

DATA NOTE

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

Orzenil Bonfim Silva-Junior^{1,2}, Dario Grattapaglia^{1,2}, Evandro Novaes³ and Rosane G. Collevatti^{4,*}

¹EMBRAPA Recursos Genéticos e Biotecnologia, EPqB, Brasília, DF. 70770–910, Brazil, ²Programa de Ciências Genômicas e Biotecnologia – Universidade Católica de Brasília, SGAN 916 Modulo B, Brasília, DF 70790-160, Brazil, ³Escola de Agronomia, Universidade Federal de Goiás, CP 131. Goiânia, GO. 74001–970, Brazil and ⁴Laboratório de Genética and Biodiversidade, Instituto de Ciências Biológicas, Universidade Federal de Goiás. Goiânia, GO. 74001–970, Brazil

*Correspondence address. Rosane Garcia Collevatti, Instituto de Ciências Biológicas, Universidade Federal de Goiás, 74001–970, Goiânia, GO, Brasil. Tel: +55 62 3521-1729; E-mail: rosanegc68@hotmail.com

Abstract

Background: *Handroanthus impetiginosus* (Mart. ex DC.) Mattos is a keystone Neotropical hardwood tree widely distributed in seasonally dry tropical forests of South and Mesoamerica. Regarded as the “new mahogany,” it is the second most expensive timber, the most logged species in Brazil, and currently under significant illegal trading pressure. The plant produces large amounts of quinoids, specialized metabolites with documented antitumorous and antibiotic effects. The development of genomic resources is needed to better understand and conserve the diversity of the species, to empower forensic identification of the origin of timber, and to identify genes for important metabolic compounds. **Findings:** The genome assembly covers 503.7 Mb (N50 = 81 316 bp), 90.4% of the 557-Mbp genome, with 13 206 scaffolds. A repeat database with 1508 sequences was developed, allowing masking of ~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 messenger RNA 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. **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.

Keywords: heterozygous genome; RNA-Seq; transposable elements; quinoids; *Bignoniaceae*

Received: 28 June 2017; Revised: 27 September 2017; Accepted: 30 November 2017

© The Author(s) 2017. Published by Oxford University Press. This is an Open Access article distributed under the terms of the Creative Commons Attribution License (<http://creativecommons.org/licenses/by/4.0/>), which permits unrestricted reuse, distribution, and reproduction in any medium, provided the original work is properly cited.

Data Description

Context

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

Species of *Handroanthus* and *Tabebuia* have virtually no genomic tools and resources, beyond a handful of 21 microsatellites [17] with their known caveats for more sophisticated genetic analyses in the areas of population genomics and evolution [18]. Whole-genome sequencing has now become accessible to a point that efforts to develop improved genomic resources for such species are possible and warranted. We built a preliminary assembly of the nuclear genome of a single individual of *Handroanthus impetiginosus* based on short reads and longer mate-pair DNA sequence data to provide the necessary framework for the development of genomic resources to support multiple genomic and genetic analyses of this keystone Neotropical hardwood tree regarded as the “new mahogany.” It is the second most expensive timber and the most logged species in Brazil [19], exported largely to North America for residential decking and currently under significant illegal trading pressure. Additionally, the tree produces large amounts of natural products such as those of quinoid systems (1,4-anthraquinones, 1,4-naphthoquinones, and 1,2-furanonaphthoquinones), specialized metabolites with promising antitumor, anti-inflammatory, and antibiotic effects [20, 21]. The high pressure of logging and illegal trading on this species with a notable ecological keystone status urges conservation efforts of existing populations.

Methods

Sample collection and sequencing

DNA of a single adult tree of *H. impetiginosus* (UFG-1) (Fig. 1) was extracted using Qiagen DNeasy Plant Mini kit (Qiagen,



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

DK). Flow cytometry was used to check the genome size of tree UFG-1, indicating a genome size of (557 ± 39) Mb/1C (Fig. S1) consistent with published estimates [22]. Total RNA from shoots of 5 seedlings and from the differentiating xylem of the adult tree (UFG-1) was extracted using Qiagen RNeasy Plant Mini kit (Qiagen, DK) and pooled for RNA sequencing. DNA and RNA sequencing was performed at the High-Throughput Sequencing and Genotyping Center of the University of Illinois Urbana-Champaign. The following libraries were generated for sequencing: (1) 2 shotgun genomic libraries of short fragments (300 bp and 600 bp) from tree UFG-1, (2) 1 shotgun library from combined pools of 5 RNA samples tagged with a single index sequence. Paired-end sequencing, 2×150 nt, was performed in 2 lanes of an Illumina HiSeq 2500 instrument (Illumina, CA, USA). Three additional mate-pair libraries (fragment lengths of 4 kb to 5.5 kb, 8 kb to 10 kb, and 15 kb to 20 kb) for UFG-1 were also sequenced in 2 lanes of an Illumina HiSeq 2000 instrument (2×101 bp). This long-range sequence resource was used to generate the final genome assembly for annotation. A complete overview of the genome assembly and annotation pipeline is provided (Fig. S2).

Genome assembly using short paired-end and mate-pair sequencing data

Short reads and mate-pair reads were stripped of sequencing adapters using *Fastq-mcf* [23]. Reads that mapped to a database containing mitochondrial and chloroplast genomes of plants with *Bowtie1* (option `-v 3 -a -m 1`) [24] were discarded. Mate-pair reads were inspected using a Perl script (*TrimAdaptor.pl*), and sequences that did not contain the circularization adaptor were discarded. By using the filtered short reads, *Jellyfish2* (*Jellyfish*, [RRID:SCR.005491](#)) [25] and *GenomeScope* [26] were applied to obtain estimates of the *H. impetiginosus* genome size, repeat fraction, and heterozygosity prior to the assembly. *ALLPATHS-LG* (*ALLPATHS-LG*, [RRID:SCR.010742](#)) [27] was used for *de novo* assembly of the sequence data from both paired-end and mate-pair data, with default options, in a stepwise strategy for error correction of reads, handling of repetitive sequences, and use of mate-pair libraries.

Transposable elements and repetitive DNA

Repetitive elements were detected and annotated on the genome assembly with the *RepeatModeler de novo* repeat family identification and modeling package (*RepeatModeler*, [RRID:SCR.015027](#)) [28]. Using *RECON*, *RepeatScout*, and *Tandem Repeat Finder*, repetitive sequences were detected in the scaffolds longer than 10 kb using a combination of similarity-based and *de novo* approaches. The TE sequences were evaluated using modeling capabilities of the *RepeatModeler* program, with default settings, to compare the TE library against the entire assembled sequences and to refine and classify consensus models of putative interspersed repeats. A complementary analysis intended to augment the number of TE sequences classified according to current criteria [29] was performed using the *PASTEC* program [30]. *RepeatMasker Open-4.0* (*RepeatMasker*, [RRID:SCR.012954](#)) [31] was used with the sequences from the *de novo* repetitive element library to annotate the interspersed repeats and to detect simple sequence repeats (SSRs) on the genome assembly.

Protein-coding gene annotation

Protein-coding gene annotation was performed with a pipeline that combines RNA-Seq assembled transcript and protein alignments to the reference with *de novo* prediction methods (Fig. S2). RNA-Seq reads were screened for the presence of adapters, which were removed using *Fastq-mcf* [23]. *Trimmomatic* (*Trimmomatic*, [RRID:SCR.011848](#)) [32] was used to (1) remove low-quality, no-base-called segments (Ns) from sequencing reads; (2) scan the read with a 4-base sliding window, cutting when the average quality per base dropped below 15; and (3) remove reads shorter than 32 bp after trimming. Trimmed reads mapped to mitochondrial, chloroplast, and ribosomal sequences from plants with *Bowtie1* (options `-v 3 -a -m 1`) [24] were also removed. Transcript *de novo* assemblies were performed using *SOAP-Transdenovo* [33] and *Trinity de novo* [34] from the processed reads. The assemblies were concatenated and used as input to *EvidentialGene* [35], a comprehensive transcriptome pipeline to identify likely complete coding regions and their proteins in the final, combined, transcriptome assembly. Gene modeling was carried out using standard procedures and tools described, for instance, in Schmutz et al. [36]. In summary, a genome-guided transcriptome assembly of *H. impetiginosus* was performed with the JGI PERTRAN RNA-Seq Read Assembler pipeline

[37] using both the RNA-Seq trimmed reads and sequences from the *de novo* transcript assembly. Loci were identified by the assembled transcript alignments using *BLASTX* [38] and *EXONERATE* [39] alignments of peptide sequences to the repeat-soft-masked genome using *RepeatMasker* [40], based on a transposon database developed as part of this genome assembly annotation. Known peptide sequences included manually curated datasets for plant species available from UniProtKB/Swiss-Prot [41] and sequences available from Phytozome [1], version 11, for *Arabidopsis thaliana*, *Oryza sativa*, *Erythranthe guttata*, *Solanum lycopersicum*, *Solanum tuberosum*, *Populus trichocarpa*, and *Vitis vinifera*. Gene structures were predicted by homology-based predictors, *FGENESH++*, *FGENESH.EST* [42, 43], and *GenomeScan* [44]. Gene predictions were improved by Program to Assemble Spliced Alignment (*PASA*, [RRID:SCR.014656](#)) [45], including adding Untranslated Regions (UTRs), correcting splicing, and adding alternative transcripts. *PASA*-improved gene model peptides were subjected to peptide homology analysis with the above-mentioned proteomes to obtain Cscore values and peptide coverage. Cscore is the ratio of the peptide Basic Local Alignment Search Tool for Proteins (BLASTP) score to the mutual best hit BLASTP score, and peptide coverage is the highest percentage of peptide aligned to the best homolog. A transcript was selected if its Cscore value was greater than or equal to 0.5 and its peptide coverage was greater than or equal to 0.5 or if it had transcript coverage but the proportion of its coding sequence overlapping repeats was less than 20%. For gene models where greater than 20% of the coding sequence overlapped with repeats, the Cscore value was required to be at least 0.9 and homology coverage was required to be at least 70% to be selected. Selected gene models were then subjected to classification analysis using *InterProScan 5* (*InterProScan*, [RRID:SCR.005829](#)) [46] for PFAM domains, PANTHER, Enzyme Commission Number (EC), and KEGG categories. Gene ontology annotation was obtained, where possible, from Interpro2GO and EC2GO mappings.

