

---

**Supplementary information**

---

**The *Taxus* genome provides insights into  
paclitaxel biosynthesis**

---

In the format provided by the  
authors and unedited

1 Supplementary information

2  
3 The *Taxus* genome provides insights into  
4 paclitaxel biosynthesis

5  
6 Xingyao Xiong<sup>1,2,9</sup>, Junbo Gou<sup>2,9</sup>, Qinggang Liao<sup>2,9</sup>, Yanlin Li<sup>1,3,9</sup>, Qian Zhou<sup>2,4</sup>, Guiqi Bi<sup>2</sup>,  
7 Chong Li<sup>2</sup>, Ran Du<sup>2</sup>, Xiaotong Wang<sup>2</sup>, Tianshu Sun<sup>2</sup>, Lvjun Guo<sup>5</sup>, Haifei Liang<sup>2</sup>, Pengjun Lu<sup>2</sup>,  
8 Yaoyao Wu<sup>2</sup>, Zhonghua Zhang<sup>6</sup>, Dae-Kyun Ro<sup>2,7</sup>, Yi Shang<sup>8</sup>, Sanwen Huang<sup>2\*</sup>, Jianbin Yan<sup>2\*</sup>

9 <sup>1</sup>College of Horticulture, Hunan Agricultural University, Changsha 410128, China.

10 <sup>2</sup>Shenzhen Branch, Guangdong Laboratory for Lingnan Modern Agriculture, Shenzhen Key  
11 Laboratory of Agricultural Synthetic Biology, Genome Analysis Laboratory of the Ministry of  
12 Agriculture and Rural Affairs, Agricultural Genomics Institute at Shenzhen, Chinese Academy  
13 of Agricultural Sciences, Shenzhen 518124, China.

14 <sup>3</sup>Engineering Research Center for Horticultural Crop Germplasm Creation and New Variety  
15 Breeding, Ministry of Education, Changsha 410128, China.

16 <sup>4</sup>Peng Cheng Laboratory Artificial Intelligence Research Center No.2, Xingke 1st Street,  
17 Nanshan, Shenzhen, 518055, China.

18 <sup>5</sup>MOE Key Laboratory of Bioinformatics, Tsinghua-Peking Joint Center for Life Sciences,  
19 School of Life Sciences, Tsinghua University, Beijing 100084, China.

20 <sup>6</sup>College of Horticulture, Qingdao Agricultural University, Qingdao 266109, China.

21 <sup>7</sup>Department of Biological Sciences, University of Calgary, Calgary, AB, T2N1N4, Canada.

22 <sup>8</sup>The AGISCAAS-YNNU Joint Academy of Potato Sciences, Yunnan Normal University,  
23 Kunming, China.

24 <sup>9</sup>These authors contributed equally to this work.

25 \*Correspondence to: jianbinlab@caas.cn (J.Y.); huangsanwen@caas.cn (S.H.)

26  
27

28  
29  
30

31 **This PDF file includes:**

32

33 1. Supplementary Methods

34 2. Additional Supplementary Figures 1, 2, 3, 4, 5, 6

35 3. Supplementary Tables 1-4, 6-8, 27

36 4. Supplementary References 1-38

## 37 **1. Supplementary Methods**

### 38 **RNA sequencing (RNA-seq) for genome structure annotation and expression analysis**

39 Samples of *T. chinensis* var. *mairei* from one male plant and one female plant collected in Hunan  
40 Province (GPS: N 113° 89' 55", E 28° 26' 32") were used for transcriptome sequencing. Eight  
41 tissues (female strobilus (mfs), leaf (mfl), bark of stem (mfb), and root (mfr) from the female  
42 plant, and male strobili (mtf), leaf (mtl), bark of stem (mtb), and root (mtr) from the male plant)  
43 were collected individually. Two half-sib cell lines (HC, LC) were induced from the embryo of  
44 the female plant mentioned above; the taxoid contents of the lines were different from each  
45 other<sup>1</sup>. The total RNA of these materials was isolated independently using TRIzol<sup>®</sup> reagent  
46 (Invitrogen). The integrity and purity of total RNA were assessed using horizontal agarose  
47 electrophoresis and an Agilent Bioanalyzer 2100 (Agilent Technologies, INC). The NEB Next  
48 Ultra™ RAN Library Prep Kit (NEB, USA) was used to construct the sequencing library. All  
49 RNA-seq libraries were sequenced using the Illumina 2500 platform. The details of the RNA-seq  
50 data of each sample are given in Supplementary Table 25.

