq 20$   
46  
47 537 kbp. This value is higher than the reported values for *Quercus lobata* v0.5 (53%), *Quercus lobata*  
48  
49 538 v1.0 (51%) and *Quercus rubra* (48%), even if those assemblies had higher NG50/N50 scaffold  
50  
51 539 lengths. Finally, contigs termini analysis has found virtually no allelic variants located at contigs  
52  
53 540 ends, suggesting that interruption of continuity and contiguity in the assembly is not related to  
54  
55 541 differences between haplotypes. This genome assembly for *Handroanthus impetiginosus* will thus  
56  
57 542 be useful for variant calling, one of the main future objectives for generating this resource.  
58  
59 543

60 544 **Genome-guided exploration of specialized metabolism genes of quinoid systems.** Aside from its  
61  
62 545 high valued wood, *H. impetiginosus* and other Ipê species are also known for their medicinal  
63  
64  
65

1  
2  
3  
4 546 effects. Extracts from its bark and wood have many ethnobotanical uses: against cancer, malaria,  
5  
6 547 fevers, trypanosomiasis, fungal and bacterial infections and stomach disorders [84, 85]. The wood  
7  
8 548 extracts have also been demonstrated to have anti-inflammatory effects [86] [87]. The main  
9  
10 549 bioactive components isolated from the Pink Ipê is Lapachol and its products [88], which are  
11  
12 550 naphthoquinones derived from the o-succinylbenzoate (OSB) pathway [89]. Lapachol is also  
13  
14 551 responsible for the well-known high resistance of the Ipê wood against rotting fungi and insects  
15  
16 552 [90]. In addition, naphthoquinones are aromatic substances with ecological importance for the  
17  
18 553 interaction of plants with other plants, insects and microbes [89]. Given their medicinal and  
19  
20 554 biological relevance, we have searched the *H. impetiginosus* annotated genes for the enzymes  
21  
22 555 involved in the biosynthesis of naphthoquinones. . By searching for the KEGG identifiers of these  
23  
24 556 enzymes (e.g. K01851) in the InterPro annotation results, we found all the important known  
25  
26 557 enzymes that lead to the biosynthesis of lapachol (Figure 6). Unfortunately, however, the last two  
27  
28 558 steps of the lapachol biosynthesis pathway still constitute unidentified enzymes [89]. For  
29  
30 559 comparative purposes, we downloaded the annotation file of five other species from the  
31  
32 560 Phytozome database. The number of *H. impetiginosus* genes encoding for the enzymes of each  
33  
34 561 step in the pathway is comparable to the numbers found in other species. However, three  
35  
36 562 exceptions were found. *H. impetiginosus* has five genes encoding the enzyme that converts  
37  
38 563 chorismate to isochorismate, the first step in the o-succinylbenzoate (OSB) pathway. Two other  
39  
40 564 steps where *H. impetiginosus* were found to have relatively more genes are the ones that lead to  
41  
42 565 the synthesis of 1,4-Dihydroxy-2-naphtoyl-CoA and of 2-Phytyl-1,4-naphthoquinone. The  
43  
44 566 availability of sequences for these genes may open new avenues for biotechnological products  
45  
46 567 and for a better understanding of their ecological roles.  
47

568

#### 569 **RE-USE POTENTIAL**

570 We have reported an unfinished genome assembly for *Handroanthus impetiginosus*, a highly  
571  
572 571 valued, ecologically keystone tropical timber and a species rich in natural products. The  
573  
574 572 fragmentation of this preliminary assembly might be still be limiting for deeper insights of whole-  
575  
576 573 genome comparative analyses or studies of genome evolution [91], although we think that such  
577  
578 574 studies may be carried out using this assembly at least at the gene-level of gene-family level.  
579  
580 575 Nevertheless, the broad validation performed provides a useful genomic resource for genetic and  
581  
582 576 functional analysis including, but not limited to, downstream applications such as variant calling,  
583  
584 577 molecular markers development and functional studies. Extensive documentation of quality  
585

1  
2  
3  
4 578 throughout the assembly process was provided showing that acceptable continuity was reached  
5  
6 579 and that the fragmentation of the final sequence mostly derived from loss of information on high-  
7  
8 580 copy families of long interspersed repeats or the presence of low-copy segmental duplications  
9  
10 581 likely recently evolved with higher sequence similarity to the consensus sequence. Certainly, there  
11  
12 582 are still inaccuracies at the base and assembly level but all efforts were made to deliver results to  
13  
14 583 end user with the appropriate documentation, making this initial read set, sequence and  
15  
16 584 annotations as a primary and reliable starting grounds for further improvement.  
17

585

18 586 We have documented in detail the main features of the reported assembly. The total assembly  
19  
20 587 size of scaffolds with  $\geq 2$  kbp in length is 90% of the flow cytometry determined genome size, we  
21  
22 588 believe a remarkable accomplishment given the anticipated difficulties in assembling such a  
23  
24 589 repetitive and highly heterozygous diploid genome based exclusively on short-read sequencing.  
25  
26 590 The percentage of base pairs in scaffolds with  $\geq 20$  kbp is 83% (461 Mbp of 557 Mbp) of the  
27  
28 591 empirically determined genome size, which corresponds to 92% of the assembled total size (461  
29  
30 592 Mbp of 503 Mbp). Using 20 kbp as an approximate value of longest plant gene length, this result  
31  
32 593 shows that 60% of the assembly is accessible for reliable gene annotation. Furthermore the  
33  
34 594 N50/NG50 (41 kbp/34 kbp) contig length is longer than 30 kbp, which has been suggested to be  
35  
36 595 an adequate minimum threshold for high utility of a genome assembly [79]. The percentage of  
37  
38 596 documented gaps in scaffolds is only 5.3% and the few misassembled signatures present in the  
39  
40 597 assembly were fully documented based on acceptable metrics such as fragment coverage  
41  
42 598 distribution error (FCD error). Less than 8% (2,473 of 31,668) of annotated gene models were  
43  
44 599 found to overlap contigs ends, indicating that very few are likely to be interrupted in this  
45  
46 600 unfinished assembly. No allelic variants were found at contigs ends, suggesting that interruption  
47  
48 601 of continuity and contiguity in the assembly is not related to differences between haplotypes,  
49  
50 602 therefore providing a valuable resource for variant calling and functional analysis. Over 86%  
51  
52 603 (27,380 of 31,668) of the gene models represented in the assembly have external evidential  
53  
54 604 support measured by Pasa-validated EST alignments from RNA-Seq or high-coverage alignments  
55  
56 605 with known plant proteins (>90% coverage). Furthermore, 80% (25,369 of 31,668) of transcripts  
57  
58 606 have conceptual translation that contain protein domain annotation, excluded those associated  
59  
60 607 to TEs. Finally, a summary of BUSCO analysis indicates that the detected number of plant single  
61  
62 608 copy orthologs represents 90% of the searched profiles (867 of 956) while only 6% is missing and  
63  
64 609 3% is fragmented.  
65

610

611 This is the first well-curated genome for a Neotropical forest tree and the first one reported for a  
612 member of the Bignoniaceae family. Besides expanding the opportunities for comparative  
613 genomic studies by including an overlooked taxonomic family, the availability of this genome  
614 assembly will foster functional studies with new targets and allow the development and  
615 application of robust and far-reaching sets of genome-wide SNP genotyping tools to support  
616 multiple population genomics analyses in *H. impetiginosus* and related species of the Tabebuia  
617 Alliance. This group includes several of the most ecologically and economically important timber  
618 species of the American tropics. Going beyond the species-specific significance of these results,  
619 this study paves the way for developing similar genomic resources for other Neotropical forest  
620 trees of equivalent relevance. This in turn will open exceptional prospects to empower a higher-  
621 level understanding of the evolutionary history, species distribution and population demography  
622 of the still largely neglected forest trees of the mega diverse tropical biomes. Furthermore, this  
623 genome assembly provides a new resource for advances in the current integration between  
624 genomics, transcriptomics and metabolomics approaches for exploration of the enormous  
625 structural diversity and biological activities of plant-derived compounds.

626

#### 627 **AVAILABILITY OF SUPPORTING DATA**

628

629 Sequences for the genome and assembly along with gene content annotation as well as the raw  
630 sequencing reads have been deposited into GenBank, BioProject PRJNA324125. This Whole  
631 Genome Shotgun (WGS) project has been deposited at DDBJ/ENA/GenBank under the accession  
632 NKXS00000000. The version described in this paper is version NKXS01000000. BioSample for WGS  
633 is SAMN05195323 and corresponding SRA run accessions are SRR3624821 - SRR3624825.  
634 BioSample for RNA-Seq is SAMN07346903 with SRA run accession SRR5820886. Perl script that  
635 automated the read set from mate-pair sequencing preprocessing (TrimAdaptor.pl) was uploaded  
636 to GigaDB under permission of the original authors at the High-Throughput Sequencing and  
637 Genotyping Center Unit of the University of Illinois Urbana-Champaign. Summary outputs for  
638 main analysis in this research were made available also through GigaDB.

