1					
2	DR XIAOLAN ZHAO (Orcid ID : 0000-0001-7130-9388)				
3	MISS XIAODI HU (Orcid ID : 0000-0001-9704-2474)				
4					
5					
6	Article type : Original Article				
7					
8					
9	Chromosome-level assembly of Neolamarckia cadamba genome provides insights				
10	into the evolution of cadambine biosynthesis				
11	Xiaolan Zhao ^{1,6} , Xiaodi Hu ^{4,6} , Kunxi OuYang ^{1,6} , Jing Yang ^{1,2,6} , Qingmin Que ¹ ,				
12	Jianmei Long ¹ , Jianxia Zhang ¹ , Tong Zhang ¹ , Xue Wang ¹ , Jiayu Gao ¹ , Xinquan Hu ¹ ,				
13	Shuqi Yang ¹ , Lisu Zhang ¹ , Shufen Li ³ , Wujun Gao ³ , Benping Li ⁴ , Wenkai Jiang ⁴ ,				
14	Erik Nielsen ^{1,5} , Xiaoyang Chen ¹ *, Changcao Peng ¹ *				
15	¹ State Key Laboratory for Conservation and Utilization of Subtropical				
16	Agro-bioresources, South China Agricultural University, Guangzhou, 510642, China.				
17	Guangdong Key Laboratory for Innovative Development and Utilization of Forest				
18	Plant Germplasm, College of Forestry and Landscape Architecture, South China				
19	Agricultural University, Guangzhou, 510642, China.				
20	² School of Chinese Medicinal Resource, Guangdong Pharmaceutical University,				
21	Guangzhou 510006, China				
22	³ College of Life Sciences, Henan Normal University, Xinxiang 453007 China				
23	⁴ Novogene Bioinformatics Institute, Building 301, Zone A10 Jiuxianqiao North 13				
24	Road, Chaoyang District, Beijing 100083, China				
25	5 Department of Molecular, Cellular, and Developmental Biology, University of				
26	Michigan, Ann Arbor, Michigan 48109, USA				
27	⁶ These authors contributed equally to this article.				
	This is the author manuscript accepted for publication and has undergone full peer review but has not been through the copyediting, typesetting, pagination and proofreading process, which may lead to differences between this version and the <u>Version of Record</u> . Please cite this article as <u>doi: 10.1111/TPJ.15600</u>				

1	Correspondence:	Changcao	Peng	(ccpeng@scau.edu.cr	n), Xiaoyang	Chen
2	(xychen@scau.edu	ı.cn)				
3						
4						
5						
6						
7						
8						
9						
10						
11						
12						
13						
14						
15	Chromosome-lev	el assembly o	of <i>Neola</i>	<i>marckia cadamba</i> gen	ome provides in	sights
16		into the evo	lution o	f cadambine biosynth	esis	
17						
18	Significance state	ment: Neola	marckia	cadamba, Coffea cane	phora and Ophic	orrhiza
19	pumila are three c	losely-related	l species	in the Rubiaceae fami	ily which produc	e quite
20	different major s	pecialized m	etabolite	es including cadambi	ne, caffeine (a	purine
21	alkaloid) and cam	ptothecin (an	ti-cancer	MIAs). This chromos	some-level assen	nbly of
22	N. cadamba genon	ne and cadam	ıbine syn	thetic pathway analysi	is can not only fa	cilitate
23	the development	of tools for	or enhai	ncing bioactive prod	uctivity by me	tabolic
24	engineering but a	lso accelerat	te the u	nderstanding of the e	evolutionary hist	ory of
25	specific metabolic	pathways.				

26 Summary

Neolamarckia cadamba (Roxb.), a close relative of *Coffea canephora* and *Ophiorrhiza pumila*, is an important traditional medicine in Southeast Asia. Three major glycosidic
 monoterpenoid indole alkaloids (MIAs), cadambine and its derivatives
 This article is protected by copyright. All rights reserved

3β-isodihydrocadambine and 3β-dihydrocadambine accumulated in the bark and 1 exhibit antimalarial, antiproliferative, antioxidant, anticancer 2 leaves, and anti-inflammatory activity. Here, we report a chromosome-scale N. cadamba genome, 3 with 744.5 Mb assembled into 22 pseudochromosomes with contig N50 and scaffold 4 N50 of 824.14 Kb and 29.20 Mb, respectively. Comparative genomic analysis of N. 5 cadamba with C. canephora revealed that N. cadamba underwent a relatively recent 6 whole-genome duplication (WGD) event after diverging from C. canephora which 7 contributed to the evolution of MIA biosynthetic pathway. We determined the key 8 intermediates of cadambine biosynthetic pathway and further showed that NcSTR1 9 catalyzed the synthesis of strictosidine in N. cadamba. A new component, 10 epoxystrictosidine (C27H34N2O10, m/z 547.2285), was identified in the cadambine 11 biosynthetic pathway. Combining genome-wide association study (GWAS), 12 population analysis, multi-omics analysis, and metabolic gene cluster prediction, this 13 study will shed light on the evolution of MIA biosynthetic pathway genes. This N. 14 cadamba reference sequence will accelerate the understanding of the evolutionary 15 16 history of specific metabolic pathways and facilitate the development of tools for enhancing bioactive productivity by metabolic engineering in microbes or by 17 molecular breeding in plants. 18

Keywords: genome; Cadambine biosynthesis, strictosidine synthase, *Neolamarckia cadamba*, medicinal plant, evolution

21 Introduction

22 The evergreen tropical tree Neolamarckia cadamba (Roxb.) Bosser (Rubiaceae), commonly known as Kadamba or Kodom, belongs to the Rubiaceae family which is 23 the fourth largest family of angiosperms consisting of more than 660 genera and 24 11,000 species (Razafimandimbison, 2002, Robbrecht and Manen, 2006). The 25 26 Rubiaceae family is also noted for the production of important plant alkaloids, which 27 includes well-known plant species such as Coffea canephora and Ophiorrhiza pumila (De Luca et al., 2014, Kai et al., 2015, Sadre et al., 2016, Tran et al., 2018, Rai et al., 28 2021). N. cadamba was called a 'miracle tree' by World Forestry Congress (WFC) in 29 This article is protected by copyright. All rights reserved

1972 because its considerable economic potential as a fast-growing timber wood and 1 traditional medicinal resource in tropical and subtropical regions (Dwevedi et al., 2 2014, Pandey and Negi, 2016). The stems, bark and leaves of N. cadamba have been 3 widely used to treat a number of diseases such as diabetes, anaemia, stomatitis, 4 leprosy, cancer and a variety of infectious diseases in Southeast Asia (Pandey and 5 Negi, 2016). Although, there is lack of in-depth studies to clarify the active 6 metabolites responsible for various pharmacological activities attributed to N. 7 *cadamba*, recent studies pinpoint that its three major glycosidic monoterpenoid indole 8 alkaloids (MIAs), cadambine and its derivatives 3β-isodihydrocadambine and 9 3β-dihydrocadambine, exhibit antimalarial, antiproliferative, antioxidant, anticancer 10 and anti-inflammatory activity (Chandel et al., 2014, Dwevedi et al., 2014, Chandel et 11 al., 2017, Yuan et al., 2020). 12

Since its original discovery in N. cadamba (syn. Anthocephalus chinensis), 13 cadambine was also found to accumulate in Uncaria spp., Haldina cordifolia and 14 Emmenopterys henryi species in the Rubiaceae family (Handa et al., 1983, Wu et al., 15 2013, Wang et al., 2019, Chen et al., 2020). The structure of cadambine was firstly 16 determined by Handa and co-workers in 1983 and its biosynthesis was deduced to be 17 18 derived from strictosidine, in which C-18 is cyclized to N-4 with an ether bridge linking C-3 and C-19 (Handa et al., 1983). Most MIAs originates from the common 19 precursor 3-a(S)-strictosidine formed by stereospecific condensation of the indole 20 metabolite tryptamine and the end product secologanin in the iridoid (also called 21 secoiridoid) branch (De Luca et al., 2014). However, it remains unknown whether this 22 seco-iridoid pathway exists in N. cadamba. Therefore, determination of the key 23 intermediate strictosidine and characterization of functional strictosidine synthase 24 (STR) is the prerequisite for elucidation of cadambine biosynthetic pathway in N. 25 cadamba. 26

Here we report a chromosome-level genome assembly of *N. cadamba* (2n = 44chromosomes) obtained by a combination of Illumina and PacBio data platform. We used a Hi-C map (Burton *et al.*, 2013) to cluster the majority of the assembled contigs onto 22 pseudochromosomes. The *N. cadamba* genome was compared with *C. canephora* and twelve other available plant genomes to investigate whole-genome

duplication (WGD) events and expansion/contraction of gene families. We further 1 determined the key intermediate strictosidine, the known MIAs 3a-dihydrocadambine 2 and cadambine and a new component epoxystrictosidine by mass spectrometry and 3 NMR spectra, and characterized the first 'Pictet-Spenglerase' NcSTR1 in N. 4 cadamba. A total of 112 N. cadamba accessions collected from Southeast Asia were 5 sequenced to discover more loci and candidate genes for cadambine biosynthesis 6 based on genome-wide association study (GWAS) and population-level analysis. This 7 study revealed the evolution of cadambine biosynthesis in N. cadamba and will likely 8 provide additional insight into plant specialized metabolites. 9

10 **RESULTS**

11 Genome Sequencing, Assembly, and Annotation

The genome size of N. cadamba (2n = 2x = 44 chromosomes) was estimated to 12 be ~ 754 Mb with 0.69% heterozygosity and 54.29% repetition based on the k-mer 13 14 distribution analysis (Figure S1 & Table S1). We obtained a total of 92.97 Gb subreads (123.25X) generated from PacBio Sequel platform, plus another 52.30 Gb 15 reads (69.33X) from Illumina platform, 184.45 Gb 10X Genomics data (244.52X) and 16 17 89.59Gb Hi-C data (118.81X) (Table S2). The Falcon software package (Chin et al., 2016) was used for initial assembly of the PacBio reads, then polished error-corrected 18 with both PacBio and Illumina reads. 10X Genomics data was used to anchor contigs 19 into scaffolds. Finally, with the LACHESIS, the assembled scaffolds were anchored 20 to 22 pseudochromosomes based on Hi-C data (Figure 1a). A high-quality 21 chromosome-level genome assembly of N. cadamba was obtained with a total length 22 of 744.5 Mb, a contig N50 of 824.14 Kb and a scaffold N50 of 29.20 Mb (Table 1 23 &Figure 1a). The total length of the assembly was 744.5 Mb, which represents 24 25 98.7% of the estimated genome size.

To evaluate the assembly quality, we first mapped the Illumina reads back to the scaffolds, with a mapping rate of 98.11% and a coverage rate of 94.43%, respectively (**Table S3**). Second, we also evaluated the assembly using 1614 Benchmarking

Universal Single Copy Orthologs (BUSCO) genes from embryophyta (Simao et al., 1 2015) and 248 highly conserved core eukaryotic genes (CEGs) (Parra et al., 2007), 2 which showed that 1563 genes (96.8%) were annotated and 243 genes (98.0%) were 3 identified in our assembly, respectively (Table S4). Third, the reads of RNA-seq data 4 from 24 samples were mapped to the genome assembly using Hisat2 (Kim et al., 5 2015a, Kim et al., 2019). The alignment rate of 23 of these RNA samples was over 6 95%, with the remaining sample aligning at 89.18% (Table S5). Taken together, these 7 results indicated that the assembly of N. cadamba has high accuracy and 8 completeness. 9

Using a combination of ab initio and evidence-based methods, we predicted a 10 total of 35,461 protein-coding genes with an average gene length of 3,489 base pairs 11 and an average of 4.7 exons per gene (Table 1). Approximately 96.4% of the genes 12 with shared homology to NR, Swiss-Prot, KEGG and InterPro databases known genes 13 were functionally annotated (Table S6). In addition, we performed homology 14 searches and annotated non-coding RNA (ncRNA) genes (Table S7), yielding 666 15 16 transfer RNA (tRNA) genes, 1,642 ribosomal RNA (rRNA, 18S, 28S, 5.8S and 5S) genes, 2,701 small nuclear RNA (snRNA) genes and 1,053 microRNA (miRNA) 17 genes. 18

Repeat annotation revealed that 52.9% (394.1 Mb) of the assembled *N. cadamba*genome comprises of transposable elements (TEs) (Figure1b&Table S8).
Retrotransposons were found to be the dominant class of repeat elements (48.2%),
while DNA transposons account for 2.89% of the genome.