### 51 **Genome survey by *k*-mer analysis**

52 Analysis of *k*-mer frequency was adopted to estimate the genome size by using WGS Illumina  
53 reads. The *k*-mer frequency was calculated by Jellyfish (v 2.0.0)<sup>2</sup>, and the genome size was  
54 estimated based on the formula (Total *k*-mer counts)/(Peak depth).

### 55 **Long-read variant calling for genome assessment**

56 The HiFi reads were advantageous for variant calling because they can span repetitive or other  
57 problematic regions. HiFi reads of the *Taxus* leaves (diploid) were aligned to the *Taxus* assembly  
58 genome to assess the quality. Calling the variants were used to check the coverage of the genome,  
59 we aligned HiFi sequencing reads with Minimap2<sup>3</sup> and called SNP (single nucleotide  
60 polymorphism) and InDel with bcftools (v1.9)<sup>4</sup>. Then, the SVs (structure variants) were called  
61 using Sniffles (v1.0.12) (<https://github.com/petabi/sniffles>).

### 62 **Analysis of the insertion time of LTR elements**

63 The candidate LTR-RTs were identified using LTR\_Finder (v1.02)<sup>5</sup> (parameters: -D 2000 -d  
64 1000 -L 3500 -l 100 -p 20 -C -M 0.8) and LTRharvest<sup>6</sup> (parameters: -similar 90 -vic 10 -seed 20  
65 -seqids yes -minlenltr 100 -maxlenltr 7000 -mintsd 4 -maxtsd 6 -motif TGCA -motifmis 1). Only  
66 full-length LTR-RTs were retained. In addition, false positives were discarded by the  
67 LTR\_retriver pipeline<sup>7</sup>. After removing non-LTR repeat elements and tandem repeats, the full-  
68 length LTR retrotransposons with optimal frames were translated into amino acids. Their  
69 functional domains were predicted using HMMER<sup>8</sup> based on the Pfam database. Paralogs of the  
70 RT domain, specific for the Copia and Gypsy superfamilies, were then detected by the functional  
71 domain orders.

72 The RT protein sequences were aligned using MUSCLE (v3.8.31)<sup>9</sup>, and maximum likelihood  
73 trees were built using FastTree<sup>10</sup> for the Copia and Gypsy superfamilies. The two ends of these  
74 LTR retrotransposons were aligned using MUSCLE<sup>5,6</sup>, and the nucleotide distance (D) was  
75 estimated using the Kimura two-parameter (K2p) criterion<sup>11</sup>, which was implemented in the  
76 distmat program of the EMBOSS package (v6.6.0)<sup>12</sup>. The rates of nucleotide substitution ( $\mu$ )  
77 were obtained from a previous study<sup>13</sup>. The insertion times of an LTR retrotransposon were  
78 calculated using  $T = D/2\mu$ .

## 79 **Gene expression level analysis**

80 RNA-seq reads were aligned to the Taxus genome using HISAT2 (v2.1.0)<sup>14</sup> with default  
81 parameters. The sequence alignment files generated by HISAT2 were subsequently inputted to  
82 StringTie (v1.3.5)<sup>15</sup> to generate counts of uniquely mapped reads of annotated genes from the  
83 reference annotation file. The genes that differed significantly between two tissues were  
84 identified with count-based methods in R packages edgeR (v3.18.1)<sup>16</sup>. Benjamini-Hochberg  
85 correction was used to correct for multiple comparisons (with a false discovery cut-off < 0.05).  
86 The read counts of genes were transformed to FPKM values using fragments mapping per  
87 kilobase of exon per million fragments mapped. The FPKM values were used to represent the  
88 expression level of genes. Whereas, the genes with high sequence similarity were difficult to  
89 separate the reads mapping on these genes based on the methods above. To detect the high  
90 similarity sequence genes, we adjusted the threshold (the min acceptable alignment score  
91 parameter --score-min set as L,0, 0 in HISAT2) to no mismatch and get the alignment files. The

92 expression level of genes was calculated using StringTie (v1.3.5)<sup>15</sup> using the alignment file as  
93 input files. Then, the relative low nucleotide differences distributed on the exons could use to  
94 distinguish the different expression level