639

#### 640 **List of abbreviations**

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  
4 641 BLASTP, Basic Local Alignment Search Tool for Proteins; BLAT, BLAST-like alignment tool; CDS,  
5  
6 642 coding DNA sequence; EC, Enzyme Commission Number; EST, Expressed Sequence Tag; GATK,  
7  
8 643 Genome Analysis Toolkit; GO, Gene Ontology; LINE, Long Interspersed Nuclear Elements; LTR,  
9  
10 644 Long Terminal Repeats; MBH, Mutual Best Hit; MITE, Miniature Inverted-Repeat Transposable  
11  
12 645 Elements; mRNA, messenger RNA; PASA, Program to Assemble Spliced Alignment; REAPR,  
13  
14 646 Recognition of Errors in Assemblies using Paired Reads; SINE, Short Interspersed Nuclear  
15  
16 647 Elements; SNP, Single Nucleotide Polymorphism; SSPACE, SSAKE-based Scaffolding of Pre-  
17  
18 648 Assembled Contigs after Extension; TE, transposable element.

19 649

20 650 **Ethics approval**

21  
22 651 Not applicable

23  
24 652

25 653 **Consent for publication**

26  
27 654 Not applicable

28  
29 655

30  
31 656 **Competing interests**

32  
33 657 The authors declare that they have no competing interests.

34  
35 658

36 659 **Funding and acknowledgements**

37  
38 660 This work was supported by competitive grants from CNPq to RGC (project no. 471366/2007-2  
39  
40 661 and Rede Cerrado CNPq/PPBio project no. 457406/2012-7), to EN (CNPq Proc. 476709/2012-1)  
41  
42 662 and to DG (PRONEX FAP-DF Project Grant "NEXTREE" 193.000.570/2009). RGC and DG have been  
43  
44 663 supported by productivity grants from CNPq, which we gratefully acknowledge. OBSJr has been  
45  
46 664 supported by an EMBRAPA doctoral fellowship and was an Affiliate Researcher at Lawrence  
47  
48 665 Berkeley National Laboratory (LBNL), Berkeley CA, at the time of this research. OBSJr thanks to  
49  
50 666 DM Goodstein and the members of the Phytozome team at the LBNL/Joint Genome Institute (JGI)  
51  
52 667 for their valuable help and support in working with the JGI pipelines for genomic research. WE  
53  
54 668 also thank Dr. Gabriela Ferreira Nogueira and André Luis X. de Souza for their help with flow  
55  
56 669 cytometry analysis.

57  
58 670

59 671 **Authors' contributions**

60  
61  
62  
63  
64  
65

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

677

678

## REFERENCES

679

680

681

682

683

684

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

700

701

702

703

704

705

706

707

708

709

710

711

712

713

714

715

1. Goodstein DM, Shu SQ, Howson R, Neupane R, Hayes RD, Fazo J, Mitros T, Dirks W, Hellsten U, Putnam N *et al*: **Phytozome: a comparative platform for green plant genomics**. *Nucleic Acids Research* 2012, **40**(D1):D1178-D1186.
2. Kang YJ, Lee T, Lee J, Shim S, Jeong H, Satyawani D, Kim MY, Lee SH: **Translational genomics for plant breeding with the genome sequence explosion**. *Plant Biotechnology Journal* 2016, **14**(4):1057-1069.
3. Bevan M, Walsh S: **The Arabidopsis genome: A foundation for plant research**. *Genome Research* 2005, **15**(12):1632-1642.
4. Morrell PL, Buckler ES, Ross-Ibarra J: **Crop genomics: advances and applications**. *Nature Reviews Genetics* 2012, **13**(2):85-96.
5. Varshney RK, Glaszmann JC, Leung H, Ribaut JM: **More genomic resources for less-studied crops**. *Trends in Biotechnology* 2010, **28**(9):452-460.
6. Myburg AA, Grattapaglia D, Tuskan GA, Hellsten U, Hayes RD, Grimwood J, Jenkins J, Lindquist E, Tice H, Bauer D *et al*: **The genome of *Eucalyptus grandis***. *Nature* 2014, **510**(7505):356-362.
7. Tuskan GA, DiFazio S, Jansson S, Bohlmann J, Grigoriev I, Hellsten U, Putnam N, Ralph S, Rombauts S, Salamov A *et al*: **The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray)**. *Science* 2006, **313**(5793):1596-1604.
8. Neale DB, Wegrzyn JL, Stevens KA, Zimin AV, Puiu D, Crepeau MW, Cardeno C, Koriabine M, Holtz-Morris AE, Liechty JD *et al*: **Decoding the massive genome of loblolly pine using haploid DNA and novel assembly strategies**. *Genome Biology* 2014, **15**(3).
9. Nystedt B, Street NR, Wetterbom A, Zuccolo A, Lin YC, Scofield DG, Vezzi F, Delhomme N, Giacomello S, Alexeyenko A *et al*: **The Norway spruce genome sequence and conifer genome evolution**. *Nature* 2013, **497**(7451):579-584.
10. Moghe G, Last R: **Something old, something new: Conserved enzymes and the evolution of novelty in plant specialized metabolism**. *Plant Physiology* 2015:pp.00994.02015.
11. Stone R: **Lifting the Veil on Traditional Chinese Medicine**. *Science* 2008, **319**(5864):709-710.
12. Chappell J, DellaPenna D, O'Connor S: **Specific Aims for Medicinal Plant Genomics Resource**. *Medicinal Plants Genomics Resource* 2017.
13. Brousseau L, Tinaut A, Duret C, Lang T, Garnier-Gere P, Scotti I: **High-throughput transcriptome sequencing and preliminary functional analysis in four Neotropical tree species**. *BMC Genomics* 2014, **15**(1):238.
14. Olsson S, Seoane-Zonjic P, Bautista Ro, Claros G, González-Martínez S, Scotti I, Scotti-Saintagne C, Hardy O, Heuertz M: **Development of genomic tools in a widespread**

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  
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

- 716 **tropical tree, *Symphonia globulifera* L.f.: a new low-coverage draft genome, SNP and**  
717 **SSR markers. *Molecular Ecology Resources* 2017, 17(4):614-630.**
- 718 15. Cadena-González A, Sorensen M, Theilade I: **Use and valuation of native and**  
719 **introduced medicinal plant species in Campo Hermoso and Zetaquirá, Boyacá,**  
720 **Colombia. *Journal of Ethnobiology and Ethnomedicine* 2013, 9(1):23.**
- 721 16. Bodker G, Bhat KKS, Burley J, Vantomme P: **Medicinal plants for forest conservation**  
722 **and health care. *Food and Agriculture Organization of the United Nations* 1997.**
- 723 17. Braga AC, Reis AMM, Leoi LT, Pereira RW, Collevatti RG: **Development and**  
724 **characterization of microsatellite markers for the tropical tree species *Tabebuia aurea***  
725 **(*Bignoniaceae*). *Molecular Ecology Notes* 2007, 7(1):53-56.**
- 726 18. Liu B, Shi Y, Yuan J, Hu X, Zhang H, Li N, Li Z, Chen Y, Mu D, Fan W: **Estimation of**  
727 **genomic characteristics by analyzing k-mer frequency in de novo genome projects.**  
728 *arXiv:13082012* 2013.
- 729 19. Schulze M, Grogan J, Uhl C, Lentini M, Vidal E: **Evaluating ipe (*Tabebuia*, *Bignoniaceae*)**  
730 **logging in Amazonia: Sustainable management or catalyst for forest degradation?**  
731 *Biological Conservation* 2008, 141(8):2071-2085.
- 732 20. Inagaki R, Ninomiya M, Tanaka K, Watanabe K, Koketsu M: **Synthesis and Cytotoxicity**  
733 **on Human Leukemia Cells of Furonaphthoquinones Isolated from *Tabebuia* Plants.**  
734 *Chemical & Pharmaceutical Bulletin* 2013, 61(6):670-673.
- 735 21. Park BS, Kim JR, Lee SE, Kim KS, Takeoka GR, Ahn YJ, Kim JH: **Selective growth-inhibiting**  
736 **effects of compounds identified in *Tabebuia impetiginosa* inner bark on human**  
737 **intestinal bacteria. *Journal of Agricultural and Food Chemistry* 2005, 53(4):1152-1157.**
- 738 22. Collevatti RG, Dornelas MC: **Clues to the evolution of genome size and chromosome**  
739 **number in *Tabebuia* alliance (*Bignoniaceae*). *Plant Systematics and Evolution* 2016,**  
740 **302(5):601-607.**
- 741 23. Aronesty E: **Comparison of sequencing utility programs. *The Open Bioinformatics***  
742 *Journal* 2013, 7:1-8.
- 743 24. Langmead B, Trapnell C, Pop M, Salzberg SL: **Ultrafast and memory-efficient alignment**  
744 **of short DNA sequences to the human genome. *Genome Biology* 2009, 10(3).**
- 745 25. Marçais G, Kingsford C: **A fast, lock-free approach for efficient parallel counting of**  
746 **occurrences of k-mers. *Bioinformatics* 2011, 27(6):764-770.**
- 747 26. Vurture GW, Sedlazeck FJ, Nattestad M, Underwood CJ, Fang H, Gurtowski J, Schatz MC: **GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics***  
748 *btx153* 2017.
- 749 27. Gnerre S, MacCallum I, Przybylski D, Ribeiro F, Burton J, Walker B, Sharpe T, Hall G, Shea  
750 T, Sykes S *et al*: **High-quality draft assemblies of mammalian genomes from massively**  
751 **parallel sequence data. *Proceedings of the National Academy of Sciences* 2011,**  
752 **108(4):1513-1518.**
- 753 28. Flutre T, Duprat E, Feuillet C, Quesneville H: **Considering Transposable Element**  
754 **Diversification in De Novo Annotation Approaches. *Plos One* 2011, 6(1).**
- 755 29. Wicker T, Sabot F, Hua-Van A, Bennetzen JL, Capy P, Chalhoub B, Flavell A, Leroy P,  
756 Morgante M, Panaud O *et al*: **A unified classification system for eukaryotic**  
757 **transposable elements. *Nature Reviews Genetics* 2007, 8(12):973-982.**
- 758 30. Hoede C, Arnoux S, Moisset M, Chaumier T, Inizan O, Jamilloux V, Quesneville H:  
759 **PASTEC: An Automatic Transposable Element Classification Tool. *Plos One* 2014, 9(5).**
- 760 31. Smit AFA, Hubley R, Green P: **RepeatMasker Open-4.0 (2013-2015).** 2015.
- 761 32. Bolger AM, Lohse M, Usadel B: **Trimmomatic: a flexible trimmer for Illumina sequence**  
762 **data. *Bioinformatics* 2014, 30(15):2114-2120.**
- 763