23

24 Evolution of the *N. cadamba* genome and comparative genomic analysis

To classify gene families in the *N. cadamba* genome, OrthoMCL (Li *et al.*, 2003)
was used to infer proteins from all 14 plant species (**Table S9**), generating a total of
32,185 orthologous gene families and 402 single-copy orthologous shared across 14
species. Phylogenetic analysis revealed that *N. cadamba* diverged from *Solanum tuberosum* (Solanaceae), *O. pumila* (Rubiaceae) and *C. canephora* (Rubiaceae) at
This article is protected by copyright. All rights reserved

around 80.4 MYA, 41.7 MYA and 31.0 MYA, respectively (Figure 1c). Compared 1 with C. canephora, 40 gene families underwent expansion in N. cadamba (Figure 2 1c). Functional annotation of these expanded genes demonstrated that 161 gene 3 ontology (GO) terms and 10 KEGG pathways were significantly enriched (FDR 4 cut-off < 0.05) which were involved in defense and sugar metabolism. Among 5 defense response functions, there is a clear expansion of G protein binding site 6 resistance and defense response genes in the N. cadamba genome. Another interesting 7 expanded gene family is GH1 β - glucosidase (BGLU) family, which are not only 8 involved in starch and sucrose metabolism, but also play a role in the biosynthesis of 9 phenylpropanoid and secondary metabolite (Ketudat and Esen, 2010, Xia et al., 10 2012). These expanded gene families may reflect the rapid growth, synchronized 11 metabolite synthesis, and specific adaptations to tropical environments for N. 12 cadamba (Figure 1c). 13

Synteny analysis revealed that two peaks (4DTv distance = ~ 0.17 and 4DTv 14 distance = ~ 0.51) were observed in the *N*. *cadamba* genome (Figure 1d). All gene 15 pairs showed a shallow peak at 0.51, likely reflecting a gamma triplication event 16 (whole-genome triplication, WGT- γ) occurred ~70 MYA in core eudicots (Paterson et 17 al., 2004). The 4DTv distribution also recovered the WGT-y in C. canephora, 18 consistent with the previous findings (Denoeud et al., 2014, Hu et al., 2019). Another 19 peak at 0.17 indicated that N. cadamba underwent a relatively recent WGD event 20 after diverging from C. canephora (Figure 1d), which was further supported by the 21 distribution of the synonymous substitution rate (Ks) (Figure S2). Intergenomic 22 co-linearity analysis demonstrated a 2:1 syntenic relationship between N. cadamba 23 24 and C. canephora, and 186 syntenic blocks were identified in N. cadamba by comparing to the C. canephora genome (Figure 1e). KEGG pathway analysis 25 revealed that the duplicated genes from the recent WGD were enriched with terms 26 such as "glucose metabolism" and "terpene synthase" (Table S10). 27

In this study, we also conducted a positive selection analysis using the genomic
sequences of *N. cadamba* and three close relatives. A total of 443 genes were possibly
under positive selection (PSGs, p-value<0.01, FDR < 0.05) using the branch-site
This article is protected by copyright. All rights reserved

model of the PAML software (Yang, 2007). KEGG functional classification of the
443 PSGs (Table S11) showed that the associated categories included "N-Glycan
biosynthesis", "Glycolysis/Gluconeogenesis", "Starch and sucrose metabolism " and
"Plant-pathogen interaction ".

5 Characterization of cadambine and the key intermediates

We firstly confirmed the structure of cadambine by Q-TOF LC-MS/MS (Figure 2a; 6 Figure S3b) (Chandel et al., 2012, Chandel et al., 2017) and NMR spectra (Figure 7 S3e, S3f) (Handa et al., 1983). The standard tryptamine (Figure 2a; Figure S3a), 8 9 3a-dihydrocadambine (Figure 2b; Figure S3c) (Takayama et al., 2003) and epoxystrictosidine (Fig 2b; Figure S3d) were also analyzed by high resolution mass 10 spectrometry. Secondly, the extracts from the leaves and barks of N. cadamba were 11 analyzed, this analysis revealed that N. cadamba accumulated tryptamine (Figure 2c, 12 2f) in the leaves, and strictosidine (Figure 2c, 2g), epoxystrictosidine (Figure 2d, 2i), 13 3a-dihydrocadambine (Figure 2d, 2h) and cadambine (Figure 2e, 2j) in the bark. 14 There were two compounds resolved with the same of m/z values (547.2285) (Figure 15 2h, 2i), one of which was matched with the peak at 12.30 min as standard 16 3a-dihydrocadambine ($C_{27}H_{34}N_2O_{10}$), another with the peak at 13.57 min as standard 17 epoxystrictosidine $(C_{27}H_{34}N_2O_{10})$ (Figure 2d). Notably, the key intermediate 18 strictosidine was detected. MS/MS spectrum of the ion at m/z 531 (Figure 2g) 19 displayed the same fragmentation pattern such as m/z 369 (C₂₁H₂₄N₂O₈) (m/z 20 $531 \rightarrow 369$, loss of glucose with 162 Da), m/z 514 (C₂₇H₃₁NO₉) (m/z 531 \rightarrow 514, loss of 21 NH₃ with 17 Da) and m/z 352 ($C_{21}H_{21}NO_8$) (m/z 531 \rightarrow m/z 352, loss of glucose and 22 NH₃). This identical fragmentation pattern is consistent with the MIA strictosidine 23 previously reported in Strychnos peckii (Santos et al., 2020). 24

Screening of Neolamarckia strictosidine synthase gene candidates and functional identification of NcSTR1

Previous orthogene-based analysis showed that the copy number of STR family
expanded in two Rubiaceae species (35 copies and 28 copies in *N. cadamba* and
coffee, respectively) compared to *C. roseus*, *T. cacao* and *A. thaliana*. Of the 35
putative *NcSTRs* (Figure 3b), seven *NcSTRs* were grouped with a conserved clade
This article is protected by copyright. All rights reserved

containing characterized strictosidine synthase gene CRO T006099 (Kellner et al., 1 2015) (Figure S4), whereas only five NcSTRs (named as NcSTR1, 13,14, 22, 27) 2 share over 45% aa identities to CrSTR. Sequence alignment of these predicted 3 NcSTRs with CrSTR1 (Pressnitz et al., 2018), RsSTR1 (Ma et al., 2006) and 4 OpSTR1 (Eger et al., 2020) demonstrated that most previously identified active site 5 residues, such as Cys89/Asn91/ Cys101 in RsSTR1, Trp145/Tyr147 in OpSTR1, and 6 the essential active site glutamate residue (Glu309 in RsSTR1, Glu301 in OpSTR1, 7 Glu315 in CrSTR1) were highly conserved in NcSTR1, NcSTR13 and NcSTR14 8 (Figure 3a). Moreover, NcSTR1 (evm.model.Contig69.90) physically clustered with 9 NcTDC2 (evm.model.Contig69.91) in N. cadamba genome (Figure 6a), and the 10 predicted NcTDC2 protein shares 71 and 84% sequence identities with the functional 11 tryptophan decarboxylases of C. roseus (De Luca and Cutler, 1987) and O. pumila 12 (Yamazaki et al., 2003a), respectively, thus, we firstly examined the catalytic 13 activities of the expressed NcSTR1 enzyme. 14

The coding sequence of NcSTR1 was expressed in Escherichia coli and recombinant 15 16 His-tagged proteins were purified by affinity chromatography (Figure S5). The purified NcSTR1 was mixed with excess tryptamine and secologanin (Figure 3c), a 17 reaction product of strictosidinewas produced and detected by HPLC-MS/MS. Two 18 ion pairs, 531.10 ->352.20 (blue curve) and 531.10 ->514.25 (pink curve), were 19 selected as parameters for multiple reaction monitoring (MRM). With the catalyzation 20 of NcSTR1, the peak area of the resulting strictosidine increased after 48 hours 21 reaction comparing with 24 hours (Figure 3c). The control without NcSTR1 did not 22 have the relative peak. These results indicated that NcSTR1 could catalyze the 23 synthesis of strictosidine in N. cadamba. 24

Integrated transcriptome and metabolome analysis for Cadambine biosynthetic gene discovery

The detection of strictosidine in Neolamarckia bark indicated that *N. cadamba* possesses a same upstream pathway to the formation of strictosidine as that of *C. roseus* (Figure 4a). The gene lists involved in MVA/ MEP and shikimate/indole

pathway in *N. cadamba* were obtained by the preliminary screening (Table S12), and 1 those with more than 75% aa identity to Catharanthus functional genes were selected 2 as candidate biosynthetic genes (Table S 13). The predicted genes involved in the 3 seco-iridoid biosynthesis pathway (GES, G8H, GOR, IS, IO, 7-DLGT, 7-DLH, LAMT 4 and SLS) in N. cadamba (Table S13) only contained the orthologues with more than 5 65% as identity to functional genes from C. roseus and O. pumila (Salim et al., 2013, 6 Brown et al., 2015, Kai et al., 2015, Kellner et al., 2015). We found that most genes 7 related to strictosidine synthesis (i.e., GPPS, GES, DLGT, LAMT, SLS and TDC) 8 9 underwent the recent WGD event after diverging from C. canephora (Ks<1 in Table S 14), indicating that the recent WGD event were important to the evolution of 10 cadambine biosynthesis. 11

As for the downstream pathway after strictosidine synthesis in N. cadamba, the 12 biogenesis of the dihydrocadambines may involve the epoxidation of the strictosidine 13 followed by internal opening of the epoxide, closing the seven-membered ring and 14 producing 3a-dihydrocadambine. After that, 3a-dihydrocadambine would be 15 catalyzed by hydrolases to form cadambine (Figure 4a and Figure S6). In 16 consistence with this prediction, we found that the epoxystrictosidine 17 (C27H34N2O10, m/z 547.2285) accumulated in Neolamarckia bark (Figure 2f). 18 Therefore, the biosynthetic enzymes with monooxygenase activity, hydrolase 19 activities, acyltransferase activity and glucosidase activities, such as cytochrome 20 P450s (CYP), squalene epoxidases (SQE), oxidosqualene cyclases (OSCs), 21 zeaxanthin epoxidase (ZEP), soluble epoxide hydrolase (SEH), violaxanthin 22 de-epoxidase (VDEs), BAHD and serine-carboxypeptidase-like acyltransferases 23 (ACT), strictosidine glucosidase (SG) etc., are proposed to serve as the possible 24 organizers in the cadambine synthetic pathway (Ma et al., 2005, Xia et al., 2012, 25 Leonelli et al., 2017, Almeida et al., 2018, Carqueijeiro et al., 2018a, Carqueijeiro et 26 al., 2018b, Qu et al., 2018, Zheng et al., 2019, Shang and Huang, 2020, Tatsis et al., 27 28 2017). In order to identify more biosynthetic gene candidates, we further conducted RNA-seq analysis, GWAS and population analysis and gene cluster prediction. 29 RNA-seq analysis with 16 different tissues of N. cadamba revealed that the 30

31 putative seco-iridoid biosynthetic genes had two distinctive co-expression patterns

(Figure 4b), one set included NcGES1- NcIS1-NcIO1 -NcDLGT1/3 which showed 1 high expression in bark, bud and young leaves, and another set, NcLAMT1/2 and 2 NcSLS4 which exhibited moderate expression in most of the collected tissues except 3 for fruit, old leaves, roots and xylem. Moreover, NcLAMT1/2 and NcSLS4 4 co-expressed with NcTDC2 (Figure 4b and Figure S7), indicating that the synthesis 5 of secoiridoid was coordinated with indole precursors needed for cadambine 6 production. In addition, NcSTR1/ 14/ 19/ 29/ 35 and NcSQE1/6 exhibited 7 co-expression pattern with that of NcLAMT1/2 and NcTDC2 (Figure 3b, 4b and 8 Figure S8). 9

To assess the regulatory processes that control the accumulation of strictosidine, 10 tryptamine and cadambine in different organs, five tissues (bark, bud, young leaves, 11 old leaves and fruits) from the individual tree sampled for genome sequencing were 12 used as the source of transcripts and metabolites. The transcripts were quantified by 13 RNA-seq expression analysis (FPKM) and the metabolites contents (µg) per gram 14 fresh weight of five tissues were determined using a liquid chromatography-mass 15 spectrometry-based method. We calculated the Pearson correlations among the 16 contents of strictosidine, tryptamine, cadambine and gene transcript levels, 17 respectively. The correlation analysis (coefficient >0.9 as the cutoff) showed that the 18 expression level of NcSTR1 was strongly associated with strictosidine content, and 19 that NcSTR14 (evm.model.Contig341.115), NcSLS3 (evm.model.Contig60.399) with 20 cadambine content (Table S15). Two NcCYPs (evm.model.Contig208.74, 21 evm.model.Contig 708.57) and NcOSC1 (evm.model.Contig207.232) were also 22 identified to be strongly associated with cadambine content in the network. Both 23 NcCYPs were assigned into CYP72A subfamily which contained numerous essential 24 components in MIA-producing plants (Figure S9) (Urlacher and Girhard, 2019, 25 Zheng et al., 2019). In addition, the expression level of twelve transcription factors 26 (TFs) were highly correlated with cadambine accumulation pattern. Of these, 27 evm.model.Contig 28.407 putatively encoded a basic helix-loop-helix (bHLH) TF 28 with 50% and 45% aa identity to CrBIS1 and CrBIS2, respectively (Van Moerkercke 29 et al., 2015, 2016). 30

GWAS and population identification of genes potentially related to cadambine biosynthesis

To gain further insights into cadambine synthetic pathway, we conducted a 3 genome-wide association study (GWAS) using 112 individuals collected from 27 4 populations distributed widely in Southeast Asia (Table S16). We identified 5 5,786,667 high-quality SNPs based on the reference genome. To establish SNP-based 6 7 identification of genotype-phenotype associations, we examined the contents of strictosidine, tryptamine and cadambine in the extracts from the bark and young 8 leaves of 112 Neolamarckia accessions by high resolution mass spectrometry. The 9 association between each SNP and the phenotype was conducted in a mixed linear 10 model using GEMMA (Zhou and Stephens, 2014). GWAS revealed that the candidate 11 genomic loci responsible for Neolamarckia strictosidine and cadambine accumulation 12 were associated with NcSTR3/ 4/13, NcTDC2 and most of the putative genes in the 13 seco-iridoid biosynthesis pathway (Table S17). GWAS also revealed new 14 15 biosynthetic gene candidates including NcSQE7, NcSEH7, NcCPR1 and six NcCYPs putatively involved in the epoxidation and the hydrolysis of epoxide (Meijer et al., 16 1993, Zheng et al., 2019) were associated with the significant loci responsible for 17 cadambine accumulation (Table S18) . 18

In this study, we also conducted population-level analysis of genetic variation of *N. cadamba* using 31 accessions with high-, medium-, or low- cadambine level. To observe the divergence among the three groups at the genomic level, we constructed a neighbor-joining phylogenetic tree and performed principal component analysis (PCA) of members of clade 3 (**Figure 5a and 5b**). We observed similar results in phylogenetic tree and PCA analyses.