Data Validation and Quality Control

Global properties of the *H. impetiginosus* tree genome from the unassembled reads

Sequencing of the *H. impetiginosus* tree genome generated c. 599 million reads, comprising 73 Gbp of sequence data. This represents nearly $\times 132$ the expected sequence coverage. After removal of adaptors, followed by standard error correction and trimming with *ALLPATHS-LG*, with default options, c. 46 Gbp of data was found useful for the assembly process, yielding sequencing coverage of $\times 82$ ($\times 63$ from the fragments libraries and $\times 19$ from the mate-pair libraries). The estimated physical coverage was $\times 400$ based on the observed fragment size distributions (Table S1). *ALLPATHS-LG* k-mer spectrum frequency analysis (at $K = 25$) on useful reads, error-corrected reads, estimated a haploid genome size of 540 968 531 bp, a repeat fraction of 38.0%, and a single nucleotide polymorphism (SNP) rate of 1/88 bp (1.14%). An alternative analysis of the k-mer frequencies using *GenomeScope* [26] produced a haploid genome size estimate of 503 748 072 bp, repetitive content of 36.6%, and an SNP rate of 1/60 bp (1.65%). Both estimates (Fig. 2A) are consistent with the flow cytometry estimates and in line with the expectations regarding the heterozygous content of the *H. impetiginosus* genome, a predominantly outcrossed tree [47]. Sequencing errors caused an extreme peak at $k = 1$ in the k-mer frequency distribution. Both k-mer histograms display 2 distinct peaks comprising the largest area of each histogram

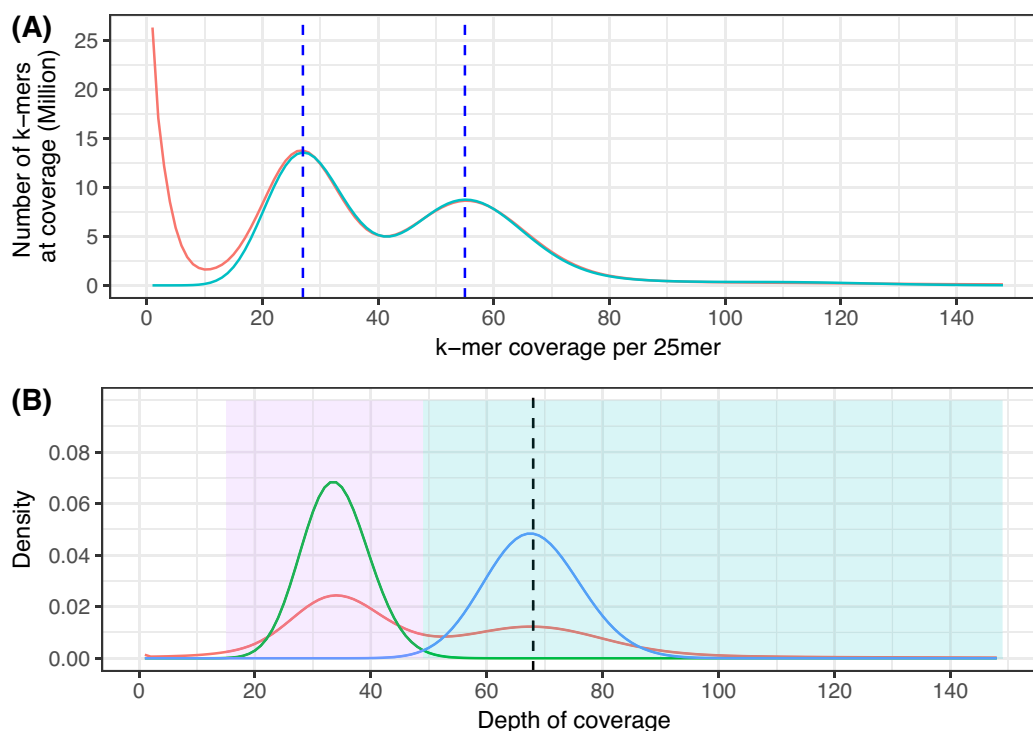


Figure 2: Depth of coverage analysis. (A) Histograms of k-mer frequencies in the filtered read data for $k = 25$ (red) and GenomeScope modeling equation on *H. impetiginosus* (blue). The x-axis shows the number of times a k-mer occurred (coverage). The vertical dashed dark blue lines correspond to the mean coverage values for unique heterozygous k-mers (left peak) and unique homozygous k-mers (right peak). (B) Density plot of read depth based on mapping all short fragment reads back to the assembled scaffolds (red). Left peak (at depth = $\times 34$) corresponds to regions where the assembler created 2 distinct scaffolds from divergent putative haplotypes. The right peak (at depth = $\times 67$) contains scaffolds from regions where the genome is less variable, allowing the assembler to construct a single contig combining homologue sequences. Histograms of Poisson modeling for read depth in the assembly (green, $\lambda = 34$; blue, $\lambda = 67$) are shown.

at depths 27 and 55. The bimodal distributions characterize the expected behavior for k-mer frequencies of a heterozygous diploid genome, as seen, for example, in the recently reported oak genome [48]. In the right homozygous peak (at $K = 55$), k-mers are shared between the 2 homologous chromosomes. The left or heterozygous peak, with half the k-mer depth of the homozygous peak, contains k-mers that are unique to each haplotype due to heterozygosity. The difference in height between these peaks (heterozygous/homozygous ratio) is a measure of the heterozygosity within the genome, which is 1.65% according to the GenomeScope modeling equation.

Genome assembly

State-of-the-art haploid genome assembler pipelines from short-read ALLPATHS-LG [27] and SOAPdenovo2 (SOAPdenovo2, RRID:SCR.014986) [49] were considered for an initial evaluation on the dataset of reads. Two relatively new algorithms specifically developed for *de novo* assembly of heterozygous genomes, MaSuRCA (MaSuRCA, RRID:SCR.010691) [50] and PLATANUS (PLATANUS, RRID:SCR.015531) [51], were also attempted as alternatives to the other 2 assemblers designed for genomes of low heterozygosity. Reads were first preprocessed and error corrected using the algorithms provided by each assembler. PLATANUS was set to run, but after 10 weeks it did not produce any result in an Intel(R) Xeon(R) server with 64×7560 2.27-GHz CPUs, 256 GB RAM, except for the k-mer count table on the input trimmed reads. After 9 week-long runtimes in an Intel(R) Xeon(R) server with 64×7560 2.27-GHz CPUs, 512 GB RAM, MaSuRCA successfully completed the generation of the super-reads

from the trimmed reads, but the process was aborted on the overlap-correction process in the Celera Assembler due to excessive CPU usage. SOAPdenovo2 ran very fast (3 days) but produced an assembly with a total scaffold size of 860 Mbp. Analysis with SOAPdenovo2 was run with different k-mer sizes, from 31 to 71, step of 10, but none of them produced a reasonable assembly size in view of the expected size estimated by flow cytometry and the k-mer frequency. ALLPATHS-LG was therefore used to assemble the genome with default options. The short reads from fragmented libraries were error-corrected using default settings (K-mer size of 24, ploidy of 2), fragment-filled, and assembled into initial unipaths (k-mer size of 96, ploidy of 2). Jumping reads from the mate-pair libraries were then aligned to the unipaths and all alignments were processed in a seed-extension strategy with junction point recognition within the read aimed to remove invalid and duplicate fragments to perform error correction and initial scaffolding. This initial process produced an assembly graph that was turned into scaffolds by analyzing branch points in the graph topology. This late process converted single-base mismatches into ambiguous base codes at branch. It also flattened some other structural features of the assembly including short indels. The contig assembly comprised 109 064 sequences of a length of 500 bp or longer with total length of 466 314 780 bp. Genome assembly after scaffolding comprised 57 815 scaffolds of length 1 kbp or longer with a total length of 610 091 865 bp and N50 of 57 Kbp. The fraction of bases captured in gaps was 23.9%, and the rate of ambiguous bases for all bases captured in the assembly was 0.24%. This assembly was only slightly larger in size (<10%) than the empirically determined genome size using flow cytometry [22].

Alternative scaffold and gap-filling

Although the ALLPATHS-LG performance was good in recovering the expected genome size in the assembled contigs, there was a high fraction of the bases captured in gaps in the scaffolds (~one-fourth of the total genome assembly). *De novo* assembly algorithms applied to moderate to high levels of heterozygosity cannot match the performance achieved in assemblies of homozygous genomes, especially at the contig assembly level [52]. We thus used the assembled contigs to perform an alternative scaffolding step with SSAKE-based Scaffolding of Pre-Assembled Contigs after Extension (SSPACE, [RRID:SCR.005056](#)) [53] using the error-corrected short fragment reads and the jumping reads. In this approach, genome assembly comprised 16 090 scaffolds of a length of 1 kbp or longer with a total length of 577 446 088 bp and N50 of 95 Kbp, respectively. The fraction of bases captured in gaps dropped from 23.9% to 18.9% in contrast to ALLPATHS-LG scaffolding, totaling 109 533 288 bp. The rate of ambiguous bases for all bases captured in the assembly dropped from 0.24% to 0.13%. All preprocessed reads were reused in an attempt to close the intra-scaffold gaps using the GapCloser (GapCloser, [RRID:SCR.015026](#)) [54] algorithm. Genome assembly after gap-filling was 586 206 884 bp in 15 671 scaffolds of a length of 1 kbp or longer, and only 20 583 469 bp (3.51% of the genome assembly) remained in 24 907 gaps. The N50 of scaffolds of a length of 1 kbp or longer, with gaps, was 97 344 Kb (L50 = 1792). Sequences longer than 20 kb were assembled in only 6791 scaffolds, totaling 538 102 146 bp, ~97% of the genome size estimated from flow cytometry (557 Mb).

Evaluation of accuracy of the genome assembly

A subset of fragments and jumping read pairs (~ $\times 15$ sequencing coverage each) was used to uncover inaccuracies in the genome assembly. Scaffolds with identified errors were broken or flagged for inspection. Recognition of Errors in Assemblies using Paired Reads (REAPR) [55] was used to test each base of the genome assembly looking for small local errors (such as a single base substitution, and short insertions or deletions) and structural errors (such as scaffolding errors), located by means of changes to the expected distribution of inferred sequencing fragments from the mapped reads using SMALT v0.7.6 [56]. REAPR reported that only 343 588 027 (~60%) bases in the assembly should be free of errors, with 5476 reported (1658 within contigs, 3818 over gaps) in the remaining 242 618 857 bp. The most frequent (~92%) type of inaccuracy reported was *Perfect.cov* and *Link*. *Perfect.cov* means low coverage of perfect uniquely mapping reads while *Link* describes situations in which reads map elsewhere in the assembly. The recognition of this inaccuracy at the base pair level should thus reflect the repetitive nature of the genome, as inferred from the k-mer frequency spectra analysis (~36%–38% of repeats). Besides the base pair inaccurate calls due to repeats, other structural problems in the assembly were identified based on sequence coverage differences from the expected fragment size distribution, and the program used this information to break these. Given the high heterozygosity and divergence between haplotypes on this diploid genome sequence, homologous sequences can assemble separately or merge. Moreover, unresolved repeat structures in the assembly might also contribute heavily to this issue. Structural errors in REAPR were likely called at the boundaries of these regions. The final genome assembly after REAPR breaks had 19 319 sequences of a length of 1 kbp or longer, with 576 829 188 bp. The N50 size of scaffolds dropped from 97 344 Kb (L50 = 1792) to 71 491 bp (L50

= 2379). The number of remaining gaps in the assembly was 21 417, totaling 30 066 113 bp (5.05%).