- 1  
2  
3  
4 764 33. Xie Y, Wu G, Tang J, Luo R, Patterson J, Liu S, Huang W, He G, Gu S, Li S *et al*:  
5 765 **SOAPdenovo-Trans: de novo transcriptome assembly with short RNA-Seq reads.**  
6 766 *Bioinformatics* 2014, **30**(12):1660-1666.  
7 767 34. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, Adiconis X, Fan L,  
8 768 Raychowdhury R, Zeng QD *et al*: **Full-length transcriptome assembly from RNA-Seq**  
9 769 **data without a reference genome.** *Nature Biotechnology* 2011, **29**(7):644-U130.  
10 770 35. Gilbert D: **EvidentialGene: mRNA Transcript Assembly Software.** *EvidentialGene :*  
11 771 *Evidence Directed Gene Construction for Eukaryotes* 2013.  
12 772 36. Schmutz J, McClean P, Mamidi S, Wu A, Cannon S, Grimwood J, Jenkins J, Shu S, Song Q,  
13 773 Chavarro C *et al*: **A reference genome for common bean and genome-wide analysis of**  
14 774 **dual domestications.** *Nature Genetics* 2014, **46**(7):707-713.  
15 775 37. Shu S, Goodstein DM, Hayes D, Mitros T, Rokhsar D: **JGI Plant Genomics Gene**  
16 776 **Annotation Pipeline.** *SciTech Connect* 2017.  
17 777 38. Gish W, States D: **Identification of protein coding regions by database similarity search.**  
18 778 *Nature Genetics* 1993, **3**(3):266-272.  
19 779 39. Slater GS, Birney E: **Automated generation of heuristics for biological sequence**  
20 780 **comparison.** *Bmc Bioinformatics* 2005, **6**.  
21 781 40. **RepeatMasker Open-4.0.** <http://www.repeatmasker.org>  
22 782 [<http://www.repeatmasker.org>]  
23 783 41. UniProt Consortium: **UniProt: a hub for protein information.** *Nucleic Acids Research*  
24 784 2014, **43**(D1):D204-D212.  
25 785 42. Salamov AA, Solovyev VV: **Ab initio gene finding in Drosophila genomic DNA.** *Genome*  
26 786 *Research* 2000, **10**(4):516-522.  
27 787 43. Solovyev V, Kosarev P, Seledsov I, Vorobyev D: **Automatic annotation of eukaryotic**  
28 788 **genes, pseudogenes and promoters.** *Genome Biology* 2006, **7**.  
29 789 44. Yeh RF, Lim LP, Burge CB: **Computational inference of homologous gene structures in**  
30 790 **the human genome.** *Genome Research* 2001, **11**(5):803-816.  
31 791 45. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK, Hannick LI, Maiti R, Ronning CM,  
32 792 Rusch DB, Town CD *et al*: **Improving the Arabidopsis genome annotation using**  
33 793 **maximal transcript alignment assemblies.** *Nucleic Acids Research* 2003, **31**(19):5654-  
34 794 5666.  
35 795 46. Jones P, Binns D, Chang HY, Fraser M, Li W, McAnulla C, McWilliam H, Maslen J, Mitchell  
36 796 A, Nuka G *et al*: **InterProScan 5: genome-scale protein function classification.**  
37 797 *Bioinformatics* 2014, **30**(9):1236-1240.  
38 798 47. Braga AC, Collevatti RG: **Temporal variation in pollen dispersal and breeding structure**  
39 799 **in a bee-pollinated Neotropical tree.** *Heredity* 2011, **106**(6):911-919.  
40 800 48. Sork VL, Fitz-Gibbon ST, Puiu D, Crepeau M, Gugger PF, Sherman R, Stevens K, Langley  
41 801 CH, Pellegrini M, Salzberg SL: **First Draft Assembly and Annotation of the Genome of a**  
42 802 **California Endemic Oak Quercus lobata Nee (Fagaceae).** *G3-Genes Genomes Genetics*  
43 803 2016, **6**(11):3485-3495.  
44 804 49. Luo RB, Liu BH, Xie YL, Li ZY, Huang WH, Yuan JY, He GZ, Chen YX, Pan Q, Liu YJ *et al*:  
45 805 **SOAPdenovo2: an empirically improved memory-efficient short-read de novo**  
46 806 **assembler.** *Gigascience* 2012, **1**.  
47 807 50. Zimin AV, Marçais G, Puiu D, Roberts M, Salzberg SL, Yorke JA: **The MaSuRCA genome**  
48 808 **assembler.** *Bioinformatics* 2013, **29**(21):2669-2677.  
49 809 51. Kajitani R, Toshimoto K, Noguchi H, Toyoda A, Ogura Y, Okuno M, Yabana M, Harada M,  
50 810 Nagayasu E, Maruyama H *et al*: **Efficient de novo assembly of highly heterozygous**