We also analyzed the linkage disequilibrium (LD) throughout the whole genome. LD (indicated by r2) decreased with physical distance between SNPs in all groups. The average distance of LD for each group was measured as the chromosomal distance when LD decreased to half of its maximum value. The three groups showed different extents of genome-wide LD decay, with LD decaying fastest in low cadambine

population and slowest in high (**Figure 5c**). Then, we explored the genomic regions with high divergence to try to identify genes possibly involved in the cadambine metabolism and/or accumulation. According to an empirical procedure described in a previous study (Li *et al.*, 2013), the intersection regions with the top low or high π ratios and the top high F_{ST} values between ecotype groups were identified as selective sweeps.

We calculated the fixation index (F_{ST}) between subgroups, and the genomic regions 7 with fixation index (F_{ST}) values (Weir and Cockerham, 1984) in the top 5% were 8 considered highly differentiated. We estimated the population diversification 9 10 parameters, π and θ w, and found that the overall nucleotide diversity in species with 11 high cadambine level was higher than that in species with low- or mediumcadambine level (high- cadambine level species, π =0.0023440 and θ w=0.0019451; 12 low- cadambine level species, π = 0.0023134 and θ w= 0.0018883) (Figure 5d). 13 KEGG analysis (FDR <0.05) indicated the genes harbored in these selective sweeps 14 15 were involved in cadambine biosynthesis (e.g., sesquiterpenoid and triterpenoid biosynthesis, terpenoid backbone biosynthesis, and indole alkaloid biosynthesis, 16 Figure 5e and Table S19). 17

18 Prediction of biosynthetic gene clustering and tandem duplication in *N. cadamba*

Biosynthetic gene clusters (BGCs) for a wide variety of natural products have now 19 been reported from diverse plant species (Li et al., 2021). The definition of a 20 21 metabolic gene cluster requires that it should contain genes encoding at least three 22 different types of tailoring enzymes (Nützmann et al., 2016, Schläpfer et al., 2017, Jacobowitz and Weng, 2020). We systematically mined the N. cadamba genome 23 sequence to identify all the BGCs with plantiSMASH algorithm (Kautsar et al., 24 2017). This identified 67 possible gene clusters across the 22 pseudochromosomes 25 (Table S20). Moreover, we identified additional BGCs that included at least one gene 26 involved in the seco-iridoid pathway, or NcTDC/NcSTR within physical sizes less 27 than 600 kb (Table S21). 28

A notable 12-gene cluster (NcSMASH35) within a 340 kb region in chromosome 11 1 was identified, which included functionally characterized NcSTR1, NcTDC2 and five 2 tandem duplicated CYP71Ds subfamily members, a CYP76 member and four TPSs 3 (Figure 6a). A predicted transporter gene evm.model.Contig69.96 putatively 4 encoding a major facilitator superfamily protein with moderate expression in all the 5 tissues was present in this NcTDC2/NcSTR1 gene cluster, similar to the previous 6 reports in R. stricta (Sabir et al., 2016), G. sempervirens and O. pumila (Rai et al., 7 2021). All the six NcCYPs in this cluster were derived from the recent WGD event 8 (Ks<1), and the closest match of the five NcCYP71D paralogs is C. roseus 9 tabersonine 3-oxygenase gene (CrT3O, AEX07771) with 51%-58% identity at the 10 amino acid level (Figure S9). 11

In N. cadamba, we identified two genes with 88% identity at the nucleotide level 12 (evm.model.Contig267.36, evm.model.Contig625.59) putatively encoding loganic 13 acid O-methyltransferase (LAMT). The sequence alignment of the predicted 14 NcLAMT1/2 with CrLAMT (CRO T028497) revealed that all the identified active 15 16 site residues of CrLAMT (Y159, H162, W163, P227, A241, H245, Q273, H275, P302, Q316, I320 and D359) were conserved in NcLAMT1/2 (Figure S10). Two 17 genomic regions with NcLAMT1/2 located on pseudochromosome 10 (25.31 18 Mb-25.13 Mb) and pseudochromosome 22 (3.78 Mb- 5.73 Mb) were identified, in 19 which included the gene clusters NcMCL10, NCMCL19, NcMCL23 and 20 NcSMASH67 (Figure 6b). In NcMCL10, the Neolamarckia specific 7-DLH 21 homolog evm.model.Contig188.56 (86% aa identity to CrDL7H) is located c. 600 kb 22 from NcLAMT1, and two other NcCYP72As were closely adjacent to the predicted 23 NcDL7H. Interestingly, the two cadambine biosynthetic candidates 24 (evm.model.Contig208.74, evm.model.Contig 708.57) identified in the 25 gene-to-metabolite network were also distributed in NcSMASH67 and NcMCL23, 26 respectively, implying a selection pressure favoring clustering of genes associated 27 with MIA production in N. cadamba. 28

The gene cluster NcMCL15 located on pseudochromosome 20 consists of NcIS1/2,
NcSQE1 and two NcCYP81s (Figure 6c). A syntenic region with high sequence This article is protected by copyright. All rights reserved

similarity and gene localization on pseudochromosome 16 was identified in N. 1 cadamba. Synteny analysis among O. pumila, N. cadamba, and C. canephora further 2 revealed that *N. cadamba* had underwent a WGD. This is further supported by the low 3 median Ks value (~ 0.41) of gene paralogs for most genes in NcMCL15 after 4 excluding the two key genes NcIS1 and NcIS2. NcSQE1 is located c. 300 kb from 5 NcIS1 in cluster NcMCL15, and the co-expression profiles of NcSQE1/6, 6 NcLAMT1/2, NcTDC2 and NcSTR1/14 suggests that NcSQE1/6 might be involved 7 in the epoxidation of the strictosidine. 8

Another predicted cluster (NcSMASH17) in N. cadamba is located on 9 pseudochromosome 4 and consists of NcSQE4, a MYB TFs gene with high similarity 10 to OpMYB1 (BAU61355.1, 83% aa identity) and a subtilisin-like protease (SBT) 11 gene associated with cadambine content in the gene-to-metabolite network (Table 12 S15). OpMYB1 is an R2R3- MYB repressor that acted as a negative regulator of MIA 13 production, and its overexpression in hairy roots of O. pumila resulted in reduced 14 production of camptothecin and reduced expression of OpTDC (Rosseleena et al., 15 16 2016, Ma and Constabel, 2019).

Previous studies demonstrated that specialized metabolic genes were more significantly enriched in local (tandem) duplication events as compared with whole-genome duplication events (Chae *et al.*, 2014). We noted that in the *N. cadamba* genome the presence of multiple paralogs of several predicted MIA biosynthetic genes putatively encoding 7-DLGT, IS, ISY, TDC, SLS, STR and CYPs were also significantly enriched in local (tandem) duplication events (**Figure 6a,6b** and **Figure S10**).

In addition, tandem duplication of TFs involved in specialized metabolite
biosynthesis were also detected in a wide range of eudicots (Colinas and Goossens,
2018). In *C. roseus*, three bHLH iridoid synthesis gene (*BIS*) tandem duplicates
(*CrBIS1/2/3*) exclusively transactivated the expression of the biosynthetic genes in the
seco-iridoid pathway (Van Moerkercke *et al.*, 2015, Van Moerkercke *et al.*, 2016,
Singh *et al.*, 2021). In the gene-to-metabolite network, we characterized one
homologue to *CrBIS1*, evm.model.Contig 28.407. Analysis of the genome sequence
This article is protected by copyright. All rights reserved

of N. cadamba showed that four additional bHLH TF genes clustered with 1 evm.model.Contig 28.407 in a tandem duplication within a 130-kb region on 2 pseudochromosome 4 (NcMCL24, Figure 6d). Synteny analysis further detected 3 another bHLH TFs tandem duplicates (NcMCL25) with close homology and gene 4 localization to those of NcMCL24 on pseudochromosome 8, in which consists of the 5 significant loci responsible for strictosidine accumulation in GWAS (24.25-24.28 6 Mb). The co- expression profile of some bHLH TFs in NcMCL24 and NcMCL25 7 suggesting co-regulation mechanisms in specialized metabolite biosynthesis (Figure 8 S12). In NcMCL26 (Figure S13), four NcMYC TFs were clustered in tandem order 9 within a 24-kb region and these loci were characterized to be responsible for 10 strictosidine accumulation in GWAS. Of the four MYC TFs, evm.TU.Contig184.722 11 showed the highest similarity (77% for the deduced amino acid sequence) to CrMYC2 12 (AF283507), the major activator of MeJA-responsive ORCA2/3, which in turn 13 regulated a subset of alkaloid biosynthesis genes in C. roseus (Zhang et al., 2011). 14

15 Discussion

Historically (Mehra and Bawa, 1969), and in the Chromosome Counts Database 16 (http://ccdb.tau.ac.il/home/), the Anthocephalus cadamba (Roxb.) Miq (now accepted 17 as N. cadamba (Roxb.) Bosser) gametophytic chromosome number n was 22. N. 18 cadamba was characterized as a tetraploid species in those relatively old references 19 (Mehra and Bawa, 1969, Bedi et al., 1981). As there were confusing names for N. 20 cadamba (Pandey and Negi, 2016), we first identified the sequenced N. cadamba is 21 diploid with 2n = 2x = 44 chromosomes by chromosome measurements and a pair of 22 45s rDNA probe signals using fluorescent in situ hybridization (FISH) (Figure S1B). 23 This provided a solid base for establishing a high-quality of N. cadamba Hi-C 24 libraries. 25

To date, three species in the genus *Coffea* (*C. canephora*, *C. arabica*, *C. eugenioides*), *O. pumila*, and *Gardenia jasminoides* in the family Rubiaceae have been sequenced. *N. cadamba* is the first sequenced tree species in this family. The genus *Neolamarckia* is a ditypic genus that included only two tree species: *N.* This article is protected by copyright. All rights reserved

cadamba and N. macrophyllus, and the latter is endemic to Sulawesi, Indonesia (Li et 1 al., 2018). This genus is featured by the densely globe-shaped flower clusters (the 2 origin of its Chinese name "TuanHua") with orange scent (Kareti and Subash, 2020). 3 Another representative characteristic of this genus is its rapid growth and remarkable 4 canopy with a height of 45 m and a stem diameter of 100-160 cm. Notably, N. 5 *cadamba* is the only species reported so far that accumulated high content ($\sim 0.1\%$) of 6 bioactive 3-dihydrocadambine and cadambine in the bark and leaves. Given that we 7 have established a highly efficient regeneration system of N. cadamba (Li et al., 8 2019), the genomic information presented here will greatly facilitate the elucidation of 9 cadambine synthetic pathway and the development of tools for enhancing bioactive 10 productivity by metabolic engineering in microbes or by molecular breeding in plants. 11

Furthermore, comparative genomic analysis among N. cadamba, C. canephora 12 and O. pumila will shed light on the evolutionary history of specific metabolic 13 pathways as they produce quite different major specialised metabolites including 14 cadambine, caffeine (a purine alkaloid) and camptothecin (anti-cancer MIAs). The 15 emergence of STR for strictosidine synthesis was generally considered an important 16 innovative step for the strictosidine-derived-MIA-producing plants (Rai et al., 2021). 17 However, we found that none candidate homologous (aa identity cutoff >55%) 18 corresponding to the late seco-iridoid pathway genes, NcLAMT1/NcLAMT2 and 19 NcSLS1-4, were identified in C. canephora genome. In contrast to coffea, the 20 Ophiorrhiza candidate homologous OpLAMT (Opuchr05 g0056110) and OpSLSs 21 (Opuchr02 g0012990, 22 Opuchr02 g0013060, Opuchr02 g0013090 and Opuchr02 g0017930) share 84% aa identity with NcLAMT1/NcLAMT2 and 67-85% 23 aa identities with NcSLS1-4, respectively. N. cadamba diverged ~ 41.7 Mya from O. 24 pumila and ~ 31.0 Mya from C. canephora (Figure 1a), we therefore propose that the 25 LAMT gene evolution involved a deletion following a duplication process, and that 26 the duplication of NcLAMT derived from a recent WGD event was key for evolution 27 28 of cadambine biosynthesies in N. cadamba as there is only one copy of LAMT gene in both well-known MIA-producing species, such as C. roseus and O. pumila (Figure 29 S10). Moreover, although camptothecin is found in *O pumila*, OpLAMT, OpSLS and 30 This article is protected by copyright. All rights reserved