Paired-end reads from the short fragment libraries were aligned back independently to this genome assembly using SMALT (map -r 0 -x -y 0.5; default alignment penalty scores). Per-scaffold depth of coverage was computed, regardless of mapping quality, using Genome Analysis Toolkit (GATK) DepthofCoverage. The mean read depth across the scaffolds resulted in $\times 66.45$. The mean read length of the mapped reads was 139.8 bp, and the corresponding k-mer coverage for the size of 25 was $\times 55.04$, which matches with the homozygous peak computed from the k-mer frequency distribution from the unassembled reads. The read depth frequencies are shown in Fig. 2B. The heterozygous/homozygous peak height (>1) in the distribution suggests that the assembly contains redundant copies of unmerged haplotypes due to the structural heterozygosity of the diploid genome of the species. To specifically deal with the heterozygosity, we introduced a step to, leniently, recognize and remove alternative heterozygous sequences. Sequences of scaffolds were aligned 1 vs all using the BLAST-like alignment tool (BLAT, [RRID:SCR.011919](#)) [57], and results were concatenated in a single file of alignments and sorted. Similar sequences were identified on the base of pairwise similarity using filterPSL utility from AUGUSTUS [58] with default parameters, and retaining all best matches to each single sequence queried against all others that satisfy minimal percentage of identity (minId = 92%) and minimal percentage of coverage of the query read (minCover = 80%). We considered as heterozygous redundant those scaffolds that showed pairwise similarity to exactly another sequence, and their depth of coverage fell in a Poisson distribution with parameters given by the heterozygous peak of the read depth distribution over all scaffolds ($\lambda = 34$) (Fig. 2B). The final step was to keep only 1 copy—the largest—of the heterozygous scaffolds among pairs with high similarity.

A preliminary assembly of the *H. impetiginosus* genome

At the end of the accuracy evaluation processes, the genome assembly had a total size of 503 308 897 bp, with gaps, in 13 206 scaffolds. The N50 of scaffolds of 1 kbp or longer was 80 946 bp (L50 = 1906), and the average size of the sequences was 38 118 bp. Using 20 kbp as an approximate value of longest plant gene length [59, 60], the percentage of scaffolds that equaled or surpassed this value in relation to the empirically determined genome size is 83%, which corresponds to over 92% of the assembly total size. Contigs generated by cutting scaffolds at each gap (of at least 25 base pairs, i.e., 25 or more Ns) produced an N50 of 40 064 bp (L50 = 3551) with an average sequence size of 19 765 bp. The remaining gaps comprised 26 447 057 bp (5.25% of the genome assembly) in 11 094 segments, with a size of 2384 ± 3167 bp. The total assembly size represents over 90% of the flow cytometry genome estimate (557 Mb) and should provide a good start to build a further improved reference genome assembly of the species using long-range scaffolding techniques such as whole-genome maps using either imaging methods [61] or contact maps of chromosomes based on chromatin interactions [62]. Table 1 summarizes the main statistics of the *Handroanthus impetiginosus* genome assembly with respect to the decisions made in the assembly process.

A reassessment of the assembly accuracy was carried out using REAPR on the final genome assembly. A total of 121 errors within a contig were still recognized, a much smaller number than previously annotated (1658 errors). Fig. 3A shows

Table 1: *Handroanthus impetiginosus* genome assembly statistics

Scaffold sequences	Allpaths-LG	Allpaths-LG/ Sspace/GapClose	Allpaths-LG/ Sspace/GapClose/Reapr
Number	57 815	16 090	13 206
Total size, without gaps, bp	469 049 393	565 959 143	476 867 120
Total size, with gaps, bp	614 626 609	586 542 612	503 314 177
Number > 10 Kbp	10 029	8602	8348
Number > 20 Kbp	6920	6791	6647
Number > 100 Kbp	1100	1709	1304
Number > 1 Mbp	2	0	0
Longest sequence, bp	1 844 569	979 053	558 523
Average size, bp	10 631	36 454	38 112
N50 length, bp	57 726	97 266	80 946
L50 count	2595	1792	1906
GC %	33.63	33.57	33.62

The final assembly for each step contains scaffolds of length 1 kbp or longer.

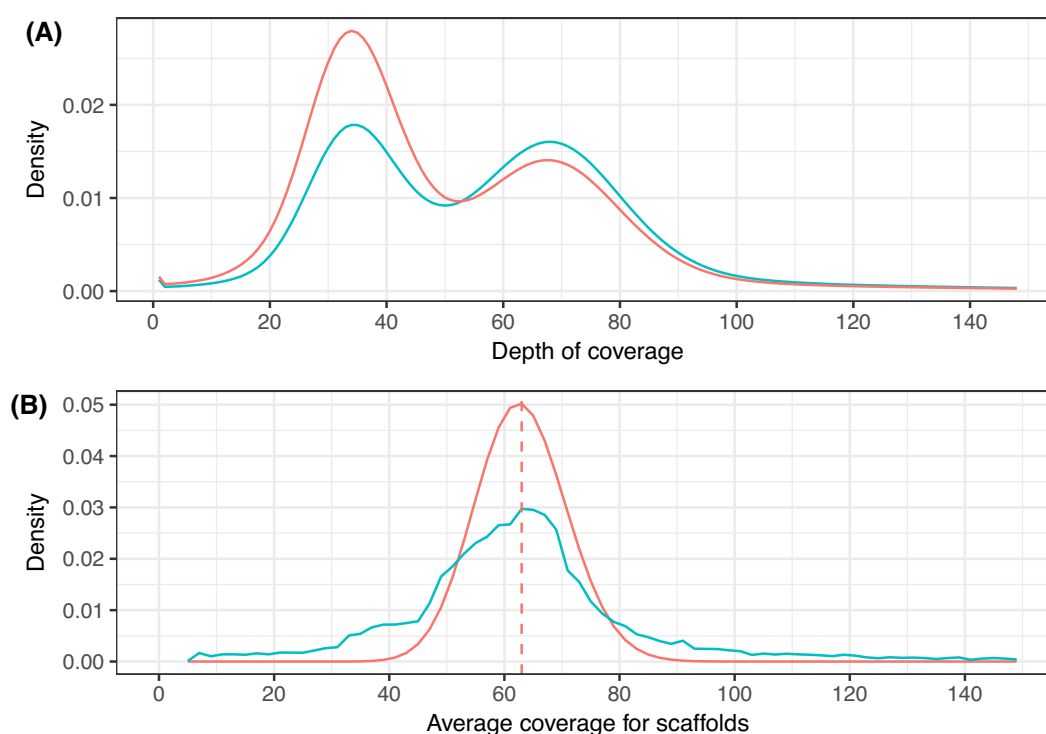


Figure 3: Depth of coverage analysis for the haplotype-reduced assembly. (A) Density plot of read depth based on mapping all short fragment reads back to the haplotype-reduced assembled sequences after identification and removal of redundant sequences due to the structural heterozygosity in the genome. (B) Density plot for average sequencing coverage per scaffold on the final assembly. The observed number of scaffolds in the final haplotype-reduced assembly and the respective read coverage (blue line) are shown in comparison with a Poisson process approximation (red line) with $\lambda = 63$, the observed average sequencing coverage in the useful read data.

the frequency distribution for the read depth computed from the paired-end read alignment to the scaffold sequences. It indicates the expected effect on the distribution in comparison with the previous more redundant assembly. The height of the heterozygous peak was successfully lowered by removing unmerged copies of the same heterozygous loci. Fig. 3B shows the relation between the observed number of scaffolds in the final assembly and their read coverage in comparison with a Poisson approximation with λ of 63, which was the observed average sequencing coverage for reads set from short fragment libraries. Loss of information due to repeat sequences is clearly a limitation of this *H. impetiginosus* assembly. Given the high rate of nonclassified consensus sequences,

we can infer that most families/subfamilies of repeats might be underrepresented.

To complement the depth of read coverage analyses, we performed additional analyses to identify the most probable causes of breaks in the assembly. We inspected contig termini defining the positions of the terminal nucleotides of each contig from the genome assembly created by cutting at each gap (of at least 1 base pair, i.e., 1 or more Ns). This analysis was developed using a protocol described elsewhere [63], and results are summarized in Fig. 4. Contig termini overlap most prominently (~50%) with regions that do not encompass any annotated feature or regions that have no depth of coverage (~15%) based on mapped reads to the assembly. It suggests that contigs end in large repeats not yet

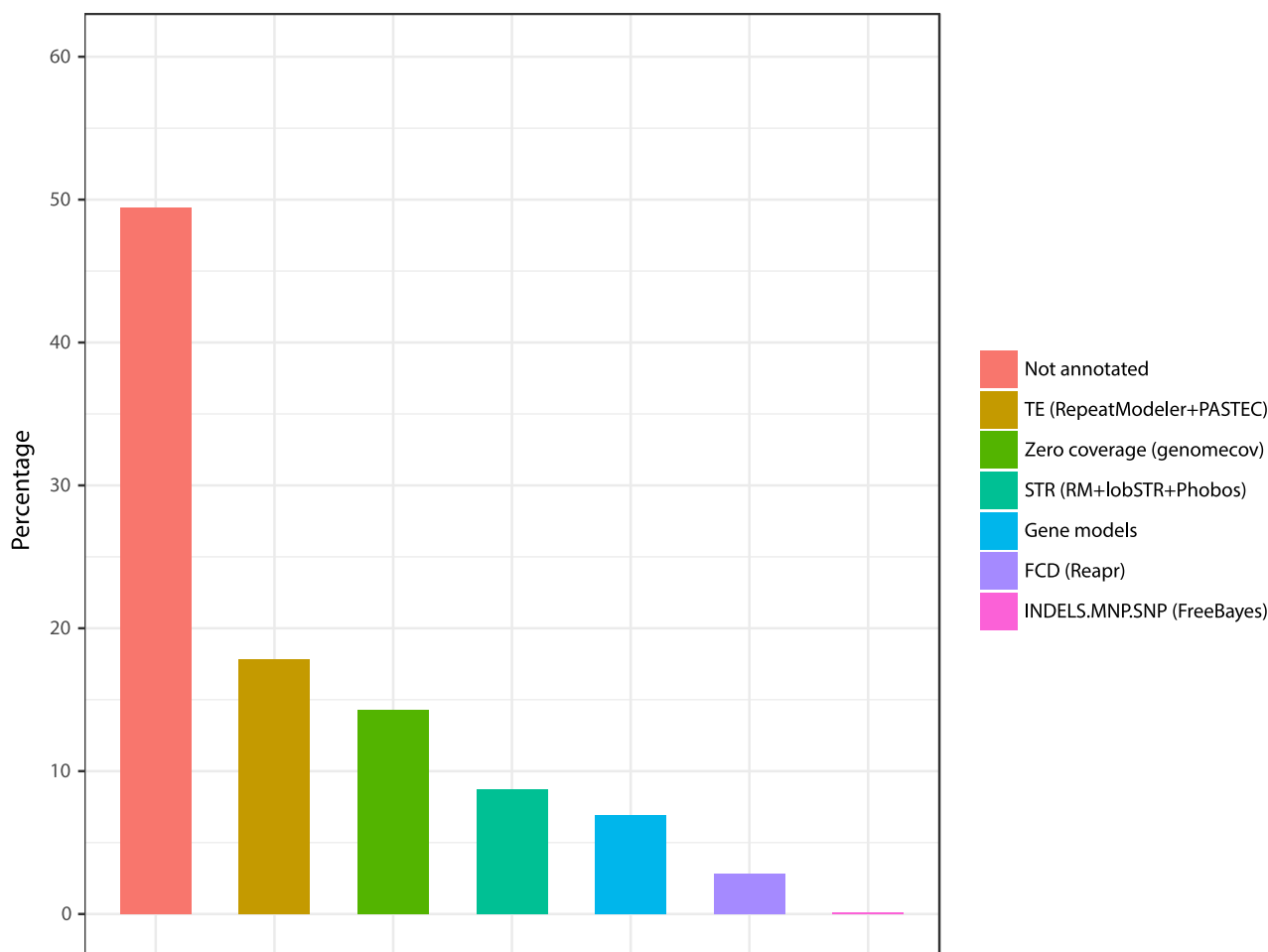


Figure 4: Contig termini analysis to investigate the possible genomic features associated with gaps in the genome assembly. Contigs were created from the genome assembly with the “cutN -n 1” command from the seqtk program, which cut at each gap (of at least 1 base pair, i.e., 1 or more Ns). The figure shows the percentage of contig termini (the position of the terminal nucleotides of each contig) intersecting with different annotations of the genome.

resolved given the inherent limitations of short-read sequence data. Another possibility is that these regions can contain low-copy young euchromatic segmental duplication with higher sequence similarity to the consensus sequence. Annotated interspersed repeats (~18%) and short tandem repeats (~9%) were the most prominently annotated features with overlap to contig ends. Less than 8% (2473 of 31 668) of annotated gene models were found to overlap contig ends, indicating that very few are likely to be interrupted in this unfinished assembly. It is a trend that was confirmed using BUSCO analysis, which reported only 3% of fragmented genes. Based on variant identification analysis with FreeBayes (FreeBayes, [RRID:SCR_010761](https://doi.org/10.1093/bioinformatics/btu170)) using read data mapped to the genome assembly, we found virtually no allelic variants located at the contigs’ end, suggesting that interruption of continuity and contiguity in the assembly is not related to differences between haplotypes.