- 1  
2  
3  
4 811 **genomes from whole-genome shotgun short reads.** *Genome Research* 2014,  
5 812 **24(8):1384-1395.**
- 6 813 52. Malinsky M, Simpson JT, Durbin R: **Trio-sga: facilitating de novo assembly of highly**  
7 814 **heterozygous genomes with parent-child trios.** *bioRxiv* 2016.
- 8 815 53. Boetzer M, Henkel CV, Jansen HJ, Butler D, Pirovano W: **Scaffolding pre-assembled**  
9 816 **contigs using SSPACE.** *Bioinformatics* 2011, **27(4):578-579.**
- 10 817 54. Nadalin F, Vezzi F, Policriti A: **GapFiller: a de novo assembly approach to fill the gap**  
11 818 **within paired reads.** *Bmc Bioinformatics* 2012, **13.**
- 12 819 55. Hunt M, Kikuchi T, Sanders M, Newbold C, Berriman M, Otto TD: **REAPR: a universal**  
13 820 **tool for genome assembly evaluation.** *Genome Biology* 2013, **14(5):R47.**
- 14 821 56. Pongstingl H, Ning ZM: **SMALT. 2010 - 2015 Genome Research Ltd** 2016.
- 15 822 57. Kent WJ: **BLAT--the BLAST-like alignment tool.** *Genome research* 2002, **12(4):656-664.**
- 16 823 58. Stanke M, Keller O, Gunduz I, Hayes A, Waack S, Morgenstern B: **AUGUSTUS: ab initio**  
17 824 **prediction of alternative transcripts.** *Nucleic Acids Research* 2006, **34:W435-W439.**
- 18 825 59. Ramírez-Sánchez O, Pérez-Rodríguez P, Delaye L, Tiessen A: **Plant Proteins Are Smaller**  
19 826 **Because They Are Encoded by Fewer Exons than Animal Proteins.** *Genomics,*  
20 827 *Proteomics & Bioinformatics* 2016, **14(6):357-370.**
- 21 828 60. Bradnam KR, Fass JN, Alexandrov A, Baranay P, Bechner M, Birol I, Boisvert S, Chapman  
22 829 JA, Chapuis G, Chikhi R *et al*: **Assemblathon 2: evaluating de novo methods of genome**  
23 830 **assembly in three vertebrate species.** *GigaScience* 2013, **2:10-10.**
- 24 831 61. Lam ET, Hastie A, Lin C, Ehrlich D, Das SK, Austin MD, Deshpande P, Cao H, Nagarajan N,  
25 832 Xiao M *et al*: **Genome mapping on nanochannel arrays for structural variation analysis**  
26 833 **and sequence assembly.** *Nature Biotechnology* 2012, **30(8):771-776.**
- 27 834 62. Ay F, Noble WS: **Analysis methods for studying the 3D architecture of the genome.**  
28 835 *Genome Biology* 2015, **16.**
- 29 836 63. Tørresen OK, Star B, Jentoft S, Reinart WB, Grove H, Miller JR, Walenz BP, Knight J,  
30 837 Ekholm JM, Peluso P *et al*: **An improved genome assembly uncovers prolific tandem**  
31 838 **repeats in Atlantic cod.** *BMC Genomics* 2017, **18(1):95.**
- 32 839 64. Wicker T, Sabot F, Hua-Van A, Bennetzen J, Capy P, Chalhoub B, Flavell A, Leroy P,  
33 840 Morgante M, Panaud O *et al*: **A unified classification system for eukaryotic**  
34 841 **transposable elements.** *Nature Reviews Genetics* 2007, **8(12):973-982.**
- 35 842 65. Morgante M, Hanafey M, Powell W: **Microsatellites are preferentially associated with**  
36 843 **nonrepetitive DNA in plant genomes.** *Nature Genetics* 2002, **30(2):194-200.**
- 37 844 66. Wang Q, Fang L, Chen J, Hu Y, Si Z, Wang S, Chang L, Guo W, Zhang T: **Genome-Wide**  
38 845 **Mining, Characterization, and Development of Microsatellite Markers in Gossypium**  
39 846 **Species.** *Scientific Reports* 2015, **5(1).**
- 40 847 67. Sonah H, Deshmukh R, Sharma A, Singh V, Gupta D, Gacche R, Rana J, Singh N, Sharma T:  
41 848 **Genome-Wide Distribution and Organization of Microsatellites in Plants: An Insight**  
42 849 **into Marker Development in Brachypodium.** *PLOS ONE* 2011, **6(6):e21298.**
- 43 850 68. Bernardi G: **Isochores and the evolutionary genomics of vertebrates.** *Gene* 2000,  
44 851 **241(1):3-17.**
- 45 852 69. Mizuno M, Kanehisa M: **Distribution profiles of GC content around the translation**  
46 853 **initiation site in different species.** *FEBS letters* 1994, **352(1):7-10.**
- 47 854 70. Amit M, Donyo M, Hollander D, Goren A, Kim E, Gelfman S, Lev-Maor G, Burstein D,  
48 855 Schwartz S, Postolsky B *et al*: **Differential GC content between exons and introns**  
49 856 **establishes distinct strategies of splice-site recognition.** *Cell reports* 2012, **1(5):543-556.**
- 50  
51  
52  
53  
54  
55  
56  
57  
58  
59  
60  
61  
62  
63  
64  
65

- 1  
2  
3  
4 857 71. Wendel JF, Greilhuber J, Dolezel J, Leitch IJ: **Plant Genome Diversity Volume 1 - Plant**  
5 858 **Genomes, their Residents, and their Evolutionary Dynamics**, vol. 1. Wien: Springer-  
6 859 Verlag; 2012.
- 8 860 72. JiaYan W, JingFa X, LingPing W, Jun Z, HongYan Y, ShuangXiu W, Zhang Z, Jun Y:  
9 861 **Systematic analysis of intron size and abundance parameters in diverse lineages.**  
10 862 *Science China Life Sciences* 2013, **56**(10):968-974.
- 11 863 73. Yu J, Yang Z, Kibukawa M, Paddock M, Passey D, Wong G: **Minimal Introns Are Not**  
12 864 **“Junk”**. *Genome Research* 2002, **12**(8):1185-1189.
- 14 865 74. Zhu J, He F, Wang D, Liu K, Huang D, Xiao J, Wu J, Hu S, Yu J: **A Novel Role for Minimal**  
15 866 **Introns: Routing mRNAs to the Cytosol**. *PLOS ONE* 2010, **5**(4):e10144.
- 16 867 75. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM: **BUSCO: assessing**  
17 868 **genome assembly and annotation completeness with single-copy orthologs.**  
18 869 *Bioinformatics* 2015, **31**(19):3210-3212.
- 20 870 76. Ye J, Fang L, Zheng H, Zhang Y, Chen J, Zhang Z, Wang J, Li S, Li R, Bolund L *et al*: **WEGO:**  
21 871 **a web tool for plotting GO annotations.** *Nucleic Acids Research* 2006, **34**(Web Server  
22 872 issue):W293-W297.
- 23 873 77. Cooper GM: **The Cell, 2nd edition. A Molecular Approach**. Sunderland (MA): Sinauer  
24 874 Associates; 2000.
- 25 875 78. Sulpice R, Trenkamp S, Steinfath M, Usadel B, Gibon Y, Witucka-Wall H, Pyl E-T, Tschoep  
26 876 H, Steinhauser MC, Guenther M *et al*: **Network Analysis of Enzyme Activities and**  
27 877 **Metabolite Levels and Their Relationship to Biomass in a Large Panel of**  
28 878 **<em>Arabidopsis</em> Accessions.** *The Plant Cell* 2010, **22**(8):2872-2893.
- 30 879 79. Hamilton JP, Robin Buell C: **Advances in plant genome sequencing.** *The Plant Journal*  
31 880 2012, **70**(1):177-190.
- 32 881 80. Barthelson R, McFarlin AJ, Rounsley SD, Young S: **Plantagora: Modeling Whole Genome**  
33 882 **Sequencing and Assembly of Plant Genomes.** *PLOS ONE* 2011, **6**(12):e28436.
- 35 883 81. Hellsten U, Wright KM, Jenkins J, Shu S, Yuan Y, Wessler SR, Schmutz J, Willis JH, Rokhsar  
36 884 DS: **Fine-scale variation in meiotic recombination in Mimulus inferred from population**  
37 885 **shotgun sequencing.** *Proceedings of the National Academy of Sciences of the United*  
38 886 *States of America* 2013, **110**(48):19478-19482.
- 40 887 82. Cruz F, Julca I, Gómez-Garrido J, Loska D, Marcet-Houben M, Cano E, Galán B, Frias L,  
41 888 Ribeca P, Derdak S *et al*: **Genome sequence of the olive tree, Olea europaea.**  
42 889 *GigaScience* 2016, **5**(1):29.
- 43 890 83. Plomion C, Aury JM, Amselem J, Alaeitabar T, Barbe V, Belser C, Berges H, Bodenes C,  
44 891 Boudet N, Boury C *et al*: **Decoding the oak genome: public release of sequence data,**  
45 892 **assembly, annotation and publication strategies.** *Molecular ecology resources* 2016,  
46 893 **16**(1):254-265.
- 48 894 84. Park B-S, Lee H-K, Lee S-E, Piao X-L, Takeoka G, Wong R, Ahn Y-J, Kim J-H: **Antibacterial**  
49 895 **activity of Tabebuia impetiginosa Martius ex DC (Taheebo) against Helicobacter pylori.**  
50 896 *Journal of Ethnopharmacology* 2006, **105**(1-2):255-262.
- 51 897 85. Gómez Castellanos R, Prieto J, Heinrich M: **Red Lapacho (Tabebuia impetiginosa)—A**  
52 898 **global ethnopharmacological commodity?** *Journal of Ethnopharmacology* 2009,  
53 899 **121**(1):1-13.
- 55 900 86. Byeon S, Chung J, Lee Y, Kim B, Kim K, Cho J: **In vitro and in vivo anti-inflammatory**  
56 901 **effects of taheebo, a water extract from the inner bark of Tabebuia avellaneda.**  
57 902 *Journal of Ethnopharmacology* 2008, **119**(1):145-152.
- 59 903 87. Koyama J, Morita I, Tagahara K, Hirai K-I: **Cyclopentene dialdehydes from Tabebuia**  
60 904 **impetiginosa.** *Phytochemistry* 2000, **53**(8):869-872.
- 61  
62  
63  
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

- 905 88. Hussain H, Krohn K, Ahmad VU, Miana GA, Green IR: **Lapachol: An overview**. *Arkivoc*  
906 2007, **2007**(2):145.
- 907 89. Widhalm J, Rhodes D: **Biosynthesis and molecular actions of specialized 1,4-**  
908 **naphthoquinone natural products produced by horticultural plants**. *Horticulture*  
909 *Research* 2016, **3**:16046.
- 910 90. Romagnoli M, Segoloni E, Luna M, Margaritelli A, Gatti M, Santamaria U, Vinciguerra V:  
911 **Wood colour in Lapacho (*Tabebuia serratifolia*): chemical composition and industrial**  
912 **implications**. *Wood Science and Technology* 2013, **47**(4):701-716.
- 913 91. Alkan C, Sajjadian S, Eichler EE: **Limitations of next-generation genome sequence**  
914 **assembly**. *Nat Meth* 2011, **8**(1):61-65.
- 915  
916

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

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

919

920

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

923

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

924

925

926

927

928 **Table 3.** The distribution of the minimal introns (53–125 bp) and the minimal-intron-containing  
 929 genes – as the number of genes with at least one minimal intron – from selected plant species in  
 930 comparison to the *H. impetiginosus* genome assembly.