OpSTR have the same functions as those of C. roseus and similarly produce loganin, 1 secologanin and strictosidine (Rai et al., 2021) also observed that the three key 2 N-methyltransferases (NMTs): xanthosine methyltransferase (CcXMT), theobromine 3 synthase (7-methylxanthine methyltransferase, CcMXMT), and caffeine synthase 4 (3,7-dimethylxanthine methyltransferase, CcDXMT) essential for caffeine 5 biosynthesis are missing in the N. cadamba genome, as the key residues Gln161, Ile 6 266, Ser 316 and Tyr 356 for CcXMT1 (Mccarthy A A and Mccarthy J G, 2007) were 7 no longer conserved in the deduced amino acid sequences of the Neolamarckia 8 9 homologous, providing a reasonable explanation for why N. cadamba lacks the first methylation steps necessary to produce caffeine from xanthosine (Figure S10). 10

Two O-β-D-glucosidases, strictosidine O-β-D-glucosidase (SG, EC 3.2.1.105) 11 and raucaffricine O-β-D-glucosidase (RG, EC 3.2.1.125), act as major components in 12 the MIA biosynthesis pathway. In C. roseus and Rauvolfia, SG follows strictosidine 13 synthase in the production of the reactive intermediate required for the diverse MIAs, 14 whereas Rauvolfia RG hydrolyzes the glucoalkaloid raucaffricine forming the 15 aglycone vomilenine, an intermediate that appears in the middle of the ajmaline 16 pathway (Barleben et al., 2007, Xia et al., 2012). From a biosynthetic point of view, it 17 is extremely rare for a glucoside such as the glucoalkaloid strictosidine to act as a 18 precursor at the beginning of biosynthetic pathway and for it to become activated by 19 deglucosylation (Barleben et al., 2007). We therefore explored whether there are 20 similar or identical enzymes to SGs or RsRG in N. cadamba. Homology searching 21 22 and protein sequence alignment displayed that all the Neolamarckia and Ophiorrhiza homologus of CrSG (Rai et al., 2021) were highly conserved in the essential active 23 site residues (Glu207, Glu416, His161 and Trp388) of RsSG (Barleben et al., 2007), 24 whereas not conservative at four critical active site residues (Thr189, His193, Tyr200, 25 Ser390) of RsRG (Figure S14), suggesting that these Neolamarckia genes are more 26 likely assigned to plant SGs. This analysis suggested that some members of 27 strictosidine O-β-D-glucosidases probably play a role in the production of the reactive 28 intermediate of the diverse MIAs in N. cadamba. 29

30

The knowledge about selective sweeps provides insights and targets for This article is protected by copyright. All rights reserved utilization of germplasm. Moreover, new genetic variation is needed to use to increase
the cadambine content. In our study, we reported the genome variation mapping of
112 accessions and we detected numerous diverse selective sweeps among ecotype
groups in association with environmental adaptability and cadambine-related traits by
analyzing genome structure diversifications between the three ecotype groups. The
information of selective sweeps associated with agriculturally important markers will
be helpful for molecular breeding in *N. cadamba*.

8

9 Methods

10 Chromosome preparation and FISH analysis

For improved characterization of the chromosomes, 45S rDNA was labeled with Chroma Tide Alexa Fluor 488–5-dUTP (Invitrogen) for FISH and theitotic metaphase spreads were prepared from meristem root tip cells of *N. cadamba* following the procedures described in the reference (Deng *et al.*, 2012) with minor modifications.

15 Plant Materials, DNA Extraction and Genome Sequencing

The young tender leaves of *N. cadamba* collected from a 7-years-old individual plant in South China Agricultural University (Guangzhou, China) were used for DNA extraction. Paired-end (PE) library (insert size 350 bp) was constructed and sequenced on Illumina Xten platform.

SMRT Bell libraries with an insert size of 20 kb were also constructed and sequenced
on PacBio Sequel platform. For 10 X genomic sequencing, DNA sample preparation,
indexing, and barcoding were done using the GemCode Instrument (10 X Genomics).

23 Genome assembly and assessment

De novo assembly of the PacBio reads was performed using FALCON
(https://github.com/PacificBiosciences/FALCON/) and FALCON-Unzip. With the
phasing information from the raw reads, it generates a subsequent set of primary
This article is protected by copyright. All rights reserved

contigs contig polished using Ouiver 1 and these sequences were (http://pbsmrtpipe.readthedocs.io/en/master/getting_started.html). The 184.45 Gb 10X 2 Genomics data were aligned to the initial assembly using BWA and scaffolding 3 approach was performed by FragScaff. After that, we used Pbjelly software to fill 4 gaps with PacBio data with the parameters: -minMatch 8 -sdpTupleSize 8 5 -minPctIdentity 75 -bestn 1 -nCandidates 10 -maxScore -500 -nproc 13 -6 noSplitSubreads. Finally, the anchorage of the genome assembly onto chromosomes 7 8 was performed by the LACHESIS pipeline.

9 To evaluate the completeness and quality of the assembly, the high-quality reads from 10 short insert size were mapped to the assembly using BWA(Li and Durbin, 2009). 11 Additionally, CEGMA (Core Eukaryotic Genes Mapping Approach) defined a set of 12 conserved protein families that occur in a wide range of eukaryotes and identified 13 their exon-intron structures in genomic sequences. We also used BUSCO 14 (Benchmarking Universal Single-Copy Orthologs) to assessment the completeness of 15 the genome assembly.

16 Repetitive elements and genes annotation

The approach combined the *de novo* and homology was used to search transposable elements (TEs) in the *N. cadamba* genome. In *de novo* approach, RepeatModeler (http://www.repeatmasker.org/RepeatModeler.html), RepeatScout, LTR-Finder and Tandem Repeats Finder (TRF) sofrwares were used. In homology prediction, RepeatMasker (version 3.3.0) and RepeatProteinMask software were used to against Repbase TE library (http://www.girinst.org/repbase) and TE protein database, respectively.

Three independent approaches, including homology alignment, *de novo* search and transcriptome prediction was applied to predict protein coding genes in the *N*. *cadamba* genome.

Homology-based gene prediction: homolog proteins sequences of *Arabidopsis thaliana*, *Populus trichocarpa*, *Solanum tuberosum* and *Coffea canephora* were

downloaded from Ensemble (http://plants.ensembl.org/index.html) and NCBI 1 (https://www.ncbi.nlm.nih.gov/) and then aligned to N. cadamba genome assembly 2 using tblastn (E-value 1e⁻⁵). 3 Ab initio gene prediction: Augustus (version 2.5.5), Genscan (version 1.0), 4 GlimmerHMM (version 3.0.1), Geneid and SNAP were used to predict coding regions 5 6 in the repeat-masked genome. Transcriptome-assisted gene prediction: Tophat (version 2.0.8)(Trapnell et al., 2009, 7 Kim et al., 2013) was used to map clean RNA-seq reads to N. cadamba genome and 8 Cufflinks (version 2.1.1)(Trapnell et al., 2012) (http://cufflinks.cbcb.umd.edu/) and 9 then used to assemble the transcripts into gene models (Cufflinks-set). 10

Gene model evidence from Homo-set, ab initio and Cufflinks-set programs were combined by EvidenceModeler (EVM) (http://evidencemodeler.github.io/) into a non-redundant set of gene structures.

14 Functional annotation protein coding genes

The predicted protein sequences were searched against six protein/function databases: InterPro, Pfam, SwissProt, NR, GO and KEGG. The InterPro and Pfam databases search were performed using InterproScan (V4.8) and HMMER (V3.1), respectively. For the other databases, BLAST searches were performed with an E value cutoff of 19 1e-05.

20 Species phylogenetic analysis

OrthoMCL(http://orthomcl.org/orthomcl/) was used to cluster paralogous and
orthologous among 14 species (*N. cadamba, Arabidopsis thaliana, C. canephora, O. pumila , Eucalyptus grandis, Elaeis guineensis, Malus domestic, Glycyrrhiza uralensis, Theobroma cacao, Oryza sativa, S. tuberosum, Populus trichocarpa, Actinidia chinensis* and *Vitis vinifera*). We obtained the similarity relations between
all species protein sequences through all-vs-all blastp with the e-value 1e-5.
MUSCLE(Edgar, 2004) (http://www.drive5.com/muscle/) was used to align all 402

single-copy gene protein sequences and make a super alignment matrix. Then, 1 RAxML(Stamatakis, 2014) 2 (http://sco.h-its.org/exelixis/web/software/raxml/index.html) was used to construct the 3 14 species phylogenetic tree with maximum likelihood method. Finally, the 4 MCMCtree program(Puttick, 2019) (http://abacus.gene.ucl.ac.uk/software/paml.html) 5 was applied to infer the divergence time based on the phylogenetic tree constructed. 6 CAFÉ 2.2(De Bie et al., 2006) was used to determin expansion and contractions of 7 orthologous gene families. 8

9 Whole genome duplication (WGD)

To identity syntenic blocks, the protein sequences from C. canephora, S. tuberosum 10 and N. cadamba were searched against themselves using blastp (E<1e-5)(Scott and 11 12 Madden, 2004). Syntenic blocks were determined by MCscanX (Wang et al., 2012) with the parameter at least five genes per block. Then we calculated the 4DTv 13 (fourfold degenerate synonymous sites of the third codons) for syntenic segments 14 from the concatenated alignments which constructed by fourfold degenerate sites of 15 all gene pairs found in each segment and plotted the distribution of the 4DTv values. 16 The synonymous substitution rate (Ks) was also calculated using MYN 17 algorithm(Wang et al., 2009) based on the Tamura-Nei Model to further confirm most 18 recent WGD. 19

20 NMR and LC -MS/ MS analyses

Fresh leaves and bark collected from the same tree for genome sequencing were ground and extracted with 70% aqueous ethanol (v/v) for 20 min. The extract was prepared for chromatograph analysis. The NMR spectra of cadambine from Song's Lab (Xishuangbanna Tropical Botanical Garden, Chinese Academy of Sciences) was characterized on a Bruker AV 400 spectrometer using tetramethylsilane (TMS, $\delta = 0$) as internal reference. The spectra was consistent with the data in the reference (Figure S3, Yuan *et al.*, 2020).

All the intermediates were detected by UPLC-Q-TOF (UPLC1290-6540B Q-TOF, 1 Agilent, USA). The mobile phase comprised with acetonitrile (A) and 0.2% formic 2 acid+10 mM ammonium formate (B). The synthesized strictosidine in vitro catalyzed 3 by NcSTR1 was analyzed by HPLC-MS/MS (LCMS-8050, Shimadzu, Japan). A 4 binary gradient elution with a flow rate of 0.3 mL min⁻¹ was performed as follows: 6 5 min 25% methanol, 5 min 90% methanol, 5.1 min 10 % methanol and 7 min 10% 6 methanol. The temperature of the column oven is 40 °C and the sample tray 4°C. The 7 standard substances included secologanin (Absin, China, 19351-63-4), tryptamine 8 (Sigma, USA, 61-54-1), 3 a -dihydrocadambine and epoxystrictosidine (Nature 9 Standard, China, 54483-84-0). 10

11 The *in vitro* catalytic assay of the recombinant strictosidine synthase NcSTR1

12 The coding sequence of evm.model.contig69.90 without start codon was amplified with a primer pairs (forward primer: 5'-CTCTTGGAAATCATGCCTCACA- 3' and 13 reverse primer: 5'- TCAGACAGAAGAAACCACTCCATTC-3') and cloned into the 14 pEASY®- Blunt E1 Expression Vector (TransGen Biotech, Beijing, China), and 15 further transformed in E. coli DH5a. The recombinant plasmids were then introduced 16 into E. coli BL21 Rosetta DE3 strains for NcSTR1 protein expression. The 17 recombinant NcSTR1 production and purification was obtained as the procedure 18 described in the reference. 19

The procedure of *in vitro* catalytic assay was modified from Pressnitz et al. (2018). 20 Generally, the catalytic reaction solution consisted of purified enzyme preparation (10 21 µg) dissolved in PIPES buffer (50 mM, PH 6.1), 5 mM secologanin and 5 mM 22 tryptamine with a total volume of 1000 µL. The mixtures were incubated on a shaker 23 24 at 35 °C for a given time (0 hr, 24 hr, 48 hr). The reaction was terminated by addition 25 of 10 M NaOH (100 µL). Then ethyl acetate (500 µL) was added to precipitate the PIPES. The collected supernatant was dried by nitrogen gas and then dissolved in 100 26 µL 50% methanol and subjected to HPLC-MS/MS analysis. 27

1 RNA-seq analysis

To capture diverse genes expression, we extracted RNA from 16 tissues of N. 2 cadamba, namely, bark (B), bud, cambium (C), young fruit (FR), old leaves (OL), phloem 3 (P), root (R), young leaves (YL), xylem (primary xylem, PX; transitional xylem, TX; 4 secondary xylem, SX), cambium (transitional cambium, TCA; secondary, SCA) and phloem 5 6 (primary phloem, PPH; transitional phloem, TPH; secondary phloem, SPH) from the first, second and fourth internodes following the RNeasy Plant Mini Kit's (Qiagen, USA) 7 protocol. Clean RNA-seq reads were mapped to N. cadamba reference genome using 8 Hisat2, then the expression level for N. cadamba genes (FPKM, and expression count 9 data) was obtained using HTSeq software(Kim et al., 2015b, Kim et al., 2019). 10