Repetitive DNA

A total of 1608 consensus sequences (average length = 773 bp, totaling 1 281 536 bp) representing interspersed repeats in the genome assembly were found. Search for domains in these sequences with similarity to known large families of genes that could confound the identification of true repeats indicated 85 false positives in the consensus library of repeats. A further 50

sequences were annotated with predicted protein domains frequently associated with protein coding genes. These 135 sequences were wiped out from the consensus library. Most of the remaining 1473 sequences (71.1%) could not find classification in the hierarchical well-known classes of transposable elements (TEs) [64], but 16.6% could be classified Class I (retrotransposons), including 3 orders: long terminal repeats (LTR; 12.8%), long interspersed nuclear elements (LINE; 1.6%), and short interspersed nuclear elements (SINE; 2.2%); 8.4% are Class II (DNA transposons). Other categories comprised nonautonomous TEs: TRIM (0.4%) and miniature inverted-repeat transposable elements (MITE; 3.5%). Unknown nonclassified sequences in the consensus library cover a wide range of sequence sizes, from 42 bp up to 5987 bp (average = 345 bp, median = 503 bp). The 1473 sequences representing interspersed repeats in the consensus repeat library were used to mask the genome with RepeatMasker. The masked fraction of the genome assembly comprised 155 348 349 bp, i.e., 30.9% of the total assembled genome of 503 Mbp. Remarkably, if we add to these ~155 Mbp the 54 Mbp of noncaptured base pairs in the assembly when considering the empirically determined genome size (557–503), the repetitive fraction of the genome approximates 37.5% (209 Mbp out of 557 Mbp). This is within the expected range (36.6%–38.0%) for the repetitive fraction of the genome estimated from the read set using k-mer profiling approaches.

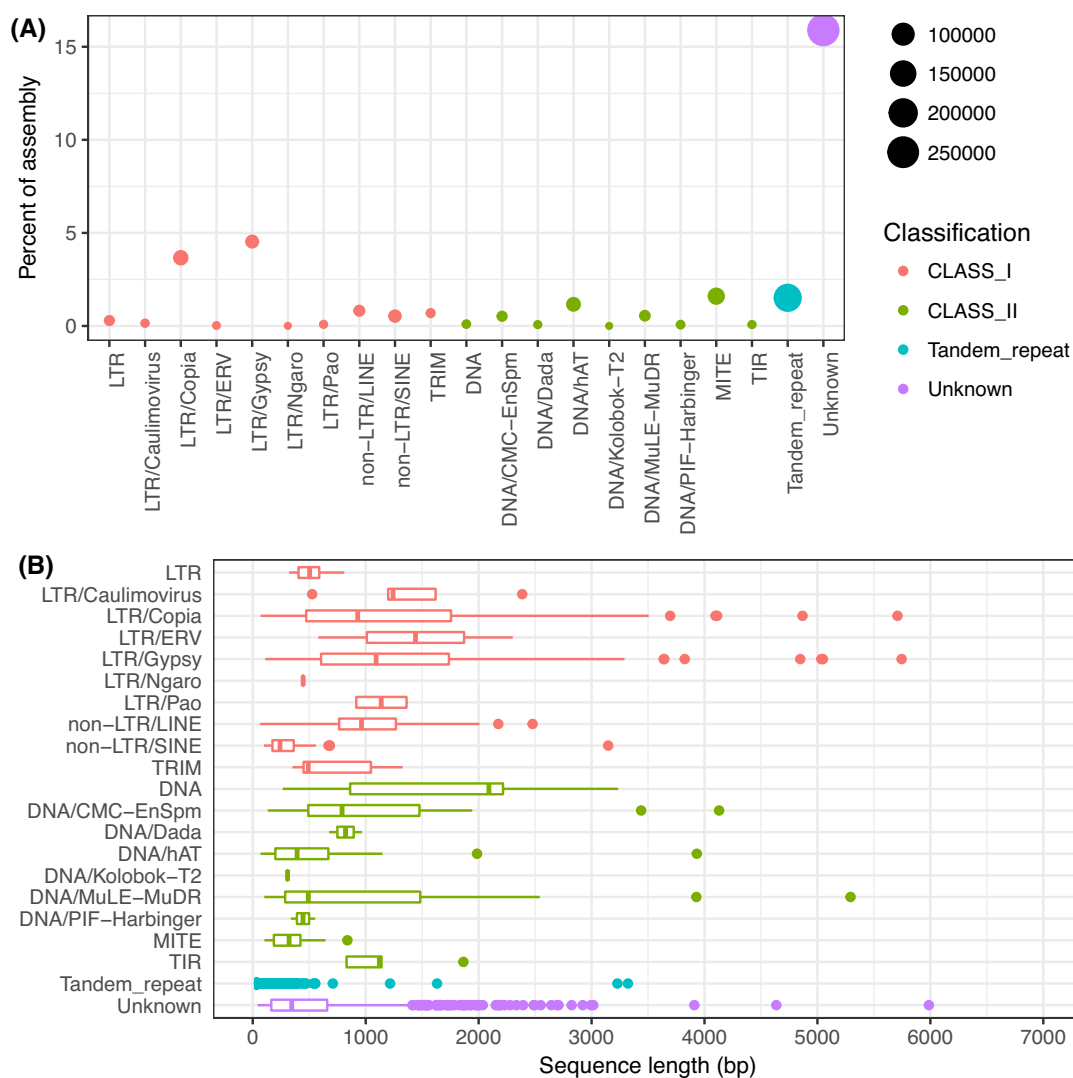


Figure 5: Repeat content of the *H. impetiginosus* genome assembly. (A) The density of interspersed and tandem repeat as percentage of the assembly. The size of the circles represents the number of copies in the assembly for each family of repeats. (B) Distribution of sizes of the consensus sequences for repeat families identified using *de novo* and homology methods for repeat characterization.

More than 50% of the masked bases in the assembly, or 80 Mbp, came from nonclassified sequences in the consensus library. In the well-known repeats, retrotransposons are the most abundant class in the assembly, comprising 50 Mbp (~one-third of the masked bases), with prominence of LTR/Gypsy (~23 Mb) and LTR/Copia (18 Mb) families of repeats. DNA transposons and nonautonomous orders of transposons masked 12 Mbp and 11 Mbp (~one-sixth of the masked bases), respectively, highlighting the prominence of DNA/hAT families of class II and MITE (Fig. 5). Simple sequence repeats (SSRs) detection using RepeatMasker identified a total of 182 115 microsatellites with a density of 2.76 kb per SSR in the genome assembly. This density corroborates the general finding that the overall frequency of microsatellites is inversely related to genome size in plant genomes [65]. This SSRs density in *H. impetiginosus* (genome size of 557 Mbp/SSR density of 362 per Mbp) is higher than in larger plant genomes such as those of maize (1115 Mbp/163 SSRs per Mbp), *S. bicolor* (738 Mbp/175 per Mbp), *G. raimondii* (761 Mbp/74.8 per Mbp) [66], but lower than densities in smaller genomes such as those of *A. thaliana* (120 Mbp/418 per Mbp), *Medicago truncatula* (307 Mbp/495 per Mbp), and *C. sativus* (367 Mbp/552 per Mbp) [67]. Different

SSR motifs ranging from 1 to 6 bp showed that the di-nucleotide repeats were the most abundant repeats, followed by the mono- (Fig. S3A). The frequency of SSR decreased with increase in motif length (Supplementary Fig. S3B), which is a trend usually observed both in monocots and dicots [67].

Transcriptome assembly and gene content annotation and analysis

A single run of Illumina HiSeq 2500 sequencing, from a pool of RNA samples, generated nearly 148 million paired-end reads. After adapter removal, trimming, and coverage normalization, 55.2 million high-quality reads (38%) were used to assemble the transcriptome using *de novo* (Trinity and SOAP-Trans-denovo transcripts combined with the EvidentialGene pipeline) and genome-guided methods (PERTRAN). The PASA pipeline was used to integrate transcript alignments to the genome assembly from these sets of sequences, generating 54 320 Expressed Sequence Tag (EST) assemblies representing putative protein-coding loci in the genome assembly. Loci were identified by the assembled transcript alignments using BLASTX [36] and

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

	Genes	Transcripts	Exons	Introns
Number	31 688	35 479	154 209	122 521
Average number/gene	–	1.12	4.87	3.87
Average length	3129	3342	285	445
N50 length	4421	4643	477	839
%GC	38.38	38.22	42.60	32.83
%N	0.43	0.43	0.00	0.29

Table 3: The distribution of the minimal introns (53–125 bp) and the minimal-intron-containing genes—as the number of genes with at least 1 minimal intron—from selected plant species in comparison with the *H. impetiginosus* genome assembly

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

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 with repeats, as described in the “Methods.” The number of predicted messenger RNA (mRNA) transcripts was 35 479.

Structural features of the gene content are shown in Tables 2 and 3. The average number of exons per gene was ~5, and its average length was 285 bp. The average number of introns per gene was ~4, and its average length was 445 bp. The GC content is significantly different between exons and introns (t test $P < 0.0001$). Coding sequences have ~43% of GC, while introns have less, with ~33% (Table 2). GC content tends to be higher in coding (exonic) than noncoding regions [68], which may be related to gene architecture and alternative splicing [69–71]. A comparison of the gene feature parameters, such as number and length (Fig. S4A), was carried out between *H. impetiginosus* and *Erythranthe guttata*, another plant in the order Lamiales (Asterids), the model plant *A. thaliana* and the model tree *P. trichocarpa* (Rosids). As depicted in the frequency histograms, the exon parameters are stable among these species (Fig. S4B). For the introns (Fig. S4C), frequency histograms have a sharp peak around 90 bp and a larger peak that is much lower in density. There is a small intron size variability from species to species in the distributions, especially for larger introns, which rarely go beyond 10 000 bp. The intron length distributions in these 4 species are similar to those observed in lineages that are late in the evolutionary time scale, such as plants and vertebrates [72]. The sharp peak in the distributions at their “minimal intron” size is supposed to affect function by enhancing the rate at which mRNA is exported from the cell nucleus [73, 74]. In the model plant *A. thaliana*, a minimal intron group was previously defined [73] as anything that lies within 3 standard deviations of the optimum peak at 89 ± 12 bp (53–125 bp). According to this definition, Table 3 summarizes the distribution of the minimal intron among genes of *H. impetiginosus* and other selected plant

species in the Asterids and Rosids lineages. We have calculated the percentages of minimal introns out of the total introns and the fraction of minimal-intron-containing genes with at least 1 minimal intron. Computed values were similar between *H. impetiginosus* and those of selected species with higher numbers of large introns (smaller minimal intron peak) but were more distinctive with those species such as *A. thaliana* and *E. guttata*, in which the number of large introns was lower (larger minimal intron peak). This is thought to be a general trend and was also observed in previous work [73]. These comparative analyses about the structural properties of the predicted genes indicate that the genome assembly of *H. impetiginosus* contains highly accurate gene structures.