931

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

932

1  
2  
3  
4 933 **Figure Legends**

5  
6 934

7 935 **Figure 1. Depth of coverage analysis.** (A) Histograms of k-mer frequencies in the filtered read  
8 936 data for  $k = 25$  (red) and GenomeScope modeling equation on *H. impetiginosus* (blue). The x-axis  
9 937 shows the number of times a k-mer occurred (coverage). The vertical dashed dark blue lines  
10 938 correspond to the mean coverage values for unique heterozygous k-mers (left peak) and unique  
11 939 homozygous k-mers (right peak). (B) Density plot of read depth based on mapping all short  
12 940 fragment reads back to the assembled scaffolds (red). Left peak (at depth = 34x) corresponds to  
13 941 regions where the assembler created two distinct scaffolds from divergent putative haplotypes.  
14 942 The right peak (at depth = 67x) contains scaffolds from regions where the genome is less variable,  
15 943 allowing the assembler to construct a single contig combining homologue sequences. Histograms  
16 944 of Poisson modeling for read depth in the assembly (green,  $\lambda = 34$ ; blue,  $\lambda = 67$ ) are  
17 945 shown.

18 946

19 947 **Figure 2. Depth of coverage analysis for the haplotype-reduced assembly.** (A) Density plot of  
20 948 read depth based on mapping all short fragment reads back to the haplotype-reduced assembled  
21 949 sequences after identification and removal of redundant sequences due the structural  
22 950 heterozygosity in the genome. (B) Density plot for average sequencing coverage per-scaffold on  
23 951 the final assembly. The observed number of scaffolds in the final haplotype-reduced assembly  
24 952 and the respective read coverage (blue line) is shown in comparison to a Poisson process  
25 953 approximation (red line) with  $\lambda = 63x$ , the observed average sequencing coverage in the  
26 954 useful read data.

27 955

28 956 **Figure 3. Repeat content of the *H. impetiginosus* genome assembly.** (A) The density of  
29 957 interspersed and tandem repeat as percent of the assembly. The size of the circles represents the  
30 958 number of copies in the assembly for each family of repeats; (B) Distribution of sizes of the  
31 959 consensus sequences for repeat families identified using *de novo* and homology methods for  
32 960 repeat characterization.

33 961

34 962 **Figure 4. Transcriptome quality assessment** (A) similarity search of *H. impetiginosus* putative  
35 963 peptides against source database of plant protein sequences using BLASTP algorithm (e-value  $1e-$   
36 964 6). Transcript count means the number of peptides of *H. impetiginosus* with best hit against the

1  
2  
3  
4 965 source database using bit-score and grouping results by taxon name. Transcript score corresponds  
5  
6 966 to the average bit-score overall hits for each group using the best hit. We ordered taxon groups  
7  
8 967 by their average bit-score overall hits and used Welch's t-test to compare the distributions of bit-  
9  
10 968 score hits between two adjacent groups with p-values <0.01 (ns = non-significant; \*\*\* significant);  
11  
12 969 (B) Completeness of the expected gene space of the genome assembly, estimated with BUSCO.  
13  
14 970 The estimates were compared with genome annotations for other lamids, *Erythranthe guttata*  
15  
16 971 and *Olea europaea*.

17 972  
18 973 **Figure 5.** Contig termini analysis to investigate the possible genomic features associated with gaps  
19  
20 974 in the genome assembly. Contigs were created from the genome assembly with the "cutN -n 1"  
21  
22 975 command from seqtk program, which cut at each gap (of at least one basepair, i.e. one or more  
23  
24 976 Ns). The figure shows the percentage of contig termini (the position of the terminal nucleotides  
25  
26 977 of each contig) intersecting with different annotations of the genome.

27 978  
28  
29 979 **Figure 6. Genes of the biosynthetic pathway of specialized quinoids.** O-succinylbenzoate (OSB)  
30  
31 980 pathway depicting the number of *H. impetiginosus* (Himp) annotated genes for the known  
32  
33 981 enzymes that lead to the biosynthesis of the naphthoquinones, including lapachol. For  
34  
35 982 comparison, it also shows the numbers of genes for the closely related *Mimulus guttatus* (Mgut),  
36  
37 983 *Solanum lycopersicum* (Slyc), for the model *Arabidopsis thaliana* (Ath), and for the tree species  
38  
39 984 *Eucalyptus grandis* (Egr) and *Populus trichocarpa* (Potri). The pathway was modified from [89].

40 985

41 986

42  
43 987 **Supplementary material**

44  
45 988

46  
47 989 **Table S1.** Summary of the sequence data generated for the genome assembly of *Handroanthus*  
48  
49 990 *impetiginosus* based on the ALLPATHS-LG algorithm.

50  
51 991

52 992 **Figure S1.** The *Handroanthus impetiginosus* (Mart. ex DC.) Mattos (syn. *Tabebuia impetiginosa*,  
53  
54 993 Bignoniaceae), tree UFG-1 whose genome was sequenced.

55  
56 994

57 995 **Figure S2.** Flow cytometry results of the sequenced tree UFG-1 of *H. impetiginosus*. Flow  
58  
59 996 cytometry estimate of the nuclear DNA content was carried out using young leaf tissue on a BD

60

61

62

63

64

65



1  
2  
3  
4 997 Accuri™ C6 Plus personal flow cytometer. *Pisum sativum* (genome size 9.09 pg/2C or ~4380  
5  
6 998 Mb/1C) was used as standard for comparison (M2). The estimate of nuclear DNA content for *H.*  
7  
8 999 *impetiginosus* (M1) averaged over 10 readings was 1.155 pg/2C or 557.3 ± 39 Mb/1C.

1000

11 1001 **Figure S3.** Overview of the analytical pipeline with the bioinformatics steps and tools employed  
12 for genome (black arrows) and transcriptome assembly (red arrows), and for gene prediction  
13 1002 and annotation (blue arrows). Bioinformatics programs are indicated in italic, blue, and the main  
14 1003 file formats in red. The input sequences are highlighted in yellow boxes and the main products  
15 1004 in green.

1005  
1006

22 1007 **Figure S4.** Distribution and characterization of simple sequence repeats in *Handroanthus*  
23 1008 *impetiginosus* genome (A) Histogram of different motifs ranging from 1 to 6 bp (B) Distribution of  
24 1009 the simple sequence repeats length detected in the genome assembly.

1010

29 1011 **Figure S5.** Comparison of the gene features parameters, such as number and length, between *H.*  
30 1012 *impetiginosus* and the other selected dicot plant across distinct lineages of Rosids (*A. thaliana* and  
31 1013 *P. trichocarpa*) and Asterids (*E. guttata* and *S. lycopersicum*). Frequency histograms are shown  
32 1014 according to the whole-genome gene content annotation for (A) the complete predicted gene  
33 1015 structure (B) exons and (C) introns. Dashed vertical lines are the average lengths for the gene  
34 1016 features.

1017

41 1018 **Figure S6.** Histograms for Gene Ontology broader term annotations in the *H. impetiginosus*  
42 1019 genome assembly. Terms for the Biological Process ontology were summarized with WEGO by the  
43 1020 second tree level setting. The Pearson Chi-Square test was applied to indicate significant  
44 1021 relationships between *H. impetiginosus* and the lamid *Erythranthe guttata* regarding the number  
45 1022 of genes (at  $\alpha \geq 5\%$ ). (A) Terms displaying remarkable relationship between the two datasets; (B)  
46 1023 terms with a significant difference between the two datasets.

1024

54 1025 **Figure S7.** Same as Figure S6 but showing comparison between numbers of genes assigned to GO  
55 1026 broader terms for *H. impetiginosus* and the lamid *Olea europaea*.