## 95 **Ortholog analysis**

96 Publicly available genome and annotated peptide sequences of *A. thaliana* (TAIR 10), *O. sativa*  
97 (IRGSP-1.0), *A. trichopoda* (AMTR1.0), and *S. moellendorffii* (v1.0) were downloaded from  
98 EnsemblPlants (<http://plants.ensembl.org/>), while the *P. abies* (v1.0b) and *G. biloba* (v1.0)  
99 sequences were downloaded from TreeGenes<sup>17</sup>. An all-to-all blast was performed using blastall<sup>18</sup>  
100 and the data were then processed to find orthologous gene groups using OrthoMCL<sup>19</sup>.  
101 Subsequently, 193 single-copy gene families (out of a total of 35,298 gene families) were found,  
102 and their protein sequences were aligned with MUSCLE<sup>20</sup>. The aligned sequences were then  
103 integrated into a single concatenated sequence for each species, and the conserved blocks were  
104 identified for further phylogenetic analysis using Gblocks<sup>21</sup>. The maximum likelihood  
105 phylogenetic tree was constructed using RAxML (Random Accelerated Maximum Likelihood)<sup>22</sup>  
106 with the GAMMAI+LG model, and the absolute rates of divergence times were estimated using  
107 r8s<sup>23</sup>. The number of expanded or contracted gene families along each branch was computed  
108 using CAFÉ<sup>24</sup> with a *p*-value cutoff of 0.05. Gene Ontology enrichment analysis was performed  
109 and visualized using clusterProfiler<sup>25</sup>.

## 110 **Expression and purification of TSs in *E. coli***

111 To obtain recombinant TSs-His6, pET28b::*TS1* and pET28b::*TS2* were transformed into *E. coli*  
112 BL21 (DE3) and streaked on LB plates with kanamycin (50 µg mL<sup>-1</sup>). A single positive colony  
113 was cultured at 37°C overnight in 5 mL of LB liquid medium with the same kanamycin  
114 concentration. The overnight cultures were used for expanding the culture in 500 mL of fresh  
115 medium and grown until the OD<sub>600</sub> reached 0.4, at which point, 1 mM isopropyl β-D-1-  
116 thiogalactopyranoside (IPTG) was added to induce expression. Expression induction was  
117 performed in a 180 rpm shaker at 16°C overnight, and harvested cells were resuspended in lysis  
118 buffer (50 mM sodium phosphate, pH 8.0 and 300 mM NaCl and 10 mM imidazole). After cell  
119 lysis by sonication and centrifugation at 4,000 g for 20 min at 4°C, the total protein supernatant

120 was loaded onto a 1.5 mL HisPur Ni-NTA resin-packed column (Thermo Scientific, U.S.A), and  
121 the recombinant TSs-His6 were finally eluted with elution buffer (50 mM sodium phosphate, pH  
122 8.0, 300 mM NaCl and 250 mM imidazole). The purified recombinant TSs-His6 were desalted  
123 into enzyme assay buffer (25 mM HEPES, pH 8.5, 10% glycerol, 5 mM DTT, 5 mM sodium  
124 ascorbate, 5 mM sodium metabisulfite, and 1 mM MgCl<sub>2</sub>) and concentrated by centrifugation  
125 with a 30-kDa cutoff Millipore concentrator. The content and purity of the recombinant TSs-  
126 His6 were determined by a BCA protein assay kit (Beyotime, China).

### 127 ***In vitro* characterization of T5aH1 and its homologous genes**

128 The ORFs of *Taxus T5aH1* (*ctg5306\_gene.3*), *T5aH2* (*ctg7747\_gene.2*), *T5aH3*  
129 (*ctg2768\_gene.2*), *55326109* (*ctg7747\_gene.3*), and *55305455* (*ctg7747\_gene.4*) were cloned  
130 from the *Taxus* cell line cDNA and inserted into the yeast expression vector pESC-His. Thus, the  
131 recombination plasmid pESC-His-*CYP725As-Flag* was obtained after positive clone screening.  
132 pESC-His-*CYP725As-Flag* was expressed in the yeast WAT11 strain<sup>26</sup>. Their expression was  
133 confirmed by western blot assay. All expression plasmids were constructed using the Hieff  
134 Clone™ One Step Cloning Kit (YEASEN, China) and primers used in this work are given in  
135 Supplementary Table 28.