To further validate the gene content annotation, we used the transcript assemblies and selected plant proteomes to inspect if these sequences could align in their entirety to the genomic sequence. Out of the 31 668 primary mRNA transcripts (considering only the longest one when isoforms were predicted) in the genome, 11 488 have 100% of their coding DNA sequence (CDS) covered by EST assemblies. The remaining 20 054 transcripts have either a minimum of 80% of their CDS covered by EST assemblies or a Cscore ≥ 0.5 . From these latter, the encoded putative peptides have excellent sequence similarity support from BLASTP comparisons with dicot species *Erythranthe guttata* (5224 genes), *Sesamum indicum* (4625 genes), potato or tomato (2777 genes), soybean (1484 genes), and the poplar tree (1424 genes), reflecting the taxonomic relationship between *H. impetiginosus* and these other related dicots. Gene model support was also found from more distantly related dicots (1826 genes) and monocots (1042 genes). Altogether, 31 048 gene models (98%) show well-supported similarity hits to other known plant protein sequences. An additional 517 predicted protein sequences did not produce hits, and 103 sequences produced ambiguous hits from nontarget species or represent possible contaminants in the assembly, such as endophytic fungi (ascomycetes, 42 sequences; basidiomycetes, 17 sequences). Fig. 6A summarizes the main finding regarding the similarity analyses with known proteins.

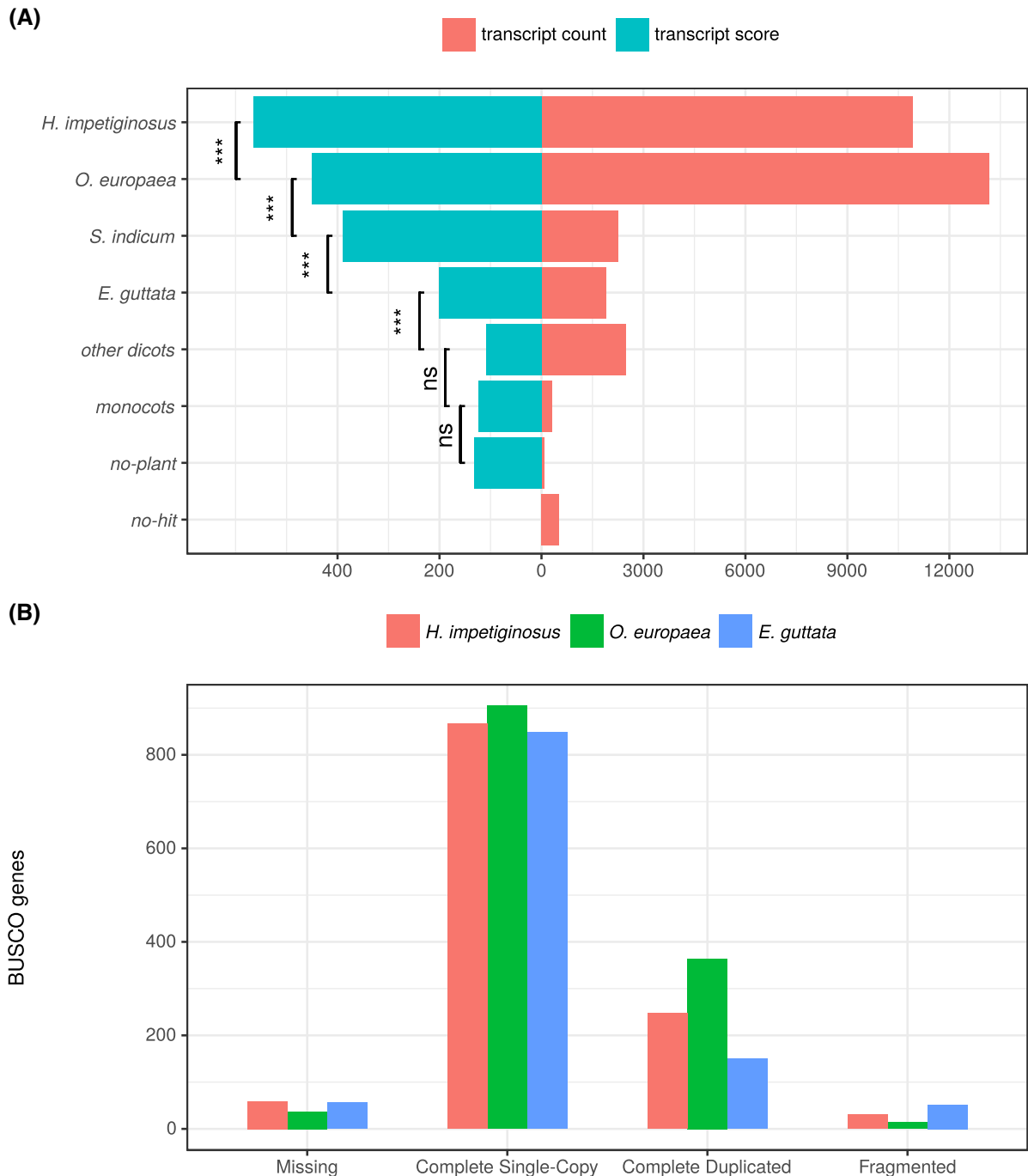


Figure 6: Transcriptome quality assessment (A) similarity search of *H. impetiginosus* putative peptides against source database of plant protein sequences using BLASTP algorithm (e -value = $1e-6$). Transcript count means the number of peptides of *H. impetiginosus* with the best hit against the source database using bit-score and grouping results by taxon name. Transcript score corresponds to the average bit-score overall hits for each group using the best hit. We ordered taxon groups by their average bit-score overall hits and used Welch's t test to compare the distributions of bit-score hits between 2 adjacent groups with P -values <0.01 (ns = nonsignificant; *significant). (B) Completeness of the expected gene space of the genome assembly, estimated with BUSCO. The estimates were compared with genome annotations for other lamids, *Erythranthe guttata* and *Olea europaea*.

BUSCO (BUSCO, [RRID:SCR.015008](https://doi.org/10.1093/bioinformatics/btu015)) [75] single-copy gene plant profiles were used to estimate completeness of the expected gene space as well as the duplicate fraction of the genome assembly. Out of the 956 profiles searched on the assembly, 59 (6.1%) were reported missing and 30 (3.1%) returned fragmented. From the profiles with a complete match to the assembly, 867 (90.7%) were reported as single-copy and 247 (25.8%) were found

to be completely duplicated. We benchmarked our results by searching the BUSCO profiles on the genomes of other lamids, *Erythranthe guttata* and *Olea europaea*. In *E. guttata*, the analysis reported a completeness level of 88% (848 single-copy profiles with a complete match), while there were 52 fragmented genes (5.4%). In *O. europaea*, the completeness level was 94% (905 complete single-copy profiles), and there were only 14 fragmented

genes (1.4%). A summary of the BUSCO analysis is presented in Fig. 6B.

Databases for gene ontology (GO) annotation are rich resources to describe the functional properties of experimentally derived gene sets. To explore relationships between the GO terms in the *H. impetiginosus* and related, well-curated genomes, we used WEGO [76] to perform a genome-wide comparative analyses among broad functional GO terms with other lamids. The *P*-value of the Pearson chi-square test was considered to indicate significant relationships between the proportions of genes of each GO term in these 2 datasets and to suggest patterns of enrichment (Figs S5 and S6). These analyses revealed several GO terms in which the proportion of genes in the 2 compared species were related. For the terms in which the comparison did not indicate a significant relationship of gene proportions between the 2 datasets, the compared GO terms suggested enrichments in *H. impetiginosus* for GO terms involved in metabolic processes and catalytic activity in comparison with *E. guttata* and *O. europaea*.

The central role of enzymes as biological catalysts is a well-studied issue related to the chemistry of cells [77]. An important feature of most enzymes is that their activities can be regulated to function properly to comply with physiological needs of the organism. We observed that the GO term for enzyme regulatory activity encompasses a higher proportion of genes in *H. impetiginosus* than in the 2 other lamids, albeit the difference did not reach significance in *E. guttata*. Research in *Arabidopsis*, a herbaceous plant, has found little connectivity between metabolites and enzyme activity [78]. In comparison with *Arabidopsis* broader GO terms, *H. impetiginosus* showed, as discussed above, enrichment for the proportion of genes assigned to the metabolic process (49.1% > 47.4%; *P* = 0.002) and catalytic activity (46.2% > 42.9%; *P* = 0). The proportion of genes for enzyme regulatory activity was also higher in *H. impetiginosus* than *A. thaliana*, though not statistically significantly (*P* = 0.083). Investigations into whether and how metabolic process and enzyme activities relate and how it could influence the known richness of metabolites for forest trees of the mega-diverse tropical biomes, particularly in the genus *Tabebuia* and *Handroanthus*, shall be an interesting issue for future molecular and chemistry studies.

Benchmarking the genome assembly of *H. impetiginosus*

Based on current standards for plant genome sequence assembly [60, 79, 80], we have provided a quality assembly of high future utility. To support functional analyses, we classified the gene models into high-confidence and low-confidence groups. Out of the 31 688 protein-coding loci annotated in the genome assembly, 28 603 (90%) produced high-confidence gene models (Supplementary File S1). This subset contains approximately the same number of genes reported in less fragmented genome assemblies for other lamids. *E. guttata* (2n = 28) reports 28 140 protein-coding genes [81]; *O. europaea* (2n = 46) has 56 349 protein-coding genes [82], but its genome has likely undergone a whole-genome duplication event. Most of *Tabebuia* and *Handroanthus* species studied so far have 2n = 40 [22]. The fraction of gene duplicates in the BUSCO analysis (see Fig. 5B) was intended to estimate the level of redundancy in the genome assembly. We benchmarked our results by searching the completed duplicated BUSCO profiles in the genomes of *E. guttata* and *O. europaea*. In the first, we found them to be 15% (150 out of 956), while in the latter the duplicated profiles were 38% (364 out of 956). In these 3

lamids, it appears that the frequency of small- and large-scale duplications, such as (paleo)polyploidy, can explain the differences in the number of annotated genes and levels of gene duplication (*E. guttata* <= *H. impetiginosus* << *O. europaea*). It suggests that the *H. impetiginosus* genome has not undergone a recent whole-genome duplication event, although a deeper analysis of this question, beyond the scope of this study, remains open.