11 Identifying genomic regions responsible for cadambine accumulation by GWAS

The strictosidine and cadambine contents of bark and leaves from 112 individuals 12 were determined by UPLC-Q-TOF (UPLC1290-6540B Q-TOF, Agilent, USA) as 13 14 described in the previous method (LC-MS/MS). GWAS analyses were performed in the mixed linear model using GEMMA (v 0.94.1)(Zhou and Stephens, 2012). The 15 Bonferroni correction was used to reduce the probability of false positives. And, P = 16 10⁻⁵ were was used as the genome-wide thresholds to screen significantly associated 17 sites. The candidate genes were obtained from 5000bp region upstream and 18 downstream of a peak value of significant region. 19

20 Reads mapping and variant calling.

The clean reads were mapped back to the assembled genome using BWA-MEM (v 0.7.8-r455) (Li and Durbin, 2009) with modified parameters: the minimum seed length was set to 32, band width for banded alignment was set to 150 and penalty for a mismatch was set to 3. PCR duplicates were removed using the rmdup function of SAMtools (v 0.1.19-44428cd)(Li *et al.*, 2009) with default parameters. SNPs were detected using SAMtools (v 0.1.19-44428cd) (Li *et al.*, 2009) and bcftools (v 1.3.1)(Narasimhan *et al.*, 2016). The obtained SNPs were further filtered when the This article is protected by copyright. All rights reserved

1 mapped reads <4, the missing data > 0.3 or the Minor Allele frequencies < 0.05

2 **Population structure and Genetic diversity analysis.**

A total of 31 samples were used to analyze the population structure of groups of N. 3 4 cadamba with different levels of cadambine content. The cadambine contents of the samples were determined by UPLC-Q-TOF (UPLC1290-6540B Q-TOF, Agilent, 5 USA) as described in the previous method (LC-MS/MS). All the filtered SNPs were 6 used to construct a neighbor-joining phylogenetic tree using TreeBeST (v 1.9.2) with 7 8 1000 bootstraps and other parameters default. The phylogenetic tree was visualized using Evolview v2 (https://www.evolgenius.info//evolview/). In addition, the 9 principal component analysis (PCA) was performed using the same SNPs set by 10 GCTA (v 1.24.2)(Yang et al., 2011). The fixation index statistic (Fst) and the 11 12 nucleotide diversity (π) were calculated using a 40 kb window in a 20 kb steps for each population by VCFtools (v 0.1.14). The selective regions were detected by the 13 top 5% of the low or high $\log_2 \pi$ ratio and the top 5% of Fst values. The linkage 14 disequilibrium (LD) of each population was estimated by calculating squared 15 correlation coefficient (r²) values between any two SNP sites in 500 kb using 16 Haploview (v 4.2)(Barrett et al., 2005). 17

18 Data availability

All raw and processed sequencing data generated in this study have been submitted to 19 the NCBI BioProject database under accession number PRJNA650253. The raw 20 genome sequencing data obtained by Illumina and PacBio platform have been 21 the **NCBI** 22 submitted to BioSample database under accession number SAMN15700858. The raw sequencing data of the resequencing data and 23 transcriptome have been submitted to the NCBI BioSample database under accession 24 number SAMN15700860 and SAMN15700859, respectively. The genome annotation 25 and assembled genome sequences can be downloaded from the webpage 26 27 (https://figshare.com/s/ed20e0e82a4e7474396b).

1 Author contributions

C.P., X.Z. and X.C. designed and managed the project. C.P., X.Z. and X.H. wrote the
manuscript with input from all authors. K.O., Q.Q., J. L. and J. Z. collected samples.
T.Z. and J.G. performed the *in vitro* catalytic assay. J.Y., T.Z. and X.H. performed
HPLC-MS/MS experiments. K.O., X.W., S.Y. and L.Z. performed RNA-seq analyses
of gene expression. S.L. and W.G. performed FISH analysis; X.H., W.J., B. L., K.O.
and X. Z. worked on sequencing and data analyzing. C.P., X.Z., X.H., E.N., X.C. and
K.O. revised the manuscript. All the authors read and approved the final manuscript.

9 Acknowledgements

10 We thank Y. Shang for critical reading of the manuscript. We thank X.S. Hu, L.Z. Gao and Q.G. Liao for helpful discussions. We thank Q. Hu, X. Zhang and S. Xiao 11 for assistance with HPLC-MS/MS experiments. This work was supported by funds 12 from the National Natural Science Foundation of China (Grant No.31470681, 13 14 31970197); the National Key Research and Development Program of China (Grant No.2016YFD0600104); Natural Science Foundation of Guangdong Province of China 15 (Grant No. 2016A030311032); Guangzhou Science and Technology Program (Grant 16 No. 201607020024) and Foundation of Young Creative Talents in Higher Education 17 of Guangdong Province (Grant No.2017KQNCX017). 18

19 Conflict of interest

20 The authors declare no competing interests.

21 22

23

24

25

26

27

1 Supplemental Figures

2 Figure S1. A. 17 K-mer analysis for estimating the genome size of *Neolamarckia cadamba* B.

3 Cytological analysis of *Neolamarckia cadamba* metaphase chromosomes by FISH using 45s

4 rDNA as probe. 45S rDNA was labeled with Chroma Tide Alexa Fluor 488 (green signal), and

5 the chromosomes were counterstained with DAPI (blue). Bars = $10 \mu m$

6 Figure S2. Estimation of synonymous substitutions per site (Ks) in N. cadamba genome

7 Figure S3 Characterization of standard compounds by Q-TOF LC-MS/MS and NMR. 8 Tryptamine, epoxystrictosidine, 3a-dihydrocadambine and cadambine used as standards (see 9 "Methods") were analyzed by Q-TOF LC-MS/MS. Cadambine was further tested by NMR to 10 verify its structural formula. NMR spectra were carried out on a Bruker AV 400 spectrometer in 11 CD₃OD using tetramethylsilane (TMS, $\delta = 0$) as internal reference. ¹H NMR(CD3OD, 400MHz) δ = 7.59 (1H, s), δ = 7.48 (1H, d), δ = 7.34 (1H, d), δ = 7.12 (1H, t), δ = 7.01 (1H, t), δ = 5.85 12 (1H, d), $\delta = 4.96$ (1H, d), $\delta = 4.81$ (1H, d), $\delta = 3.90$ (1H, m), $\delta = 3.53$ (1H, m), $\delta = 3.63$ 13 $(3H, s), \delta = 3.56 (1H, m), \delta = 3.42 (1H, m), \delta = 3.18 (1H, m), \delta = 3.03 (1H, m), \delta = 2.08$ 14 15 (2H, m), $\delta = 1.79$ (1H, m). In accordance with the results of Handa *et al.* (1983) and Xu *et al.* 16 (2011), this compound is characterized as cadambine. **a-d** MS/MS spectrum of tryptamine (a), cadambine (b), 3a-dihydrocadambine (c) and epoxystrictosidine (d). e NMR spectrum of 17 18 cadambine. f The structural formula of cadambine.

Figure S4. Phylogenetic tree of STR genes from six plant species. The genes from
Arabidopsis thaliana (AT), Catharanthus roseus (CRO), Coffea canephora (Cc), Camptotheca
acuminate (Cac), N. cadamba (evm.model) and Theobroma cacao (EOY).

Figure S 5. *In vitro* expressed fusion NcSTR1. Left, SDS-PAGE of NcSTR1 expressed in *E. coli* and the strain STR1-2 successfully expressed His-tagged NcSTR1; Middle, *in vitro* expressed
fusion NcSTR1 validated by immunoblotting; Right, SDS-PAGE of purified fusion NcSTR1 by
His-trap.

Figure S6 Chemical structures of strictosidine, 3a-dihydrocadambine cadambine and
 predicted epoxystrictosidine.

Figure S7 The expression profile of all the predicted biosynthetic genes in Shikimatepathway.

- 1 Figure S8 The expression profile of all the predicted squalene epoxidase genes (SQEs).
- 2 Figure S9. Phylogenetic tree of NcCYPs in BGCs and 52 functionally characterized plant
- **3** CYP family members involved in terpenoid biosynthesis
- 4 Figure S10. Protein alignment of CrLAMT, CcNMTs and their respective Neolamarckia
- 5 paralogues.
- 6 Figure S11 The structure of tandem duplicated NcDLGTs.
- 7 Figure S12 Heat Map of Expression Data of genes in NcMCL24 and NcMCL25.
- 8 Figure S13 Duplication of NcMYC TFs in tandem order.
- 9 Figure S14 Protein alignment of CrSGD, RsSGD, RvSGD, RsRG and their respective
- 10 Neolamarckia and Ophiorrhiza paralogues. Accession numbers: CrSGD (AAF28800.1),
- 11 RsSGD (CAC83098.1), RsRG (AAF03675.1) and RvSGD (AFI71457). Previously reported
- 12 active-site residues of RsSGD and RsRG were framed or showed with bright colors.
- 13

14 Supplemental Tables

- 15 Table S1. Survey statistic results of *Neolamarckia cadamba*.
- 16 Table S2. Sequencing data statistics of *Neolamarckia cadamba*.
- 17 Table S3. Coverage statistics of *Neolamarckia cadamba* genome
- 18 Table S4. Assessment the gene coverage rate using CEGMA and BUSCO
- 19 Table S5. Assessment the genome assembly using RNA-seq
- 20 Table S6. The statistical results of gene function annotation of *Neolamarckia cadamba* 21 genome.
- 22 Table S7. The statistical results of non-coding RNA of *Neolamarckia cadamba* genome.
- 23 Table S8. Summary of Repeat contents in *Neolamarckia cadamba* genome.
- 24 Table S9. Genes used for gene family clustering in each species.
- 25 Table S10. Enrichment analysis of duplicated genes from the recent WGD.
- 26 Table S11 Gene lists under positive selection.
- 27 Table S12 The predicted biosynthetic genes in five plant species.
- 28 Table S13. Short lists of Neolamarckia cadamba candidate genes involved in MVA/MEP,

1	shikimate/indole pathway and iridoid biosynthesis pathway and 35 putative NcSTRs.
2	Table S14 Unique gene lists underwent recent WGD event.
3	Table S 15 Gene lists in the association network of gene-to-metabolite.
4	Table S 16 The sampling sites of 112 individuals in GWAS
5	Table S17 Putative Iridoid pathway genes in Neolamarckia cadamba associated with loci
6	responsible for MIA-production in GWAS
7	Table S18 The putative Neolamarckia cadamba cytochrome P450s associated MIA -
8	production in GWAS
9	Table S19 KEGG analysis of selected genes from population analysis.
10	Table S20 The 67 biosynthetic gene clusters of Neolamarckia cadamba predicted with
11	plantiSMASH algorithm
12	Table S21 The detailed data for the predicted biosynthetic gene clusters and tandem
13	duplicates of Neolamarckia cadamba
14	
15	
16	
17	
18	
19	
20	
21	
22	
23	
24	
25	
26	
27	
28	
29	
30	

1

2

3 References

- ALMEIDA, A., DONG, A. L., KHAKIMOV, A. B., BASSARD, J. E. & MOSES, A. T. 2018. A Single
 Oxidosqualene Cyclase Produces the Seco-Triterpenoid a-Onocerin. Plant Physiology 176 (2):
 1469–84.
 BARLEBEN, L., PANJIKAR, S., RUPPERT, M., KOEPKE, J. & STÖCKIGT, J. 2007. Molecular architecture of
- 8 strictosidine glucosidase: the gateway to the biosynthesis of the monoterpenoid indole
 9 alkaloid family. *Plant Cell*, 19, 2886-97.
- BARRETT, J. C., FRY, B., MALLER, J. & DALY, M. J. 2005. Haploview: analysis and visualization of LD and
 haplotype maps. *Bioinformatics*, 21, 263-5.
- BEDI, Y. S., BIR, S. S. & GILL, B. S. 1981. Cytopalynology of Woody Taxa of Family Rubiaceae from
 North and Central India. Proceedings of the indian national science academy.part b.biological
 sciences.47 (6):708-715.
- BROWN, S., CLASTRE, M., COURDAVAULT, V. & O'CONNOR, S. E. 2015. De novo production of the
 plant-derived alkaloid strictosidine in yeast. *Proceedings of the National Academy of Sciences*, 112, 3205-3210.
- BURTON, J. N., ANDREW, A., PATWARDHAN, R. P., RUOLAN, Q., KITZMAN, J. O. & JAY, S. 2013.
 Chromosome-scale scaffolding of de novo genome assemblies based on chromatin
 interactions. *Nature Biotechnology*, 31, 1119.
- CARQUEIJEIRO, I., BROWN, S., CHUNG, K., DANG, T. T., WALIA, M., BESSEAU, S., DUGE DE
 BERNONVILLE, T., OUDIN, A., LANOUE, A., BILLET, K., MUNSCH, T., KOUDOUNAS, K., MELIN,
 C., GODON, C., RAZAFIMANDIMBY, B., DE CRAENE, J. O., GLEVAREC, G., MARC, J.,
 GIGLIOLI-GUIVARC'H, N., CLASTRE, M., ST-PIERRE, B., PAPON, N., ANDRADE, R. B.,
 O'CONNOR, S. E. & COURDAVAULT, V. 2018a. Two Tabersonine 6,7-Epoxidases Initiate
 Lochnericine-Derived Alkaloid Biosynthesis in Catharanthus roseus. *Plant Physiol*, 177,
 1473-1486.
- CARQUEIJEIRO, I., DUGE DE BERNONVILLE, T., LANOUE, A., DANG, T. T., TEIJARO, C. N., PAETZ, C.,
 BILLET, K., MOSQUERA, A., OUDIN, A., BESSEAU, S., PAPON, N., GLEVAREC, G., ATEHORTUA,
 This article is protected by copyright. All rights reserved

L., CLASTRE, M., GIGLIOLI-GUIVARC'H, N., SCHNEIDER, B., ST-PIERRE, B., ANDRADE, R. B.,
 O'CONNOR, S. E. & COURDAVAULT, V. 2018b. A BAHD acyltransferase catalyzing
 19-O-acetylation of tabersonine derivatives in roots of Catharanthus roseus enables
 combinatorial synthesis of monoterpene indole alkaloids. *Plant J*, 94, 469-484.