136 For the *in vitro* enzyme assay, microsomes of the yeast strain WAT11 expressing pESC-His-  
137 *CYP725As-Flag* were prepared as described previously<sup>27</sup>. As a control, microsomes of the  
138 WAT11 strain harboring the empty vector pESC-His were prepared. The activity of the pESC-  
139 His-*CYP725As-Flag* protein was tested in a 1 mL mixture containing 50 mM HEPES buffer (pH  
140 7.5), 100 μM taxadiene substrate, 500 μM NADPH, and 2 mg of pESC-His-*CYP725As-Flag*-  
141 containing microsomal proteins and overlaid with 500 μL of pentane. As a control, microsomes  
142 harboring the empty vector pESC-His were prepared. After incubating the mixture at 32°C for 2  
143 h, the pentane layer was separated for GC/MS analysis.

144 The previous peroxide and stability assays indicated that OCT could be oxidized from taxa-  
145 4(5),11(12)-diene by hydrogen peroxide, and OCT can spontaneously form iso-OCT at 37°C  
146 overnight<sup>28</sup>. For given OCT and iso-OCT, 100 μL of the above purified taxadienes  
147 (approximately 200 ug L<sup>-1</sup>) were concentrated by N<sub>2</sub> gas in a glass GC vial, then the residue was

148 resuspended in 500  $\mu\text{L}$  of hydrogen peroxide and shaken at 37°C and 220 rpm overnight. After  
149 reacting overnight, the taxanes were extracted with hexane and analyzed with GC-MS.

#### 150 **Transcription analysis under methyl jasmonate (MeJA) treatment**

151 The *T. chinensis* var. *mairei* LC cell lines were divided into two groups: the test group was  
152 treated with 100  $\mu\text{M}$  MeJA solution, whereas the control group was treated with solvent  
153 (ethanol). Samples were harvested after 0, 2, 4, 8, and 24 h of elicitation for gene expression  
154 analysis and taxane measurements. Experiments were carried out in triplicate. Total RNA was  
155 isolated from the MeJA-treated samples using the EASYspin Plant RNA isolation kit (Aidlab,  
156 Beijing, China), and 1  $\mu\text{g}$  of cDNA was synthesized using Hifair<sup>®</sup> III 1st Strand cDNA Synthesis  
157 SuperMix for qPCR (gDNA digester plus) (YEASEN, China). Transcript abundance was  
158 measured using a QuantStudio<sup>™</sup> 3 System with Hieff<sup>®</sup> qPCR SYBR<sup>®</sup> Green Master Mix (Low  
159 Rox Plus) (YEASEN, China). All real-time PCRs were repeated with at least two technical and  
160 three biological replicates. Mean cycle threshold (*Ct*) values were normalized using *T. chinensis*  
161 var. *mairei* actin 1 gene (*7G702435613*) as a validated reference gene. The relative gene  
162 expression value was calculated using the  $2^{-\Delta\Delta C_t}$  method. Gene-specific primers are listed in  
163 Supplementary Table 28.

#### 164 **Metabolite extraction and chromatography-mass spectrometry analysis**

165 The cell samples were extracted with a modified Wolfender method. One hundred milligrams of  
166 freeze-dried cell powder was weighed and transferred into a 2-mL centrifuge tube with 10  $\mu\text{L}$   
167 internal standards (500 ng  $\text{mL}^{-1}$  dexamethasone). Then, the samples were resuspended in 1.5 mL  
168 extraction buffer (methanol:water = 80:20, v/v), and a 30 min ultrasonic-assisted extraction  
169 process was applied. After 15 min of centrifugation at 14,000 g, the supernatants were dried in a  
170 LABCONCO CentriVap vacuum centrifugal concentrator and resuspended in 200  $\mu\text{L}$  methanol  
171 solvent (80:20, v/v).

172 GC/MS analysis was performed on an Agilent 7890B GC machine (Agilent Technologies,  
173 Waldbronn, USA) equipped with an Agilent 7000C mass selective detector at 70 eV and 1.2 mL  
174  $\text{min}^{-1}$  helium flow. One to five microliters of the sample was injected and analyzed on an Agilent  
175 HP-5MS column (5% phenyl methyl silox, 30 m  $\times$  250  $\mu\text{m}$  internal diameter, 0.25- $\mu\text{m}$  film