1027  
58  
59  
60  
61  
62  
63  
64  
65

1  
2  
3  
4 1028 **Figure S8.** Sequence length distribution from the assemblies of *H. impetiginosus* and other two  
5  
6 1029 highly heterozygous trees of the genus *Quercus*. Figure shows density plots for the size of scaffolds  
7  
8 1030 with 2 kbp or longer in the three assemblies. Contigs metrics were computed by cutting at each  
9  
10 1031 gap (of at least 25 base pair, i.e. 25 or more Ns). Scaffolds and contigs length were plotted using  
11  
12 1032 the common logarithm to respond to skewness towards large values.

13 1033

14  
15 1034 **File S1.** Evidences adopted to support protein-coding loci identification and assignment in the  
16  
17 1035 *H. impetiginosus* genome assembly. Two qualifiers – high-confidence and low-confidence – were  
18  
19 1036 added to the locus based on the reported evidences.

20 1037

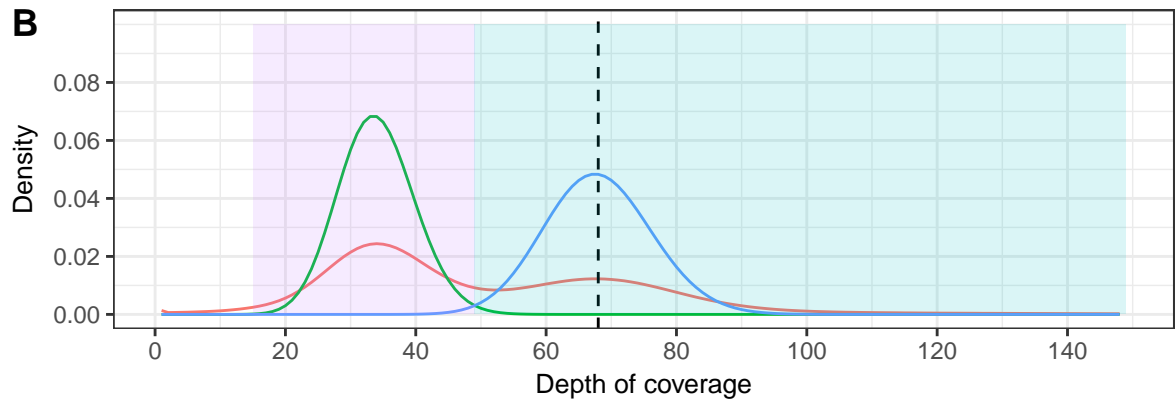
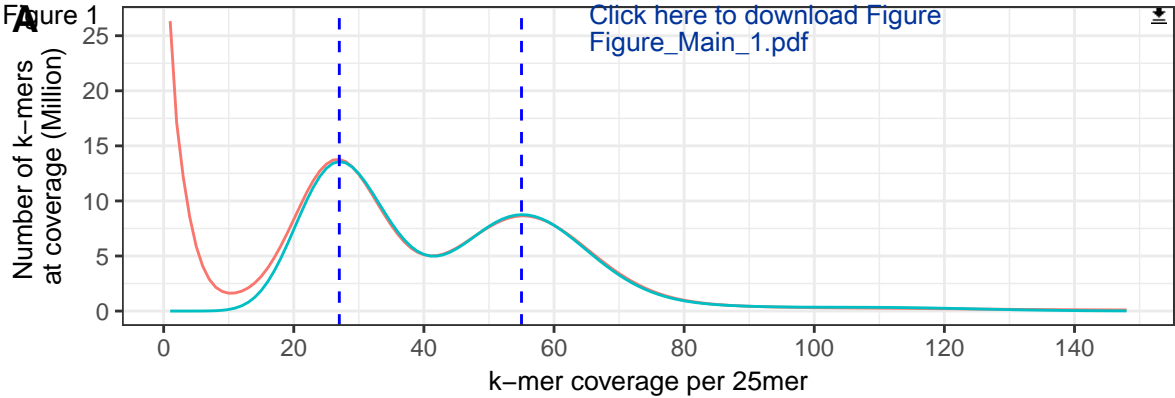
21  
22 1038 **File S2.** Genome assembly metrics from the assemblies of *H. impetiginosus* and other two highly  
23  
24 1039 heterozygous trees of the genus *Quercus*. Comparison between metrics based on the  
25  
26 1040 `assemblathon_stats` script part of the `assemblathon2-analysis` package  
27  
28 1041 (<https://github.com/ucdavis-bioinformatics/assemblathon2-analysis>). Metrics were computed  
29  
30 1042 for scaffolds with 2 kbp or longer in length. Genomic sequences in scaffolds for *Quercus lobata*  
31  
32 1043 was obtained from <https://valleyoak.ucla.edu/genomicresources/> (accessed on 9/20/2017). For  
33  
34 1044 *Quercus rubra*, genomic sequences in scaffolds were downloaded from the ENA (European  
35  
36 1045 Nucleotide Archive) repository, accessions LN776247-LN794156.