- 5 CHAE, L., KIM, T., NILO-POYANCO, R. & RHEE, S. Y. 2014. Genomic signatures of specialized
 6 metabolism in plants. *Science*, 344, 510-513.
- CHANDEL, M., KUMAR, M., SHARMA, U., SINGH, B. & KAUR, S. 2017. Antioxidant, Antigenotoxic and
 Cytotoxic Activity of Anthocephalus cadamba (Roxb.) Miq. Bark Fractions and their
 Phytochemical Analysis using UPLC-ESI-QTOF-MS. *Combinatorial chemistry & high throughput screening*, 20, 760-772.
- 11 CHANDEL, M., SHARMA, U., KUMAR, N., SINGH, B. & KAUR, S. 2012. Antioxidant activity and 12 identification of bioactive compounds from leaves of *Anthocephalus cadamba* by 13 ultra-performance liquid chromatography/electrospray ionization quadrupole time of flight 14 mass spectrometry. *Asian Pacific Journal of Tropical Medicine*, 5, 977-985.
- CHANDEL, M., SHARMA, U., KUMAR, N., SINGH, B. & KAUR, S. 2014. In vitro studies on the
 antioxidant/antigenotoxic potential of aqueous fraction from Anthocephalus cadamba bark.
 Perspectives in Cancer Prevention-Translational Cancer Research. Springer.
- CHEN, Z., TIAN, Z., ZHANG, Y., FENG, X., LI, Y. & JIANG, H. 2020. Monoterpene indole alkaloids in
 Uncaria rhynchophlly (Miq.) Jacks chinensis and their chemotaxonomic significance.
 Biochemical Systematics and Ecology, 91, 104057.
- 21 CHIN, C.-S., PELUSO, P., SEDLAZECK, F. J., NATTESTAD, M., CONCEPCION, G. T., CLUM, A., DUNN, C.,
- O'MALLEY, R., FIGUEROA-BALDERAS, R. & MORALES-CRUZ, A. 2016. Phased diploid genome
 assembly with single-molecule real-time sequencing. *Nature methods*, 13, 1050.
- COLINAS, M. & GOOSSENS, A. 2018. Combinatorial Transcriptional Control of Plant Specialized
 Metabolism. *Trends Plant Sci*, 23, 324-336.
- DE BIE, T., CRISTIANINI, N., DEMUTH, J. P. & HAHN, M. W. 2006. CAFE: a computational tool for the
 study of gene family evolution. *Bioinformatics*, 22, 1269-71.
- De LUCA, V. D. & CUTLER, A. J. 1987. Subcellular Localization of Enzymes Involved in Indole Alkaloid
 Biosynthesis in *Catharanthus roseus*. *Plant Physiology*, 85, 1099-1102.
- 30 De LUCA, V. D., SALIM, V., THAMM, A., MASADA, S. A. & YU, F. 2014. Making iridoids/secoiridoids and This article is protected by copyright. All rights reserved

1

2

monoterpenoid indole alkaloids: progress on pathway elucidation. *Current Opinion in Plant Biology*, 19, 35-42.

DENG, C., QIN, R., GAO, J., CAO, Y., LI, S., GAO, W. & LU, L. 2012. Identification of sex chromosome of
spinach by physical mapping of 45s rDNAs by FISH. *Caryologia*, 65, 322-327.
DENOEUD, F., CARRETERO-PAULET, L., DEREEPER, A., DROC, G., GUYOT, R., PIETRELLA, M., ZHENG, C.
F., ALBERTI, A., ANTHONY, F., APREA, G., AURY, J. M., BENTO, P., BERNARD, M., BOCS, S.,
CAMPA, C., CENCI, A., COMBES, M. C., CROUZILLAT, D., DA SILVA, C., DADDIEGO, L., DE
BELLIS, F., DUSSERT, S., GARSMEUR, O., GAYRAUD, T., GUIGNON, V., JAHN, K., JAMILLOUX,
V., JOET, T., LABADIE, K., LAN, T. Y., LECLERCQ, J., LEPELLEY, M., LEROY, T., LI, L. T., LIBRADO,

P., LOPEZ, L., MUNOZ, A., NOEL, B., PALLAVICINI, A., PERROTTA, G., PONCET, V., POT, D.,
 PRIYONO, RIGOREAU, M., ROUARD, M., ROZAS, J., TRANCHANT-DUBREUIL, C., VANBUREN,
 R., ZHANG, Q., ANDRADE, A. C., ARGOUT, X., BERTRAND, B., DE KOCHKO, A., GRAZIOSI, G.,
 HENRY, R. J., JAYARAMA, MING, R., NAGAI, C., ROUNSLEY, S., SANKOFF, D., GIULIANO, G.,
 ALBERT, V. A., WINCKER, P. & LASHERMES, P. 2014. The coffee genome provides insight into
 the convergent evolution of caffeine biosynthesis. *Science*, 345, 1181-1184.

- DWEVEDI, A., SHARMA, K. & SHARMA, Y. K. 2014. Cadamba: A miraculous tree having enormous
 pharmacological implications. *Pharmacognosy Reviews*, 9, 107-113.
- EDGAR, R. C. 2004. MUSCLE: multiple sequence alignment with high accuracy and high throughput.
 Nucleic Acids Res, 32, 1792-7.
- EGER, E., SIMON, A., SHARMA, M., YANG, S., BREUKELAAR, W. B., GROGAN, G., HOUK, K. & KROUTIL,
 W. 2020. Inverted binding of non-natural substrates in strictosidine synthase leads to a
 switch of stereochemical outcome in enzyme-catalyzed Pictet–Spengler reactions. *Journal of the American Chemical Society*, 142, 792-800.
- HANDA, S. S., BORRIS, R. P., CORDELL, G. A. & PHILLIPSON, J. D. 1983. NMR spectral analysis of
 cadambine from *Anthocephalus chinensis*. *Journal of natural products*, 46, 325-330.
- HU, L., XU, Z., WANG, M., FAN, R., YUAN, D., WU, B., WU, H., QIN, X., YAN, L., TAN, L., SIM, S., LI, W.,
 SASKI, C. A., DANIELL, H., WENDEL, J. F., LINDSEY, K., ZHANG, X., HAO, C. & JIN, S. 2019. The
 chromosome-scale reference genome of black pepper provides insight into piperine
 biosynthesis. *Nature Communications*, 10, 4702.
- 30 JACOBOWITZ, J. R. & WENG, J.-K. 2020. Exploring uncharted territories of plant specialized This article is protected by copyright. All rights reserved

1	metabolism in the postgenomic era. Annual review of plant biology, 71, 631-658.
2	KAI, G., WU, C., GEN, L., ZHANG, L., CUI, L. & NI, X. 2015. Biosynthesis and biotechnological production
3	of anti-cancer drug Camptothecin. Phytochemistry Reviews, 14, 525-539.
4	KARETI, S. R. & SUBASH, P. 2020. In silico exploration of anti-Alzheimer's compounds present in
5	methanolic extract of Neolamarckia cadamba bark using GC-MS/MS. Arabian Journal of
6	Chemistry, 13, 6246-6255.
7	KAUTSAR, S. A., SUAREZ DURAN, H. G., BLIN, K., OSBOURN, A. & MEDEMA, M. H. 2017. plantiSMASH:
8	automated identification, annotation and expression analysis of plant biosynthetic gene
9	clusters. Nucleic acids research, 45, W55-W63.
10	KELLNER, F., KIM, J., CLAVIJO, B. J., HAMILTON, J. P., CHILDS, K. L., VAILLANCOURT, B., CEPELA, J.,
11	HABERMANN, M., STEUERNAGEL, B. & CLISSOLD, L. 2015. Genome - guided investigation of
12	plant natural product biosynthesis. The Plant Journal, 82, 680-692.
13	KETUDAT CAIRNS, J. R. & ESEN, A. 2010. β-Glucosidases. <i>Cell Mol Life Sci,</i> 67, 3389-405.
14	KIM, D., LANGMEAD, B. & SALZBERG, S. L. 2015a. HISAT: a fast spliced aligner with low memory
15	requirements. Nature methods, 12, 357-360.
16	KIM, D., LANGMEAD, B. & SALZBERG, S. L. 2015b. HISAT: a fast spliced aligner with low memory
17	requirements. Nat Methods, 12, 357-60.
18	KIM, D., PAGGI, J. M., PARK, C., BENNETT, C. & SALZBERG, S. L. 2019. Graph-based genome alignment
19	and genotyping with HISAT2 and HISAT-genotype. Nature biotechnology, 37, 907-915.
20	KIM, D., PERTEA, G., TRAPNELL, C., PIMENTEL, H., KELLEY, R. & SALZBERG, S. L. 2013. TopHat2:
21	accurate alignment of transcriptomes in the presence of insertions, deletions and gene
22	fusions. Genome Biol, 14, (4): 1-13.
23	LEONELLI, L., BROOKS, M. & NIYOGI, K. K. 2017. Engineering the lutein epoxide cycle into Arabidopsis
24	thaliana. Proceedings of the National Academy of ences, 114, 201704373.
25	LI, H. & DURBIN, R. 2009. Fast and accurate short read alignment with Burrows-Wheeler transform.
26	Bioinformatics, 25, 1754-60.
27	LI, H., HANDSAKER, B., WYSOKER, A., FENNELL, T., RUAN, J., HOMER, N., MARTH, G., ABECASIS, G.,
28	DURBIN, R. & SUBGROUP, G. P. D. P. 2009. The Sequence Alignment/Map format and
29	SAMtools. Bioinformatics, 25, 2078-2079.
30	LI, J., ZHANG, D., OUYANG, K. & CHEN, X. 2018. The complete chloroplast genome of the miracle tree
	This article is protected by copyright. All rights reserved

1	Neolamarckia cadamba and its comparison in Rubiaceae family. Biotechnology &
2	Biotechnological Equipment, 32, 1087-1097.
3	LI, J., ZHANG, D., OUYANG, K. & CHEN, X. 2019. High frequency plant regeneration from leaf culture of
4	Neolamarckia cadamba. Plant Biotechnology, 18.1119 a.
5	LI, L., STOECKERT, C. J., JR. & ROOS, D. S. 2003. OrthoMCL: identification of ortholog groups for
6	eukaryotic genomes. <i>Genome Res,</i> 13, 2178-89.
7	Li, M., Tian, S., Jin, L., Zhou, G., Li, Y., Zhang, Y., Ma, J. 2013. Genomic analyses identify distinct
8	patterns of selection in domesticated pigs and Tibetan wild boars. Nature Genetics, 45(12),
9	1431–1438.
10	LI, Y., LEVEAU, A., ZHAO, Q., FENG, Q. & OSBOURN, A. 2021. Subtelomeric assembly of a multi-gene
11	pathway for antimicrobial defense compounds in cereals. Nature communications, 12, 2563.
12	MA, D. & CONSTABEL, C. P. 2019. MYB Repressors as Regulators of Phenylpropanoid Metabolism in
13	Plants. Trends Plant Sci, 24, 275-289.
14	MA, X., KOEPKE, J., PANJIKAR, S., FRITZSCH, G. & STOCKIGT, J. 2005. Crystal structure of vinorine
15	synthase, the first representative of the BAHD superfamily. J Biol Chem, 280, 13576-83.
16	MA, X., PANJIKAR, S., KOEPKE, J., LORIS, E. & STÖCKIGT, J. 2006. The structure of Rauvolfia serpentina
17	strictosidine synthase is a novel six-bladed β -propeller fold in plant proteins. The Plant Cell,
18	18, 907-920.
19	Mccarthy A A, & Mccarthy J G. 2007. The structure of two N-methyltransferasesfrom the caffeine
20	biosynthetic pathway.P1ant physiology,144(2): 879-889.
21	MEHRA, P. & BAWA, K. 1969. Chromosomal evolution in tropical hardwoods. Evolution, 466-481.
22	MEIJER, A. H., CARDOSO, M., VOSKUILEN, J. T., WAAL, A. D., VERPOORTE, R. & HOGE, J. 1993. Isolation
23	and characterization of a cDNA clone from Catharanthus roseus encoding
24	NADPH:cytochrome P-450 reductase, an enzyme essential for reactions catalysed by
25	cytochrome P-450 mono-oxygenases in plants. <i>The Plant Journal</i> .4 (1):47-60.
26	NARASIMHAN, V., DANECEK, P., SCALLY, A., XUE, Y., TYLER-SMITH, C. & DURBIN, R. 2016.
27	BCFtools/RoH: a hidden Markov model approach for detecting autozygosity from
28	next-generation sequencing data. Bioinformatics, 32, 1749-51.
29	NüTZMANN, H. W., HUANG, A. & OSBOURN, A. 2016. Plant metabolic clusters-from genetics to
30	genomics. New phytologist, 211, 771-789.