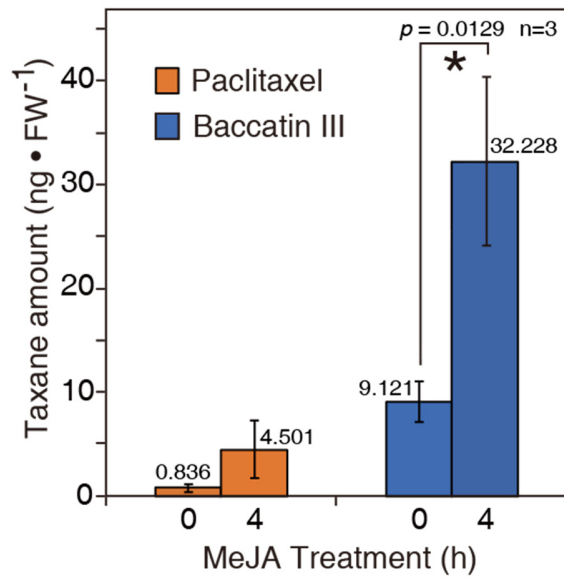


176 thickness). The oven temperature program was as follows: 45°C for 1 min, then a 10°C min<sup>-1</sup>  
177 ramp to 250°C, and a hold at 250°C for 5 min. The injection temperature was 250°C. Full mass  
178 spectra were generated for metabolite identification by scanning the mass-to-charge ratio (*m/z*)  
179 range from 40 to 350.

180 Liquid chromatography-mass spectrometry (LC-MS) analysis was carried out on an ACQUITY  
181 UPLC I-Class (Waters) coupled to a 4500 QTRAP triple quadrupole mass spectrometer (AB  
182 SCIEX) equipped with a 50 × 2.1 mm, 1.8 μm ACQUITY UPLC™ Waters Xselect HSS T3  
183 column (Waters). Ten microliters of supernatant was loaded each time and then eluted at a flow  
184 rate of 200 μL min<sup>-1</sup> with initial conditions of 40% mobile phase A (0.1% formic acid in  
185 methanol) and 60% mobile phase B (0.1% formic acid in water) followed by a 10-min linear  
186 gradient to 100% mobile phase B. The autosampler was set at 10°C. Mass spectrometry was  
187 operated separately in positive electrospray ionization mode. The analyte [M+H] was selected as  
188 the precursor ion. The quantitation mode was multiple reaction monitoring (MRM) mode using  
189 mass transitions (precursor ions/product ions). The temperature of the ESI ion source was set at  
190 500°C. The curtain gas flow was set at 20 psi, collisional activated dissociation (CAD) gas was  
191 set as the medium, and the ion spray voltage was (+) 5,500 V, with ion gases 1 and 2 set as 50  
192 psi. AB SCIEX Analyst 1.6.3 Software (Applied Biosystems) was used for data acquisition and  
193 processing.

194 2. Additional Supplementary Figures

195



196

197 **Additional Supplementary Fig.1 The accumulation of paclitaxel and baccatin III in Taxus**  
198 **cell line induced by 100 μM MeJA at 0 hours and 4 hours.** The amounts of paclitaxel and  
199 baccatin III were measured by LC-MS analysis. FW indicates fresh weight (100 mg); error bars  
200 display the standard error (n = 3 biological replicates); asterisks show significant differences  
201 with \* $p \leq 0.05$  by two-tailed Student's t-test.