36 1046

37  
38 1047

39  
40 1048

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



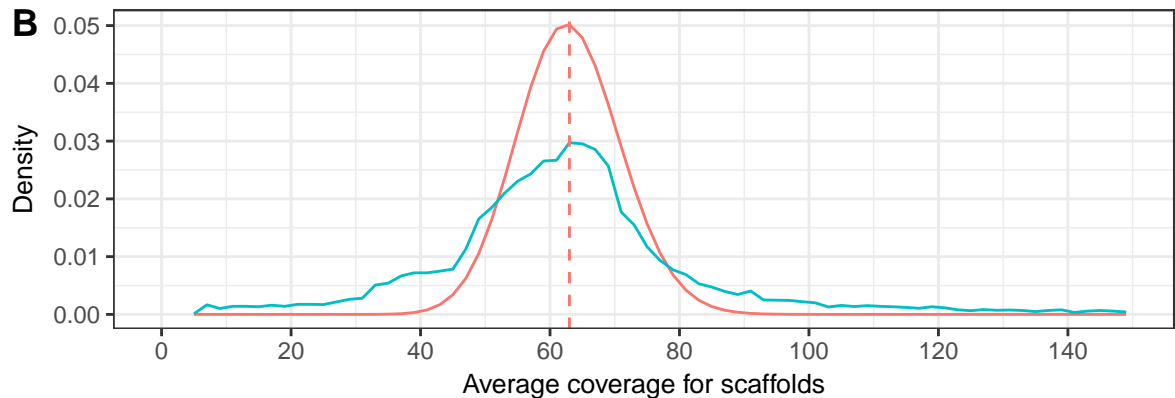
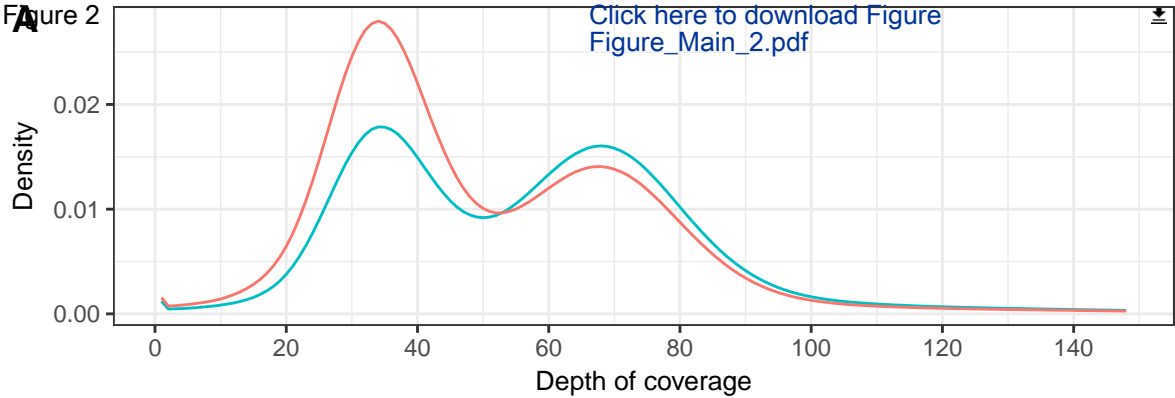
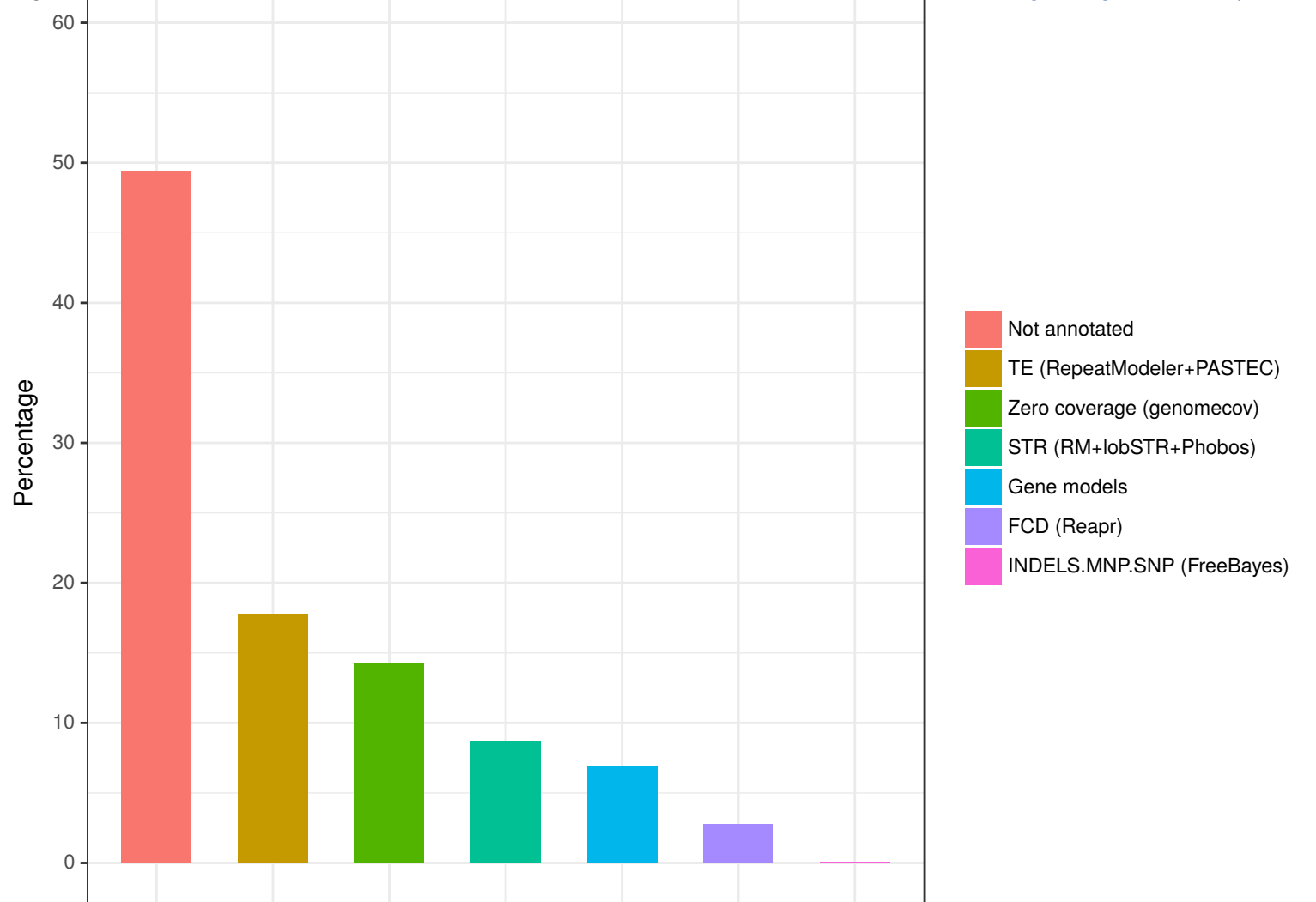
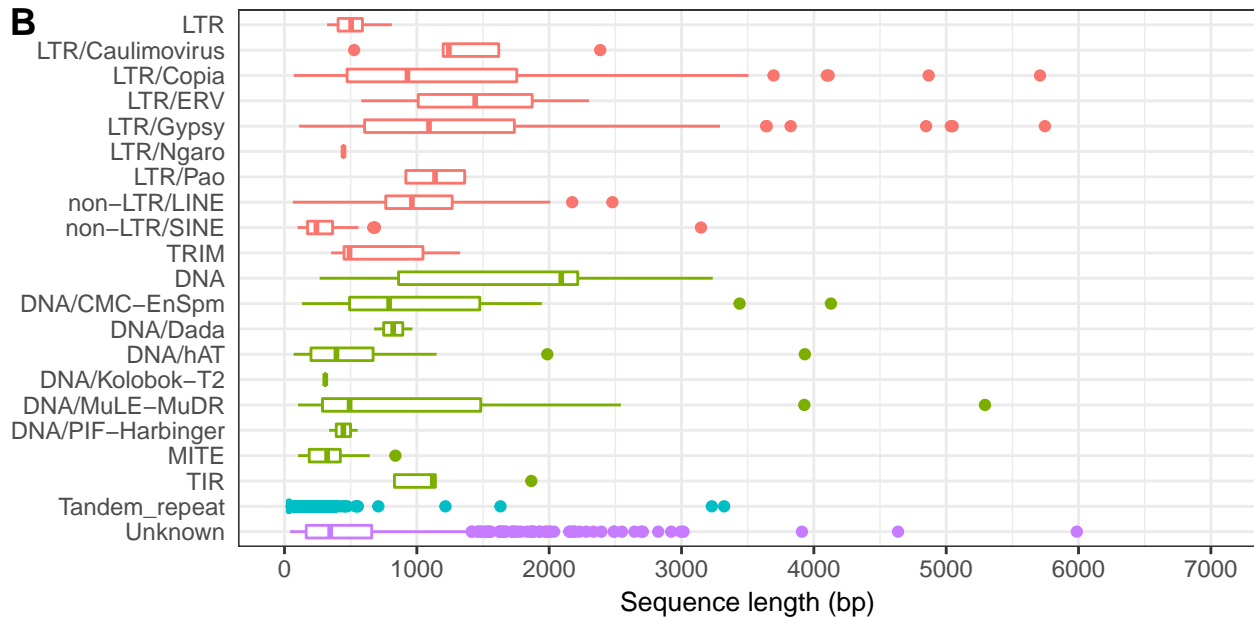
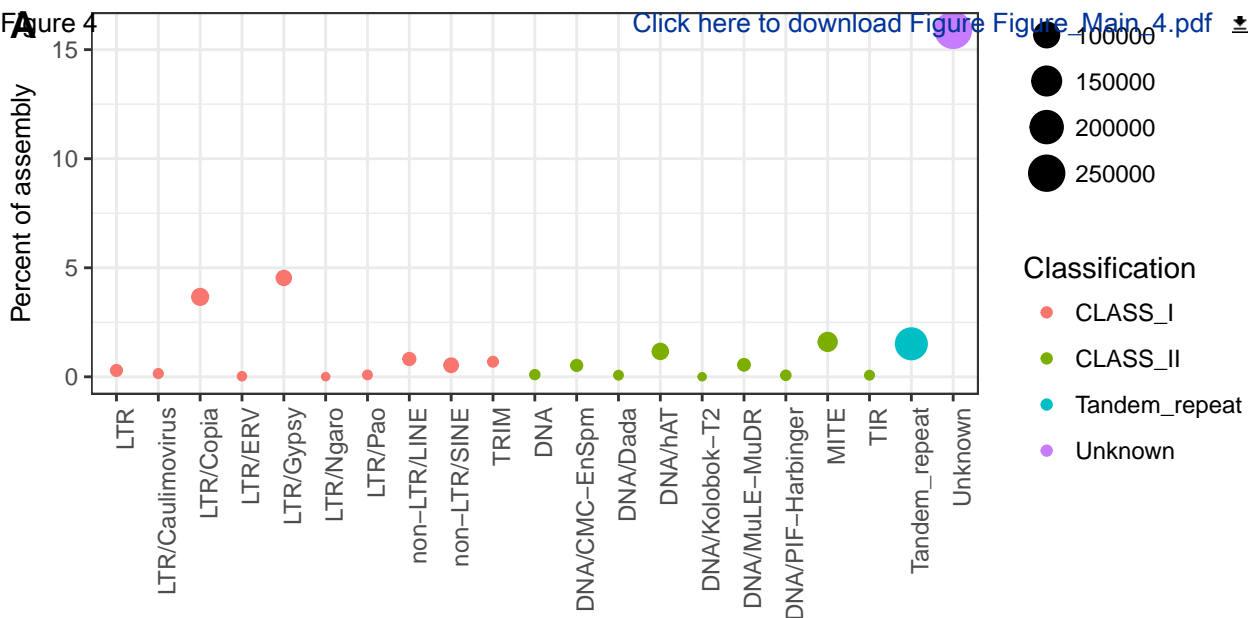