Our genome assembly metrics were benchmarked against comparable genome assemblies of other highly heterozygous forest tree genomes (File S2 and Fig. S7). The *H. impetiginosus* assembly has 503 Mbp in 13 206 scaffolds ≥ 2 kbp, representing over 90% of the flow cytometry estimated size (557 Mb). For *Quercus robur*, the assembly had 17 910 scaffolds ≥ 2 kbp with scaffolds N50 of 260 kbp, but corresponding to 1.34 Gbp, i.e., 81% larger than the expected 740-Mbp genome, which is clearly an undesirable result [83]. For *Quercus lobata*, with a genome size of 730 Mbp, 2 assemblies were provided: a haplotype-reduced assembly, with 40 158 contigs totaling 760 Mb, N50 of 95 kbp, and a more complete version for gene models, containing 94 394 scaffolds ≥ 2 kbp, totaling 1.15 Gbp, with an N50 of 278 kbp [48]. Despite our lower NG50/N50 scaffold length <100 kbp, the *H. impetiginosus* assembly has a large (60%) percentage of scaffolds ≥ 20 kbp. This value is higher than the reported values for *Quercus lobata* v0.5 (53%), *Quercus lobata* v1.0 (51%), and *Quercus rubra* (48%), even if those assemblies had higher NG50/N50 scaffold lengths. Finally, contig termini analysis has found virtually no allelic variants located at contig ends, suggesting that interruption of continuity and contiguity in the assembly is not related to differences between haplotypes. This genome assembly for *Handroanthus impetiginosus* will thus be useful for variant calling, one of the main future objectives for generating this resource.

Genome-guided exploration of specialized metabolism genes of quinoid systems

Aside from their highly valued wood, *H. impetiginosus* and other Ipê species are also known for their medicinal effects. Extracts from their bark and wood have many ethnobotanical uses: against cancer, malaria, fevers, trypanosomiasis, fungal and bacterial infections, and stomach disorders [84, 85]. The wood extracts have also been demonstrated to have anti-inflammatory effects [86] [87]. The main bioactive components isolated from the Pink Ipê are Lapachol and its products [88], which are naphthoquinones derived from the o-succinylbenzoate (OSB) pathway [89]. Lapachol is also responsible for the well-known high resistance of the Ipê wood against rotting fungi and insects [90]. In addition, naphthoquinones are aromatic substances with ecological importance for the interaction of plants with other plants, insects, and microbes [89]. 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 found all the important known enzymes that lead to the biosynthesis of lapachol (Fig. 7). Unfortunately, however, the last 2 steps of the lapachol biosynthesis pathway still constitute unidentified enzymes [89]. For comparative purposes, we downloaded the annotation file of 5 other species from the Phytozome database. 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. However, 3 exceptions were found. *H. impetiginosus* has 5 genes encoding the enzyme that converts chorismate to isochorismate, the first step in the OSB pathway. Two other steps found to have

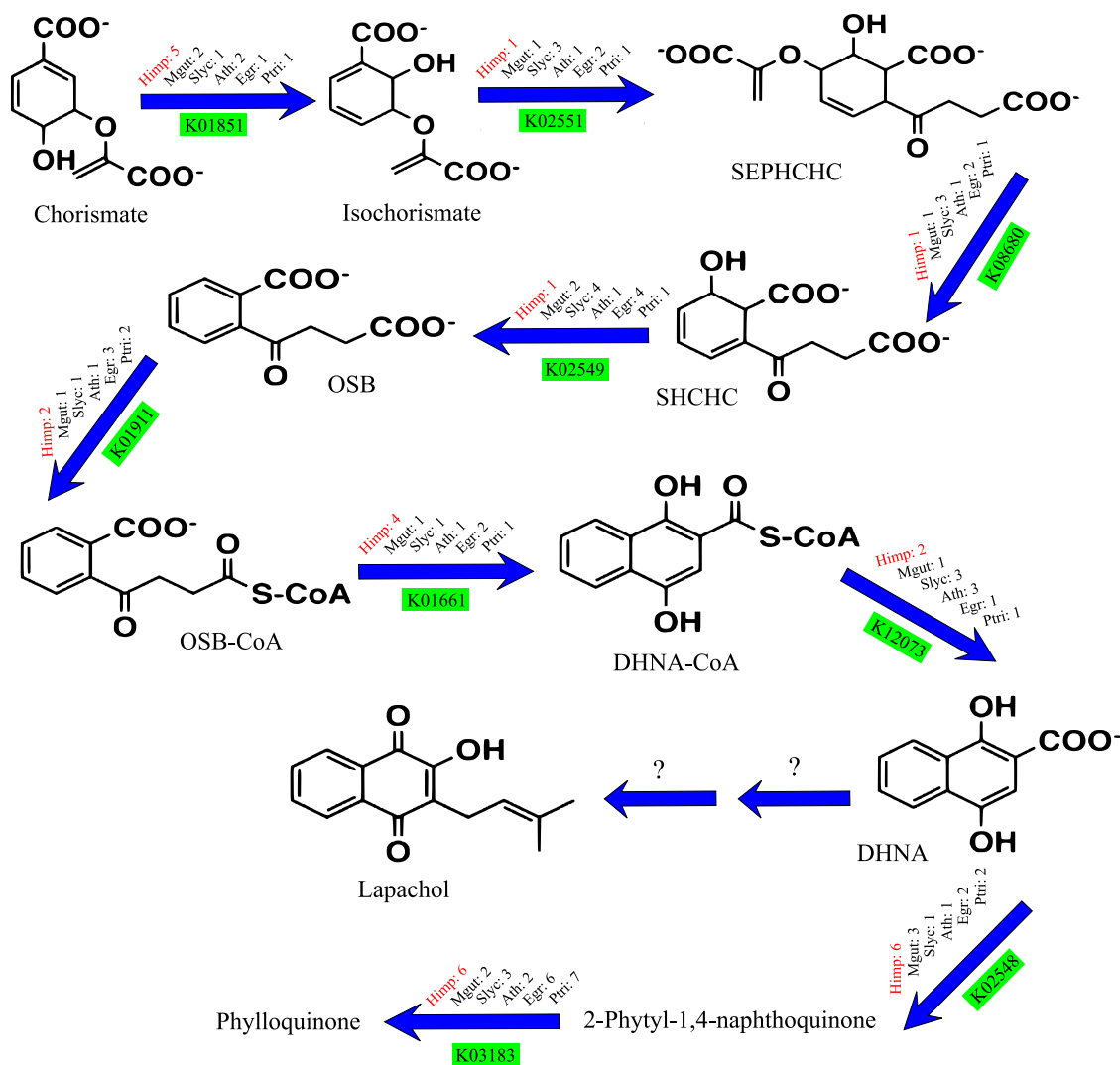


Figure 7: Genes of the biosynthetic pathway of specialized quinoids. O-succinylbenzoate (OSB) pathway depicting the number of *H. impetiginosus* (Himp) annotated genes for the known enzymes that lead to the biosynthesis of the naphthoquinones, including lapachol. For comparison, it also shows the numbers of genes for the closely related *Mimulus guttatus* (Mgut), *Solanum lycopersicum* (Slyc), for the model *Arabidopsis thaliana* (Ath), and for the tree species *Eucalyptus grandis* (Egr) and *Populus trichocarpa* (Ptri). The pathway was modified from Widhalm and Rhodes [89].

relatively more genes in *H. impetiginosus* are the ones that lead to the synthesis of 1,4-Dihydroxy-2-naphthoyl-CoA and of 2-Phytl-1,4-naphthoquinone. The availability of sequences for these genes may open new avenues for biotechnological products and for a better understanding of their ecological roles.

Re-use potential

We have reported a well-curated but still unfinished genome assembly for *Handroanthus impetiginosus*, a highly valued, ecologically keystone tropical timber and a species rich in natural products. The fragmentation of this preliminary assembly might be still be limiting for deeper insights of whole-genome comparative analyses or studies of genome evolution [91], although we think that such studies may be carried out using this assembly at least at the gene level or gene-family level. Nevertheless, the broad validation performed provides a useful genomic resource for genetic and functional analysis, including, but not limited to, downstream applications such as variant calling, molecular markers development, and

functional studies. Extensive documentation of quality throughout the assembly process was provided showing that acceptable continuity was reached and that the fragmentation of the final sequence mostly derives from loss of information on high-copy families of long interspersed repeats or the presence of low-copy segmental duplications likely recently evolved with higher sequence similarity to the consensus sequence. Certainly, there are still inaccuracies at the base and assembly levels, but all efforts were made to deliver results to the end user with the appropriate documentation, making this initial read set, sequence, and annotations a primary and reliable starting grounds for further improvement.

We have documented in detail the main features of the reported assembly. The total assembly size of scaffolds ≥ 2 kbp in length is 90% of the flow cytometry-determined genome size, a remarkable accomplishment, we believe, given the anticipated difficulties in assembling such a repetitive and highly heterozygous diploid genome based exclusively on short-read sequencing. The percentage of base pairs in scaffolds with ≥ 20 kbp is 83% (461 Mbp of 557 Mbp) of the empirically determined

genome size, which corresponds to 92% of the assembled total size (461 Mbp of 503 Mbp). Using 20 kbp as an approximate value of the longest plant gene length, this result shows that 60% of the assembly is accessible for reliable gene annotation. Furthermore, the N50/NG50 (41 kbp/34 kbp) contig length is longer than 30 kbp, which has been suggested to be an adequate minimum threshold for high utility of a genome assembly [79]. The percentage of documented gaps in scaffolds is only 5.3%, and the few misassembled signatures present in the assembly were fully documented based on acceptable metrics such as fragment coverage distribution error (FCD error). Less than 8% (2473 of 31 668) of annotated gene models were found to overlap contig ends, indicating that very few are likely to be interrupted in this unfinished assembly. No allelic variants were found at contig ends, suggesting that interruption of continuity and contiguity in the assembly is not related to differences between haplotypes, therefore providing a valuable resource for variant calling and functional analysis. More than 86% (27 380 of 31 668) of the gene models represented in the assembly have external evidential support measured by PASA-validated EST alignments from RNA-Seq or high-coverage alignments with known plant proteins (>90% coverage). Furthermore, 80% (25 369 of 31 668) of transcripts have conceptual translation that contains protein domain annotation, excluding those associated to TEs. Finally, a summary of BUSCO analysis indicates that the detected number of plant single-copy orthologs represents 90% of the searched profiles (867 of 956), while only 6% are missing and 3% are fragmented.

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

Availability of supporting data

Sequences for the genome and assembly, along with gene content annotation and the raw sequencing reads, have been deposited into GenBank, BioProject PRJNA324125. This Whole Genome Shotgun (WGS) project has been deposited at DDBJ/ENA/GenBank under the accession NKXS00000000. The version described in this paper is version NKXS01000000. BioSample for WGS is SAMN05195323, and the corresponding SRA run accessions are SRR3624821–SRR3624825. BioSample for RNA-Seq is SAMN07346903, with SRA run accession SRR5820886.

Supporting data and summary outputs for the main analyses in this Data Note are available via the GigaScience repository, GigaDB [92]. The Perl script that automated the read set from mate-pair sequencing preprocessing (TrimAdaptor.pl) was uploaded to GigaDB under permission of the original authors at the High-Throughput Sequencing and Genotyping Center Unit of the University of Illinois Urbana-Champaign.