1	PANDEY, A. & NEGI, P. S. 2016. Traditional uses, phytochemistry and pharmacological properties of
2	Neolamarckia cadamba: A review. Journal of ethnopharmacology, 181, 118-135.
3	PARRA, G., BRADNAM, K. & KORF, I. 2007. CEGMA: a pipeline to accurately annotate core genes in
4	eukaryotic genomes. <i>Bioinformatics,</i> 23, 1061-7.
5	PATERSON, A. H., BOWERS, J. E. & CHAPMAN, B. A. 2004. Ancient polyploidization predating
6	divergence of the cereals, and its consequences for comparative genomics. Proc Natl Acad
7	Sci U S A, 101, 9903-8.
8	PRESSNITZ, D., FISCHEREDER, E. M., PLETZ, J., KOFLER, C., HAMMERER, L., HIEBLER, K., LECHNER, H.,
9	RICHTER, N., EGER, E. & KROUTIL, W. 2018. Asymmetric Synthesis of (R) - 1 - Alkyl -
10	Substituted Tetrahydro - β - carbolines Catalyzed by Strictosidine Synthases. Angewandte
11	Chemie, 130, 10843-10847.
12	PUTTICK, M. N. 2019. MCMCtreeR: functions to prepare MCMCtree analyses and visualize posterior
13	ages on trees. Bioinformatics, 35, 5321-5322.
14	QU, Y., EASSON, M. E., SIMIONESCU, R., HAJICEK, J., THAMM, A. M., SALIM, V. & DE LUCA, V. 2018.
15	Solution of the multistep pathway for assembly of corynanthean, strychnos, iboga, and
16	aspidosperma monoterpenoid indole alkaloids from 19E-geissoschizine. Proceedings of the
17	National Academy of Sciences, 115, 3180-3185.
18	RAI, A., HIRAKAWA, H., NAKABAYASHI, R., KIKUCHI, S., HAYASHI, K., RAI, M., TSUGAWA, H., NAKAYA,
19	T., MORI, T. & NAGASAKI, H. 2021. Chromosome-level genome assembly of Ophiorrhiza
20	pumila reveals the evolution of camptothecin biosynthesis. Nature communications, 12,
21	1-19.
22	RAZAFIMANDIMBISON, S. G. 2002. A Systematic Revision of Breonia (Rubiaceae-Naucleeae). Annals of
23	the Missouri Botanical Garden, 89, 1-37.
24	ROBBRECHT, E. & MANEN, J. F. 2006. The Major Evolutionary Lineages of the Coffee Family
25	(Rubiaceae, Angiosperms). Combined Analysis (nDNA and cpDNA) to Infer the Position of
26	Coptosapelta and Luculia, and Supertree Construction Based on rbcL, rps16, trnL-trnF and
27	atpB-rbcL Data. A New Classi. Syst.geogr.pl, 76, 85-145.
28	ROSSELEENA, R. E., MOTOAKI, C., MIKI, K., TAKASHI, A., YOSHIMI, O., NOBUTAKA, M., MASARU, O. T.,
29	ATSUSHI, F., AMIT, R. & KAZUKI, S. 2016. An MYB transcription factor regulating specialized
30	metabolisms in Ophiorrhiza pumila. Plant Biotechnology, 33, 1-9.

1	SABIR, J. S., JANSEN, R. K., ARASAPPAN, D., CALDERON, V., NOUTAHI, E., ZHENG, C., PARK, S., SABIR,
2	M. J., BAESHEN, M. N. & HAJRAH, N. H. 2016. The nuclear genome of Rhazya stricta and the
3	evolution of alkaloid diversity in a medically relevant clade of Apocynaceae. Scientific reports,
4	6, 1-10.
5	SADRE, R., MAGALLANES-LUNDBACK, M., PRADHAN, S., SALIM, V., MESBERG, A., JONES, A. D. &
6	DELLAPENNA, D. 2016. Metabolite diversity in alkaloid biosynthesis: a multilane
7	(diastereomer) highway for camptothecin synthesis in Camptotheca acuminata. The Plant
8	<i>Cell,</i> 28, 1926-1944.
9	SALIM, V., YU, F., ALTAREJOS, J. & DE LUCA, V. 2013. Virus - induced gene silencing identifies
10	Catharanthus roseus 7 - deoxyloganic acid - 7 - hydroxylase, a step in iridoid and
11	monoterpene indole alkaloid biosynthesis. The Plant Journal, 76, 754-765.
12	SANTOS, C. L. G., K.O.G., N., COSTA E.V., SOUZA A.D.L., M.L.B., P., H.H.F., K. & F.M.A, S. 2020.
13	Molecular networking - based dereplication of strictosidine - derived monoterpene indole
14	alkaloids from the curare ingredient Strychnos peckii. Rapid Communications in Mass
15	Spectrometry, 34, e8683.
16	SCHLäPFER, P., ZHANG, P., WANG, C., KIM, T., BANF, M., CHAE, L., DREHER, K., CHAVALI, A. K.,
17	NILO-POYANCO, R. & BERNARD, T. 2017. Genome-wide prediction of metabolic enzymes,
18	pathways, and gene clusters in plants. <i>Plant physiology</i> , 173, 2041-2059.
19	SCOTT, M. G. & MADDEN, T. L. 2004. BLAST: at the core of a powerful and diverse set of sequence
20	analysis tools. Nucleic Acids Research, 32: 20-25.
21	SHANG, Y. & HUANG, S. 2020. Engineering plant cytochrome P450s for enhanced synthesis of natural
22	products: past achievements and future perspectives. Plant Communications, 1, 100012.
23	SIMAO, F. A., WATERHOUSE, R. M., IOANNIDIS, P., KRIVENTSEVA, E. V. & ZDOBNOV, E. M. 2015.
24	BUSCO: assessing genome assembly and annotation completeness with single-copy
25	orthologs. Bioinformatics, 31, 3210-2.
26	SINGH, S. K., PATRA, B., PAUL, P., LIU, Y., PATTANAIK, S. & YUAN, L. 2021. BHLH IRIDOID SYNTHESIS 3
27	is a member of a bHLH gene cluster regulating terpenoid indole alkaloid biosynthesis in
28	Catharanthus roseus. <i>Plant Direct</i> , 5, e00305.
29	STAMATAKIS, A. 2014. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large
30	phylogenies. <i>Bioinformatics,</i> 30, 1312-3.
	This article is protected by copyright. All rights reserved

1	TAKAYAMA, H., TSUTSUMI, SI., KITAJIMA, M., SANTIARWORN, D., LIAWRUANGRATH, B. & AIMI, N.
2	2003. Gluco-indole alkaloids from Nauclea cadamba in thailand and transformation of
3	3α -dihydrocadambine into the indolopyridine alkaloid, 16-carbomethoxynaufoline. Chemical
4	and pharmaceutical bulletin, 51, 232-233.
5	TATSIS, E. C., CA RQUEIJEIRO, I., THOMAS, D., FRANKE, J., DANG, T. T. T., OUDIN, A., LANOUE, A.,
6	LAFONTAINE, F., STAVRINIDES, A. K. & CLASTRE, M. 2017. A three enzyme system to
7	generate the Strychnos alkaloid scaffold from a central biosynthetic intermediate. Nature
8	Communications, 8, 316.
9	TRAN, H. T., RAMARAJ, T., FURTADO, A., LEE, L. S. & HENRY, R. J. 2018. Use of a draft genome of
10	coffee (Coffea arabica) to identify SNP s associated with caffeine content. Plant
11	biotechnology journal, 16, 1756-1766.
12	TRAPNELL, C., PACHTER, L. & SALZBERG, S. L. 2009. TopHat: discovering splice junctions with RNA-Seq.
13	Bioinformatics, 25, 1105-11.
14	TRAPNELL, C., ROBERTS, A., GOFF, L., PERTEA, G., KIM, D., KELLEY, D. R., PIMENTEL, H., SALZBERG, S.
15	L., RINN, J. L. & PACHTER, L. 2012. Differential gene and transcript expression analysis of
16	RNA-seq experiments with TopHat and Cufflinks. Nature Protocols, 7, 562-578.
17	URLACHER, V. B. & GIRHARD, M. 2019. Cytochrome P450 monooxygenases in biotechnology and
18	synthetic biology. Trends in biotechnology, 37, 882-897.
19	VAN MOERKERCKE, A., STEENSMA, P., GARIBOLDI, I., ESPOZ, J., PURNAMA, P. C., SCHWEIZER, F.,
20	MIETTINEN, K., VANDEN BOSSCHE, R., DE CLERCQ, R., MEMELINK, J. & GOOSSENS, A. 2016.
21	The basic helix-loop-helix transcription factor BIS2 is essential for monoterpenoid indole
22	alkaloid production in the medicinal plant Catharanthus roseus. Plant J, 88, 3-12.
23	VAN MOERKERCKE, A., STEENSMA, P., SCHWEIZER, F., POLLIER, J., GARIBOLDI, I., PAYNE, R., VANDEN
24	BOSSCHE, R., MIETTINEN, K., ESPOZ, J., PURNAMA, P. C., KELLNER, F., SEPPANEN-LAAKSO, T.,
25	O'CONNOR, S. E., RISCHER, H., MEMELINK, J. & GOOSSENS, A. 2015. The bHLH transcription
26	factor BIS1 controls the iridoid branch of the monoterpenoid indole alkaloid pathway in
27	Catharanthus roseus. Proc Natl Acad Sci U S A, 112, 8130-5.
28	WANG, C., WU, C., WANG, Y., XIE, C., SHI, M., NILE, S., ZHOU, Z. & KAI, G. 2019. Transcription factor
29	OpWRKY3 is involved in the development and biosynthesis of camptothecin and its
30	precursors in Ophiorrhiza pumila hairy roots. International journal of molecular sciences, 20,
	This article is protected by copyright. All rights reserved

1 3996.

2	WANG, DP., WAN, HL., ZHANG, S. & YU, J. 2009. Gamma-MYN: a new algorithm for estimating Ka
3	and Ks with consideration of variable substitution rates. <i>Biology direct</i> , 4 (1): 20-20
4	WANG, Y., TANG, H., DEBARRY, J. D., TAN, X., LI, J., WANG, X., LEE, T. H., JIN, H., MARLER, B., GUO, H.,
5	KISSINGER, J. C. & PATERSON, A. H. 2012. MCScanX: a toolkit for detection and evolutionary
6	analysis of gene synteny and collinearity. Nucleic Acids Res, 40, e49.
7	WEIR, B. S. & COCKERHAM, C. C. 1984. ESTIMATING F-STATISTICS FOR THE ANALYSIS OF POPULATION
8	STRUCTURE. Evolution, 38, 1358-1370.
9	WU, X. D., WANG, L., HE, J., LI, X. Y., DONG, L. B., GONG, X., GAO, X., SONG, L. D., LI, Y. & PENG, L. Y.
10	2013. Two new indole alkaloids from Emmenopterys henryi. Helvetica Chimica Acta, 96,
11	2207-2213.
12	XIA, L., RUPPERT, M., WANG, M., PANJIKAR, S., LIN, H., RAJENDRAN, C., BARLEBEN, L. & STOCKIGT, J.
13	2012. Structures of alkaloid biosynthetic glucosidases decode substrate specificity. ACS Chem
14	<i>Biol,</i> 7, 226-34.
15	Yamazaki, Y., Sudo, H., Yamazaki, M., Aimi, N., Saito, K., 2003. Camptothecin biosynthetic genes in
16	hairy roots of Ophiorrhiza pumila: cloning, characterization and differential expression in
17	tissues and by stress compounds. Plant Cell Physiol. 44, 395–403
18	YANG, J., LEE, S. H., GODDARD, M. E. & VISSCHER, P. M. 2011. GCTA: a tool for genome-wide complex
19	trait analysis. Am J Hum Genet, 88, 76-82.
20	YANG, Z. 2007. PAML 4: phylogenetic analysis by maximum likelihood. <i>Mol Biol Evol</i> , 24, 1586-91.
21	YUAN, HL., ZHAO, YL., QIN, XJ., LIU, YP., YU, HF., ZHU, PF., JIN, Q., YANG, XW. & LUO, XD.
22	2020. Anti-inflammatory and analgesic activities of Neolamarckia cadamba and its bioactive
23	monoterpenoid indole alkaloids. Journal of ethnopharmacology, 260, 113103.
24	ZHANG, H., HEDHILI, S., MONTIEL, G., ZHANG, Y., CHATEL, G., PRÉ, M., GANTET, P. & MEMELINK, J.
25	2011. The basic helix-loop-helix transcription factor CrMYC2 controls the
26	jasmonate-responsive expression of the ORCA genes that regulate alkaloid biosynthesis in
27	Catharanthus roseus. Plant J, 67, 61-71.
28	ZHENG, X., LI, P. & LU, X. 2019. Research advances in cytochrome P450-catalysed pharmaceutical
29	terpenoid biosynthesis in plants. Journal of experimental botany, 70, 4619-4630.