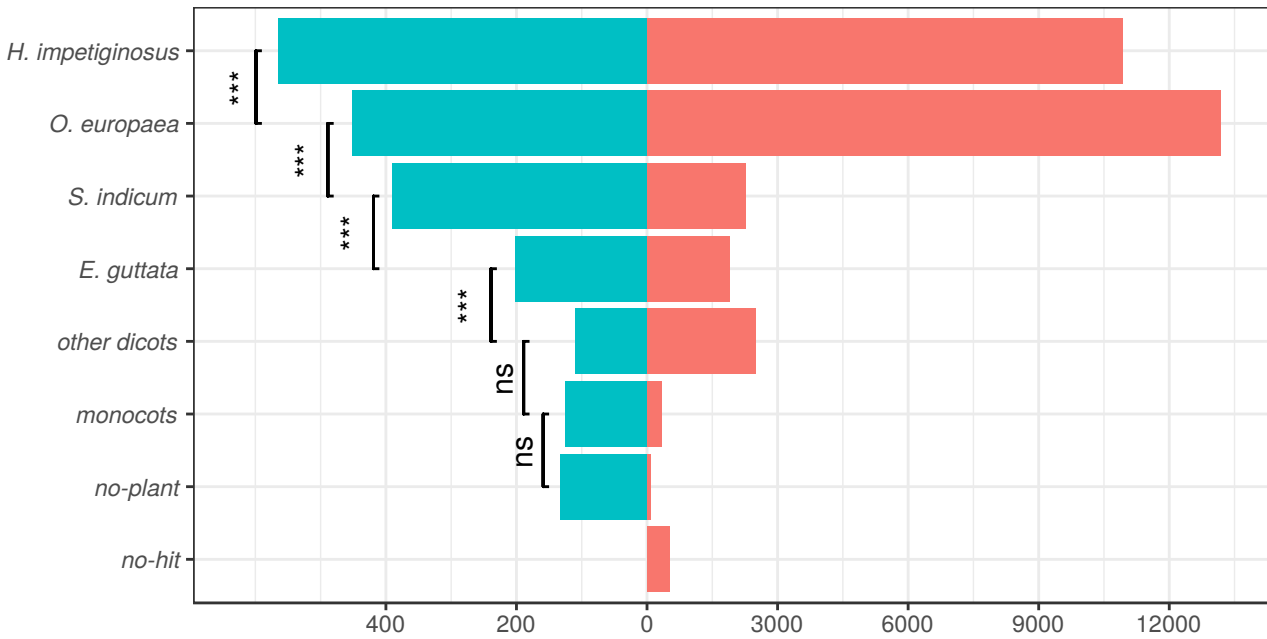
Figure 3

[Click here to download Figure Figure\\_Main\\_3.pdf](#)



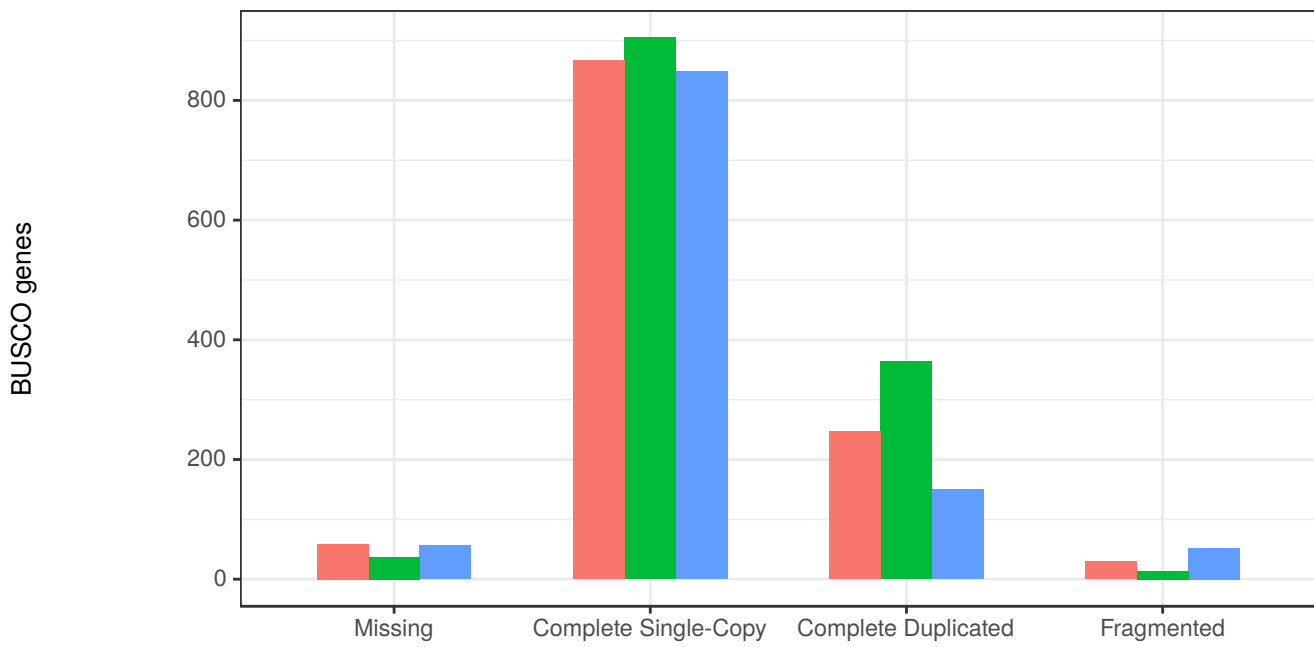


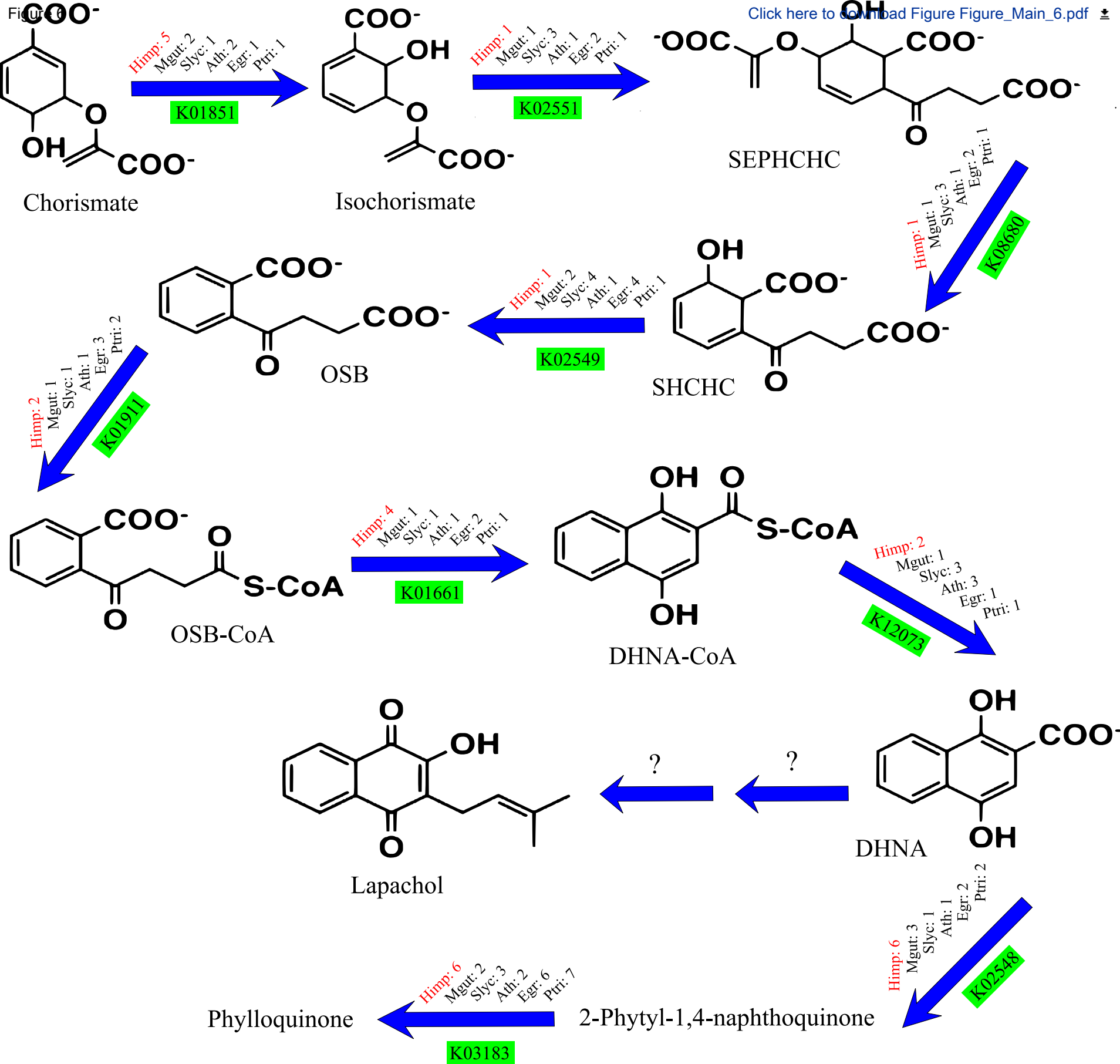
transcript count transcript score



**B**

*H. impetiginosus* *O. europaea* *E. guttata*









Click here to access/download

**Supplementary Material**

Supp\_Material\_H.impetiginosus\_GIGA-D-17-  
00159.R1.docx



Click here to access/download  
**Supplementary Material**  
Supplementary\_FileS1.xlsx





Click here to access/download  
**Supplementary Material**  
Supplementary\_FileS2.xlsx