Additional file

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

Figure S1: Flow cytometry results of the sequenced tree UFG-1 of *H. impetiginosus*. Flow cytometry estimate of the nuclear DNA content was carried out using young leaf tissue on a BD Accuri C6 Plus personal flow cytometer. *Pisum sativum* (genome size 9.09 pg/2C or ~4380 Mb/1C) was used as standard for comparison (M2). The estimate of nuclear DNA content for *H. impetiginosus* (M1) averaged over 10 readings was 1.155 pg/2C or 557.3 ± 39 Mb/1C.

Figure S2: Overview of the analytical pipeline with the bioinformatics steps and tools employed for genome (black arrows) and transcriptome assembly (red arrows), and for gene prediction and annotation (blue arrows). Bioinformatics programs are indicated in italic, blue, and the main file formats in red. The input sequences are highlighted in yellow boxes and the main products in green.

Figure S3: Distribution and characterization of simple sequence repeats in *Handroanthus impetiginosus* genome. (A) Histogram of different motifs ranging from 1 to 6 bp. (B) Distribution of the simple sequence repeat length detected in the genome assembly.

Figure S4: Comparison of the gene feature parameters, such as number and length, between *H. impetiginosus* and the other selected dicot plant across distinct lineages of Rosids (*A. thaliana* and *P. trichocarpa*) and Asterids (*E. guttata* and *S. lycopersicum*). Frequency histograms are shown according to the whole-genome gene content annotation for (A) the complete predicted gene structure, (B) exons, and (C) introns. Dashed vertical lines are the average lengths for the gene features.

Figure S5: Histograms for Gene Ontology broader term annotations in the *H. impetiginosus* genome assembly. Terms for the Biological Process ontology were summarized with WEGO using the second tree level setting. The Pearson chi-square test was applied to indicate significant relationships between *H. impetiginosus* and the lamid *Erythranthe guttata* regarding the number of genes (at $\alpha \geq 5\%$). (A) Terms displaying a remarkable relationship between the 2 datasets; (B) terms with a significant difference between the 2 datasets.

Figure S6: Same as Fig. S6 but showing comparison between numbers of genes assigned to GO broader terms for *H. impetiginosus* and the lamid *Olea europaea*.

Figure S7: Sequence length distribution from the assemblies of *H. impetiginosus* and the other 2 highly heterozygous trees of the genus *Quercus*. Figure shows density plots for the size of scaffolds 2 kbp or longer in the 3 assemblies. Contigs metrics were computed by cutting at each gap (of at least 25 base pairs, i.e., 25 or more Ns). Scaffolds and contigs length were plotted using the common logarithm to respond to skewness toward large values.

File S1: Evidence adopted to support protein-coding loci identification and assignment in the *H. impetiginosus* genome assem-

bly. Two qualifiers—high confidence and low confidence—were added to the locus based on the reported evidence.

File S2: Genome assembly metrics from the assemblies of *H. impetiginosus* and the other 2 highly heterozygous trees of the genus *Quercus*. Comparison between metrics was based on the `assemblathon_stats` script part of the `assemblathon2-analysis` package (<https://github.com/ucdavis-bioinformatics/assemblathon2-analysis>). Metrics were computed for scaffolds 2 kbp or longer in length. Genomic sequences in scaffolds for *Quercus lobata* was obtained from <https://valleyoak.ucla.edu/genomicresources/> (accessed on 20 September 2017). For *Quercus rubra*, genomic sequences in scaffolds were downloaded from the European Nucleotide Archive repository, accessions LN776247–LN794156.

Abbreviations

BLASTP: Basic Local Alignment Search Tool for Proteins; BLAT: BLAST-like alignment tool; CDS: coding DNA sequence; EC: Enzyme Commission Number; EST: Expressed Sequence Tag; GATK: Genome Analysis Toolkit; GO: Gene Ontology; LINE: long interspersed nuclear elements; LTR: long terminal repeats; MITE: miniature inverted-repeat transposable elements; mRNA: messenger RNA; PASA: Program to Assemble Spliced Alignment; REAPR: Recognition of Errors in Assemblies using Paired Reads; SINE: short interspersed nuclear elements; SNP: single nucleotide polymorphism; SSPACE: SSAKE-based Scaffolding of Pre-Assembled Contigs after Extension; TE: transposable element.

Competing interests

The authors declare that they have no competing interests.

Funding

This work was supported by competitive grants from CNPq to R.G.C. (project no. 471366/2007–2, Rede Cerrado CNPq/PPBio project no. 457406/2012–7, and Procad/Capes project # 88881.068425/2014–01), to E.N. (CNPq Proc. 476709/2012–1), and to D.G. (PRONEX FAP-DF Project Grant “NEXTREE” 193.000.570/2009). R.G.C. and D.G. have been supported by productivity grants from CNPq, which we gratefully acknowledge. O.B.S.Jr. has been supported by an EMBRAPA doctoral fellowship and was an Affiliate Researcher at Lawrence Berkeley National Laboratory (LBNL), Berkeley, California, at the time of this research.

Author contributions

O.B.S.Jr. performed sequence data analysis and genome assembly and, together with E.N., carried out transcriptome and protein-coding gene annotation. R.C. and D.G. conceived the project, collected samples, extracted genomic DNA and RNA, carried out flow cytometry analysis, and supervised the project. All authors were involved in discussions, writing, and editing. All authors read and approved the final manuscript.

Acknowledgements

O.B.S.Jr. thanks D. M. Goodstein and the members of the Phytozome team at the LBNL/Joint Genome Institute (JGI) for their valuable help and support in working with the JGI pipelines for

genomic research. We also thank Dr. Gabriela Ferreira Nogueira and André Luis X. de Souza for their help with flow cytometry analysis.

References

1. Goodstein DM, Shu SQ, Howson R et al. Phytozome: a comparative platform for green plant genomics. *Nucleic Acids Res* 2012;**40**(D1):D1178–86.
2. Kang YJ, Lee T, Lee J et al. Translational genomics for plant breeding with the genome sequence explosion. *Plant Biotechnol J* 2016;**14**(4):1057–69.
3. Bevan M, Walsh S. The Arabidopsis genome: a foundation for plant research. *Genome Res* 2005;**15**(12):1632–42.
4. Morrell PL, Buckler ES, Ross-Ibarra J. Crop genomics: advances and applications. *Nat Rev Genet* 2012;**13**(2):85–96.
5. Varshney RK, Glaszmann JC, Leung H et al. More genomic resources for less-studied crops. *Trends Biotechnol* 2010;**28**(9):452–60.
6. Myburg AA, Grattapaglia D, Tuskan GA et al. The genome of *Eucalyptus grandis*. *Nature* 2014;**510**(7505):356–62.
7. Tuskan GA, DiFazio S, Jansson S et al. The genome of black cottonwood, *Populus trichocarpa* (Torr. & Gray). *Science* 2006;**313**(5793):1596–604.
8. Neale DB, Wegrzyn JL, Stevens KA et al. Decoding the massive genome of loblolly pine using haploid DNA and novel assembly strategies. *Genome Biol* 2014;**15**:R59.
9. Nystedt B, Street NR, Wetterbom A et al. The Norway spruce genome sequence and conifer genome evolution. *Nature* 2013;**497**(7451):579–84.
10. Moghe G, Last R. Something old, something new: Conserved enzymes and the evolution of novelty in plant specialized metabolism. *Plant Physiol* 2015;**169**(3):1512–23.
11. Stone R. Lifting the veil on Traditional Chinese Medicine. *Science* 2008;**319**(5864):709–10.
12. Chappell J, DellaPenna D, O'Connor S. Specific aims for medicinal plant genomics resource. *Med Plants Genomics Resour* 2017, <http://medicinalplantgenomics.msu.edu/>. Accessed 27 September, 2017.
13. Brousseau L, Tinaut A, Duret C et al. 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 et al. Development of genomic tools in a widespread tropical tree, *Symphonia globulifera* L.f.: a new low-coverage draft genome, SNP and SSR markers. *Mol Ecol Resour* 2017;**17**(4): 614–30.
15. Cadena-González A, Sorensen M, Theilade I. Use and valuation of native and introduced medicinal plant species in Campo Hermoso and Zetaquira, Boyacá, Colombia. *J Ethnobiol Ethnomed* 2013;**9**(1):23.
16. Bodker G, Bhat KKS, Burley J et al. Medicinal plants for forest conservation and health care. *Food Agricult Org UN* 1997, <http://www.fao.org/3/a-w7261e.pdf>. Accessed 27 September, 2017.
17. Braga AC, Reis AMM, Leoi LT et al. Development and characterization of microsatellite markers for the tropical tree species *Tabebuia aurea* (Bignoniaceae). *Mol Ecol Notes* 2007;**7**(1):53–56.
18. Liu B, Shi Y, Yuan J et al. Estimation of genomic characteristics by analyzing k-mer frequency in de novo genome projects. *arXiv.org* 2013. arXiv:13082012.