ZHOU, X. & STEPHENS, M. 2012. Genome-wide efficient mixed-model analysis for association studies.This article is protected by copyright. All rights reserved

1	Nature Genetics, 44, 821-824.
2	ZHOU, X. & STEPHENS, M. 2014. Efficient multivariate linear mixed model algorithms for
3	genome-wide association studies. Nat Methods, 11, 407-9.
4	
5	
6	Table 1. Global statistics of <i>N. cadamba</i> genome assembly and annotation

	Number	Size
Genome Assembly		
Total contigs	2881	741.90 Mb
Contig N50	225	824.14 Kb
Contig N90	1023	136.8 Kb
Total scaffolds	807	744.45 Mb
Scaffold N50	11	29.20 Mb
Scaffold N90	22	24.51 Mb
Pseudochromosomes	22	744.5 Mb
Genome Annotation		
Predicted protein-coding genes	35,461	
Average gene length (bp)		3,489.6 bp
Average CDS length (bp)		1,151.7 bp
Average exons per gene	4.7	
Average exon length (bp)		245.1 bp
Average intron length (bp)		632.0 bp

- ___

This article is protected by copyright. All rights reserved

1 2 3 4 5 6 7 8 9 10 Figure 1 Hi-C map and overview of the genomic features of the N. cadamba 22 11 pseudochromosomes and evolutionary analyses. a Hi-C map of the N. cadamba 12 genome showing genome-wide all-by-all interactions. b Characteristics of the 22 13 chromosomes of N. cadamba. From outermost to innermost: (i) Circular 14 representation of the 22 chromosomes. (ii-vi) Densities of transposable element, gene, 15 16 GC, miRNA and tRNA. (vii) Densities of snRNA density and syntenic blocks. Densities were calculated in 100Kb windows. C Phylogenetic tree with 402 17 single-copy orthologs from 14 species identified by OrthoMCL to show divergence 18 times and expanded/contracted gene families. d Distribution of 4DTv distance of 19 homologous genes from N. cadamba (Nca), C. canephora (Cca) and S. tuberosum 20 (Stu). e Synteny blocks between N. cadamba, C. canephora and O. pumila. 21

22 Figure 2 Components Characterized by Q-TOF LC-MS/MS in *Neolamarckia*

23 *cadamba*. The leaf and bark extracts of *N. cadamba* were characterized by Q-TOF

LC-MS/MS (see methods). **a** and **b** LC-MS chromatograms of the four authentic

standards, tryptamine, cadambine, 3a-dihydrocadambine and expoxystrictosidine,

²⁶ respectively. **c** Detection of tryptamine from the leaves and strictosidine from the bark

- in *N. cadamba* by LC-MS. **d** and **e** Expoxystrictosidine, 3a-dihydrocadambine and
- cadambine were detected in the barks of *N. cadamba* by LC-MS. f- j MS/MS
- 29 chromatograms indicated those substances represented by the peak in the previous Fig.
- 30 c-e as tryptamine (**f**, 3.62 min), strictosidine (**g**, 13.02 min), 3α-dihydrocadambine (**h**,

This article is protected by copyright. All rights reserved

1 12.30 min), expoxystrictosidine (i, 13.57min) and cadambine (j, 12.33 min),

2 respectively.

3 Figure 3 Neolamarckia strictosidine synthase gene candidates selection and 4 functional identification of NcSTR1. a Alignment of the active domains in STRs from 5 Ophiorrhiza pumila (Op), Catharanthus roseus (Cr), Rauvolfia serpentine (Rs) and the five 6 Neolamarckia homologous with the highest similarity to CrSTR at amino acid level. Previously reported active-site residues of OpSTR (Q94LW9, Eger, 2020), CrSTR 7 8 (CAA43936.1, Petronikolou, 2018) and RsSTR (P68175.1, Ma 2006) were indicated in 9 red color letters under the aligned sequences. The essential active site glutamate residue (Glu309 in RsSTR1) were highly conserved in NcSTR1, NcSTR13 and NcSTR14. b The 10 expression profile of all the predicted 35 Neolamarckia STRs in 16 tissues. Bark (B), bud, 11 cambium (C), young fruit (FR), old leaves (OL), phloem (P), root (R), young leaves (YL), 12 13 xylem (primary xylem, PX; transitional xylem, TX; secondary xylem, SX), cambium 14 (transitional cambium, TCA; secondary, SCA) and phloem (primary phloem, PPH; transitional phloem, TPH; secondary phloem, SPH) from the first, second and fourth 15 16 internodes. The second internode of 1-year-old seedling was identified as the transition from primary growth to secondary growth. c NcSTR1 catalyzed the Pictet-Spengler reaction of 17 tryptamine with secologanin. In the *in vitro* catalytic reaction, tryptamine and secologanin 18 19 were used as substrates to synthesize strictosidine. The curve in blue and pink indicated the 20 531.10 ->352.20 and 531.10 ->514.25 ion pairs respectively, and black indicated the superimposed curve with blue and pink curve. (i) The standard of strictosidine. (ii)The 21 22 negative control without the enzyme NcSTR1 48 hours after the reaction. (iii) The NcSTR1 23 catalyzed 24 hours after the reaction. (iv) The NcSTR1 catalyzed 48 hours after the reaction. 24 (v)The catalyzed Pictet-Spengler chemical reaction by NcSTR1. 25 Figure 4 The predicted cadambine biosynthetic pathway and MIA biosynthetic gene candidates. a Proposed cadambine biosynthetic pathway in plants. Dotted line arrows indicate 26 multiple steps between intermediates. b The expression profile of all the Neolamarckia 27 28 candidate biosynthetic genes for the seco-iridoid pathway in 16 tissues. Bark (B), bud, 29 cambium (C), young fruit (FR), old leaves (OL), phloem (P), root (R), young leaves (YL), xylem (primary xylem, PX; transitional xylem, TX; secondary xylem, SX), cambium 30

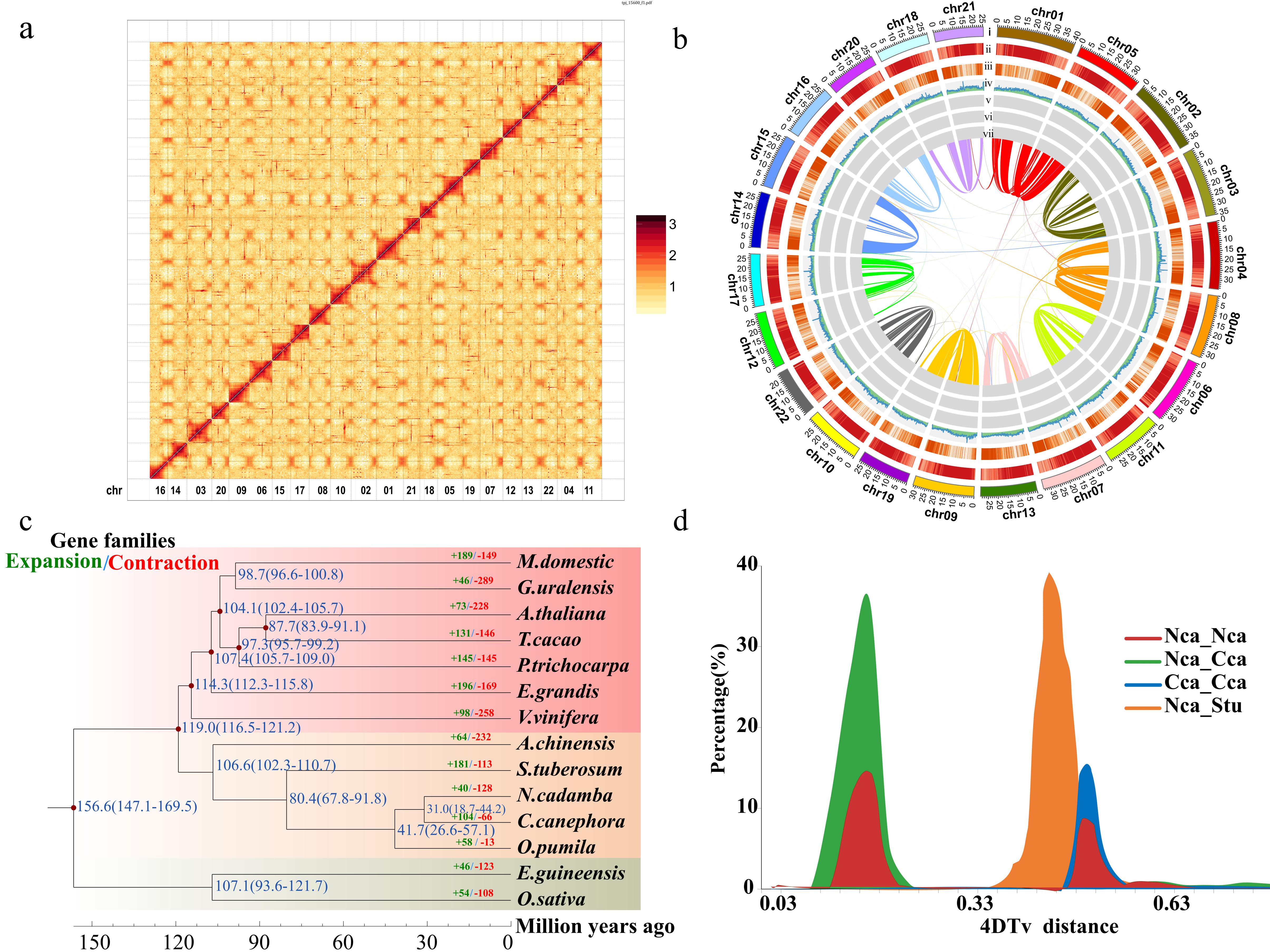
This article is protected by copyright. All rights reserved

(transitional cambium, TCA; secondary, SCA) and phloem (primary phloem, PPH; 1 transitional phloem, TPH; secondary phloem, SPH) from the first, second and fourth 2 3 internodes. Enzymes in abbreviations are: 7-DLH, 7-deoxyloganic acid hydroxylase (CYP72A224); 8-HGO, 8-hydroxy-geraniol oxidoreductase; AS, anthranilic acid synthetase; Ca10OMT, Camptotheca 4 acuminata 10-hydroxycamptothecin O-methyltransferase; CPR, NADPH-Cytochrome P450 reductase; 5 DLGT, 7-deoxyloganetic acid UDP-glucosyltransferase; G10H, Geraniol-10-hydroxylase; GES, 6 7 geraniol synthase; IO, iridoid oxidase (CYP76A26); IS, iridoid synthase; LAMT, loganic acid 8 O-methyltransferase; SGD, strictosidine beta-glucosidase; SLS, Secologanin synthetase; STR, 9 strictosidine synthase; TDC, tryptophan decarboxylase; THAS, tetrahydroalstonine synthase (Kai et al. 2015; Wu et al.2018; Rai, 2021; Qu, 2019; Yang, 2019). 10

Figure 5 Population analysis of Neolamarckia cadamba accessions based on 11 cadambine content. a A neighbor-joining phylogenetic tree of the 31 accessions 12 13 based on their cadambine content. **b** Principal component analysis of the selected 31 accessions. c Decay of LD, measured by r2, in the four groups. d Selective signals in 14 the whole genome between different ecotypes. e Manhattan plot of $\theta\pi$ -based detection 15 16 of selective sweeps identified by comparison between two groups with high and low cadambine content. KEGG enrichment candidate genes associated with cadambine 17 biosynthesis were highlighted with red arrow. 18

Figure 6 The structure and evolution of the predicted cadambine biosynthetic 19 gene clusters and tandem duplications. a Schematic diagram of the predicted gene 20 clusters and tandem duplicates in N. cadamba. NcSMASH35: TDC/STR1/CYPs; NcMCL22: 21 CYP716/NcSTR14/BGLUs; NcMCL18: GES/NcSTR13/CYPs; NcSMASH9, NcSMASH65 22 and NcMCL20: Tandem duplicated NcSTRs; NcMCL7: Tandem duplicated NcSLSs; 23 24 NcMCL9: Tandem duplicated NcTDCs. b Evolution of two genomic regions encompassing gene clusters NcMCL10, NCMCL19, NcMCL23 and NcSMASH67 and the expression 25 26 profile of the putative biosynthetic genes in these clusters. NcMCL10: NcLAMT2/Cytb5/CYP, 27 NcLAMT1/NcDL7H/CYPs, NcMCL19: NcMCL23 and 28 NcSMASH67: Tandem duplicated NcCYP72s. c Evolution of the predicted gene cluster with 29 NcSQE1/IS/ISY/CYP81. d Evolution of tandem duplicated bHLH TFs in NcMCL25 and NcMCL24. Enzymes in abbreviations are: ACT, acyltransferase; ATX, copper transport 30 This article is protected by copyright. All rights reserved

1	protein ATX family; BGLU, β -glucosidase; BHLH25, transcription factor bHLH25; CHX,
2	cation/H(+) antiporter; CPS, terpenoid cyclases; CYP, cytochrome P450s; ERF,
3	ethylene-responsive transcription factor; GES, geraniol synthase; GRF, growth-regulating factor;
4	IS, iridoid synthase; ISY, iridoid synthase paralogue; MATE, multidrug and toxin extrusion
5	protein; MCTS, malignant T-cell-amplified sequence; MFS, major facilitator superfamily
6	protein; NRAMP, metal transporter Nramp family; PRS, disease resistance protein; SEH,
7	soluble epoxide hydrolase; SLS, secologanin synthetase; SQE, squalene epoxidase; SWEET,
8	bidirectional sugar transporter; TDC, tryptophan decarboxylase; TPS, terpenoid synthase.



C