19. Schulze M, Grogan J, Uhl C et al. Evaluating ipe (*Tabebuia*, Bignoniaceae) logging in Amazonia: sustainable management or catalyst for forest degradation? *Biol Conserv* 2008;**141**(8):2071–85.
20. Inagaki R, Ninomiya M, Tanaka K et al. Synthesis and cytotoxicity on human leukemia cells of furonaphthoquinones isolated from *tabebuia* plants. *Chem Pharmaceut Bull* 2013;**61**(6):670–3.
21. Park BS, Kim JR, Lee SE et al. Selective growth-inhibiting effects of compounds identified in *Tabebuia impetiginosa* inner bark on human intestinal bacteria. *J Agricult Food Chem* 2005;**53**(4):1152–7.
22. Collevatti RG, Dornelas MC. Clues to the evolution of genome size and chromosome number in *Tabebuia alliance* (Bignoniaceae). *Plant Syst Evol* 2016;**302**(5):601–7.
23. Aronesty E. Comparison of sequencing utility programs. *Open Bioinformatics J* 2013;**7**:1–8.
24. Langmead B, Trapnell C, Pop M et al. Ultrafast and memory-efficient alignment of short DNA sequences to the human genome. *Genome Biol* 2009;**10**(3).
25. Marcais G, Kingsford C. A fast, lock-free approach for efficient parallel counting of occurrences of k-mers. *Bioinformatics* 2011;**27**(6):764–70.
26. Vurture GW, Sedlazeck FJ, Nattestad M et al. GenomeScope: fast reference-free genome profiling from short reads. *Bioinformatics* 2017;**33**(14):2202–20.
27. Gnerre S, MacCallum I, Przybylski D et al. High-quality draft assemblies of mammalian genomes from massively parallel sequence data. *Proc Natl Acad Sci U S A* 2011;**108**(4):1513–8.
28. Flutre T, Duprat E, Feuillet C et al. Considering transposable element diversification in de novo annotation approaches. *PLoS One* 2011;**6**(1).
29. Wicker T, Sabot F, Hua-Van A et al. A unified classification system for eukaryotic transposable elements. *Nat Rev Genet* 2007;**8**(12):973–82.
30. Hoede C, Arnoux S, Moisset M et al. PASTEC: an automatic transposable element classification tool. *PLoS One* 2014;**9**(5):e91929.
31. Smit AFA, Hubley R, Green P. RepeatMasker Open-4.0 (2013-2015). 2015. <http://www.repeatmasker.org>. Accessed 27 September, 2017.
32. Bolger AM, Lohse M, Usadel B. Trimmomatic: a flexible trimmer for Illumina sequence data. *Bioinformatics* 2014;**30**(15):2114–20.
33. Xie Y, Wu G, Tang J et al. SOAPdenovo-Trans: de novo transcriptome assembly with short RNA-Seq reads. *Bioinformatics* 2014;**30**(12):1660–6.
34. Grabherr MG, Haas BJ, Yassour M et al. Full-length transcriptome assembly from RNA-Seq data without a reference genome. *Nat Biotechnol* 2011;**29**(7):644–U130.
35. Gilbert D. EvidentialGene: mRNA Transcript Assembly Software. EvidentialGene: Evidence Directed Gene Construction for Eukaryotes. 2013. <http://arthropods.eugenesis.org/EvidentialGene/>. Accessed 27 September, 2017.
36. Schmutz J, McClean P, Mamidi S et al. A reference genome for common bean and genome-wide analysis of dual domestications. *Nat Genet* 2014;**46**(7):707–13.
37. Shu S, Goodstein DM, Hayes D et al. JGI plant genomics gene annotation pipeline. *SciTech Connect* 2017. <https://www.osti.gov/scitech/biblio/1241222-jgi-plant-genomics-gene-annotation-pipeline>. Accessed 27 September, 2017.
38. Gish W, States D. Identification of protein coding regions by database similarity search. *Nat Genet* 1993;**3**(3):266–72.
39. Slater GS, Birney E. Automated generation of heuristics for biological sequence comparison. *BMC Bioinformatics* 2005;**6**:31.
40. RepeatMasker Open-4.0. <http://www.repeatmasker.org>. Accessed 27 September, 2017.
41. UniProt Consortium. UniProt: a hub for protein information. *Nucleic Acids Res* 2014;**43**(D1):D204–12.
42. Salamov AA, Solovyev VV. Ab initio gene finding in *Drosophila* genomic DNA. *Genome Res* 2000;**10**(4):516–22.
43. Solovyev V, Kosarev P, Seledsov I et al. Automatic annotation of eukaryotic genes, pseudogenes and promoters. *Genome Biol* 2006;**7**(Suppl 1):S10.1–12.
44. Yeh RF, Lim LP, Burge CB. Computational inference of homologous gene structures in the human genome. *Genome Res* 2001;**11**(5):803–16.
45. Haas BJ, Delcher AL, Mount SM et al. Improving the Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids Res* 2003;**31**(19):5654–66.
46. Jones P, Binns D, Chang HY et al. InterProScan 5: genome-scale protein function classification. *Bioinformatics* 2014;**30**(9):1236–40.
47. Braga AC, Collevatti RG. Temporal variation in pollen dispersal and breeding structure in a bee-pollinated Neotropical tree. *Heredity* 2011;**106**(6):911–9.
48. Sork VL, Fitz-Gibbon ST, Puiu D et al. First draft assembly and annotation of the genome of a California endemic oak *Quercus lobata* nee (Fagaceae). *G3 (Bethesda)* 2016;**6**(11):3485–95.
49. Luo RB, Liu BH, Xie YL et al. SOAPdenovo2: an empirically improved memory-efficient short-read de novo assembler. *Gigascience* 2012;**1**(1):18.
50. Zimin AV, Marcais G, Puiu D et al. The MaSuRCA genome assembler. *Bioinformatics* 2013;**29**(21):2669–77.
51. Kajitani R, Toshimoto K, Noguchi H et al. Efficient de novo assembly of highly heterozygous genomes from whole-genome shotgun short reads. *Genome Res* 2014;**24**(8):1384–95.
52. Malinsky M, Simpson JT, Durbin R. Trio-sga: facilitating de novo assembly of highly heterozygous genomes with parent-child trios. *bioRxiv* 2016, doi: <https://doi.org/10.1101/051516>.
53. Boetzer M, Henkel CV, Jansen HJ et al. Scaffolding pre-assembled contigs using SSPACE. *Bioinformatics* 2011;**27**(4):578–9.
54. Nadalin F, Vezzi F, Policriti A. GapFiller: a de novo assembly approach to fill the gap within paired reads. *BMC Bioinformatics* 2012;**13**(Suppl 14):S8.
55. Hunt M, Kikuchi T, Sanders M et al. REAPR: a universal tool for genome assembly evaluation. *Genome Biol* 2013;**14**(5):R47.
56. Ponstingl H, Ning ZM. SMALT. 2010 - 2015 Genome Research Ltd. 2016. <http://www.sanger.ac.uk/science/tools/smalt-0>. Accessed 27 September, 2017.
57. Kent WJ. BLAT—the BLAST-like alignment tool. *Genome Res* 2002;**12**(4):656–64.
58. Stanke M, Keller O, Gunduz I et al. AUGUSTUS: ab initio prediction of alternative transcripts. *Nucleic Acids Res* 2006;**34**:W435–9.
59. Ramírez-Sánchez O, Pérez-Rodríguez P, Delaye L et al. Plant proteins are smaller because they are encoded by fewer exons than animal proteins. *Genomics Proteomics Bioinformatics* 2016;**14**(6):357–70.
60. Bradnam KR, Fass JN, Alexandrov A et al. Assemblathon 2: evaluating de novo methods of genome assembly in three vertebrate species. *Gigascience* 2013;**2**:10–10.

61. Lam ET, Hastie A, Lin C et al. Genome mapping on nanochannel arrays for structural variation analysis and sequence assembly. *Nat Biotechnol* 2012;**30**(8):771–6.
62. Ay F, Noble WS. Analysis methods for studying the 3D architecture of the genome. *Genome Biol* 2015;**16**:183.
63. Tørresen OK, Star B, Jentoft S et al. An improved genome assembly uncovers prolific tandem repeats in Atlantic cod. *BMC Genomics* 2017;**18**(1):95.
64. Wicker T, Sabot F, Hua-Van A et al. A unified classification system for eukaryotic transposable elements. *Nat Rev Genet* 2007;**8**(12):973–82.
65. Morgante M, Hanafey M, Powell W. Microsatellites are preferentially associated with nonrepetitive DNA in plant genomes. *Nat Genet* 2002;**30**(2):194–200.
66. Wang Q, Fang L, Chen J et al. Genome-wide mining, characterization, and development of microsatellite markers in gossypium species. *Sci Rep* 2015;**5**:10638.
67. Sonah H, Deshmukh R, Sharma A et al. Genome-wide distribution and organization of microsatellites in plants: an insight into marker development in brachypodium. *PLoS One* 2011;**6**(6):e21298.
68. Bernardi G. Isochores and the evolutionary genomics of vertebrates. *Gene* 2000;**241**(1):3–17.
69. Mizuno M, Kanehisa M. Distribution profiles of GC content around the translation initiation site in different species. *FEBS Lett* 1994;**352**(1):7–10.
70. Amit M, Donyo M, Hollander D et al. Differential GC content between exons and introns establishes distinct strategies of splice-site recognition. *Cell Rep* 2012;**1**(5):543–56.
71. Wendel JF, Greilhuber J, Dolezel J et al. *Plant Genome Diversity Volume 1 - Plant Genomes, Their Residents, and Their Evolutionary Dynamics*, vol. 1. Wien: Springer-Verlag; 2012.
72. JiaYan W, JingFa X, LingPing W et al. Systematic analysis of intron size and abundance parameters in diverse lineages. *Sci China Life Sci* 2013;**56**(10):968–74.
73. Yu J, Yang Z, Kibukawa M et al. Minimal introns are not “junk.” *Genome Res* 2002;**12**(8):1185–9.
74. Zhu J, He F, Wang D et al. A novel role for minimal introns: routing mRNAs to the cytosol. *PLoS One* 2010;**5**(4):e10144.
75. Simao FA, Waterhouse RM, Ioannidis P et al. BUSCO: assessing genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics* 2015;**31**(19):3210–2.
76. Ye J, Fang L, Zheng H et al. WEGO. a web tool for plotting GO annotations. *Nucleic Acids Res* 2006;**34**(Web Server issue):W293–7.
77. Cooper GM. *The Cell*, 2nd edition. A Molecular Approach. Sunderland, MA: Sinauer Associates; 2000.
78. Sulpice R, Trenkamp S, Steinfath M et al. Network analysis of enzyme activities and metabolite levels and their relationship to biomass in a large panel of *Arabidopsis* Accessions. *Plant Cell* 2010;**22**(8):2872–93.
79. Hamilton JP, Robin Buell C. Advances in plant genome sequencing. *Plant J* 2012;**70**(1):177–90.
80. Barthelson R, McFarlin AJ, Rounsley SD et al. Plantagora: modeling whole genome sequencing and assembly of plant genomes. *PLoS One* 2011;**6**(12):e28436.
81. Hellsten U, Wright KM, Jenkins J et al. Fine-scale variation in meiotic recombination in *Mimulus* inferred from population shotgun sequencing. *Proc Natl Acad Sci U S A* 2013;**110**(48):19478–82.
82. Cruz F, Julca I, Gómez-Garrido J et al. Genome sequence of the olive tree, *Olea europaea*. *Gigascience* 2016;**5**(1):29.
83. Plomion C, Aury JM, Amselem J et al. Decoding the oak genome: public release of sequence data, assembly, annotation and publication strategies. *Mol Ecol Resour* 2016;**16**(1):254–65.
84. Park B-S, Lee H-K, Lee S-E et al. Antibacterial activity of *Tabebuia impetiginosa martius* ex DC (Taheebo) against *Helicobacter pylori*. *J Ethnopharmacol* 2006;**105**(1–2):255–62.
85. Gómez Castellanos R, Prieto J, Heinrich M. Red Lapacho (*Tabebuia impetiginosa*)—a global ethnopharmacological commodity? *J Ethnopharmacol* 2009;**121**(1):1–13.
86. Byeon S, Chung J, Lee Y et al. In vitro and in vivo anti-inflammatory effects of taheebo, a water extract from the inner bark of *Tabebuia avellanaeae*. *J Ethnopharmacol* 2008;**119**(1):145–52.
87. Koyama J, Morita I, Tagahara K et al. Cyclopentene dialdehydes from *Tabebuia impetiginosa*. *Phytochemistry* 2000;**53**(8):869–72.
88. Hussain H, Krohn K, Ahmad VU et al. Lapachol: an overview. *Arkivoc* 2007;**2007**(2):145.
89. Widhalm J, Rhodes D. Biosynthesis and molecular actions of specialized 1,4-naphthoquinone natural products produced by horticultural plants. *Horticult Res* 2016;**3**:16046.
90. Romagnoli M, Segoloni E, Luna M et al. Wood colour in Lapacho (*Tabebuia serratifolia*): chemical composition and industrial implications. *Wood Sci Technol* 2013;**47**(4):701–16.
91. Alkan C, Sajjadian S, Eichler EE. Limitations of next-generation genome sequence assembly. *Nat Methods* 2011;**8**(1):61–65.
92. Silva-Junior OB, Grattapaglia D, Novaes E et al. Supporting data for “Genome assembly of the pink Ipê (*Handroanthus impetiginosus*, Bignoniaceae), a highly valued, ecologically keystone Neotropical timber forest tree.” *GigaScience Database* 2017. <http://dx.doi.org/10.5524/100379>.