202

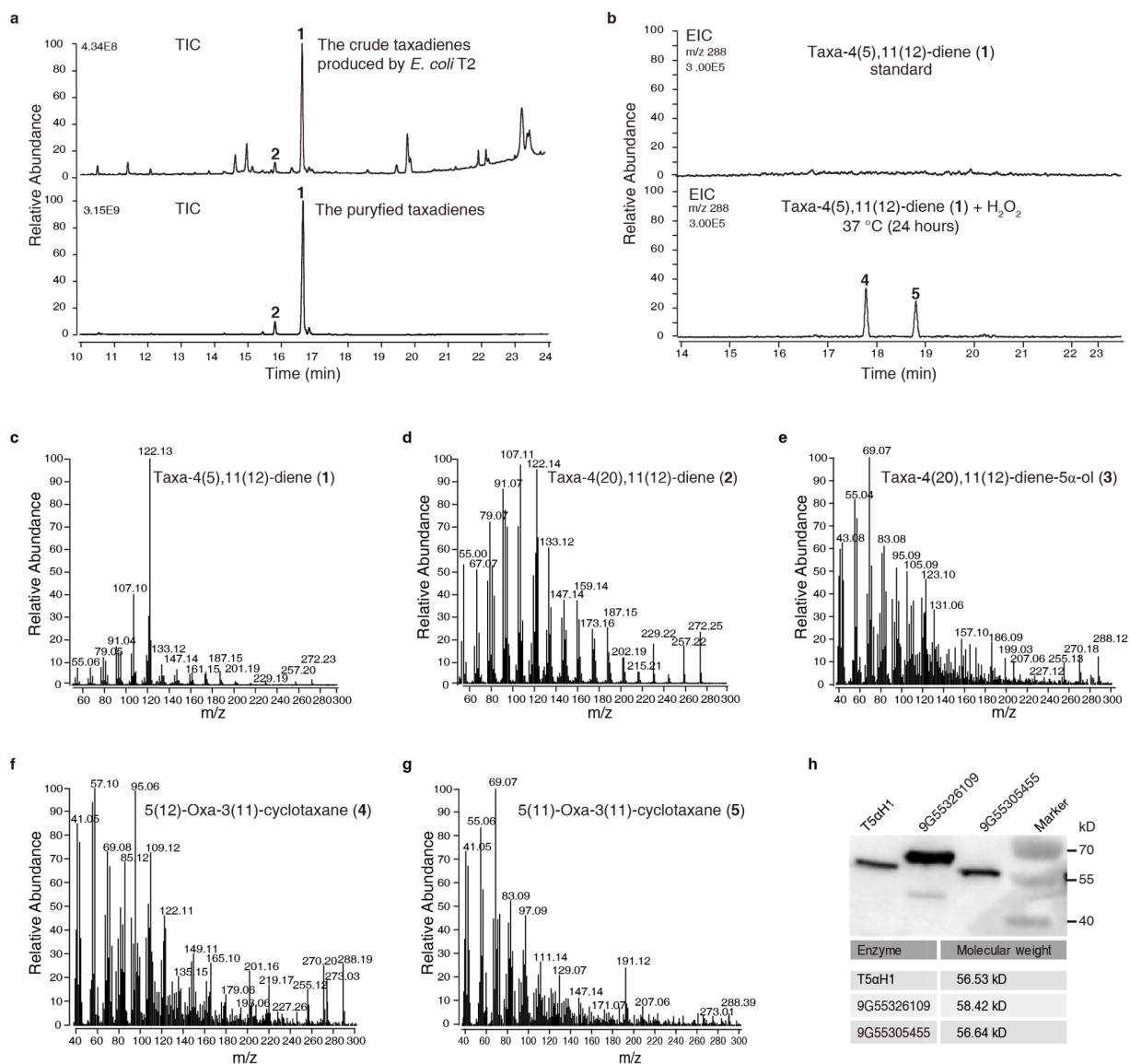


207  
208  
209  
210

T5 α H1	ATGGAAGCCCGTGTAAAGAGGACAGTTGCAAAATTTAAAGAGTCAACAGCTGGAAGCTGTCACAGTGAATTTTTTTCATTGCCCTCTCACTATTGCTGCTATTCTTCTGCTCTCTCTG	120
T5 α H2	ATGGAAGCCCGTGTAAAGAGGACAGTTGCAAAATTTAAAGAGTCAACAGCTGGAAGCTGTCACAGTGAATTTTTTTCATTGCCCTCTCACTATTGCTGCTATTCTTCTGCTCTCTCTG	120
T5 α H3	ATGGAAGCCCGTGTAAAGAGGACAGTTGCAAAATTTAAAGAGTCAACAGCTGGAAGCTGTCACAGTGAATTTTTTTCATTGCCCTCTCACTATTGCTGCTATTCTTCTGCTCTCTCTG	120
Consensus	at ggaagcccg t g t a a a g a g c a g t t g c a a a a t t t a a g a g t c a c a c a g t g g a c t g t t c c a c t g a a t c t t t t c c a t t g c c c t c t c a t c t a t g c t g g t a t t c t c t g c t t c t c t g	
T5 α H1	CTCTTCGGTCTAAAGCCGACTCCTCCCTAAACTCTCTCTGGAAATTAGGATCCCTTTTCATTGGCGAGTGGTTATATCTGAGGCTCTCGATCGAAGCTCGCTGGAGCAATTT	240
T5 α H2	CTCTTCGGTCTAAAGCCGACTCCTCCCTAAACTCTCTCTGGAAATTAGGATCCCTTTTCATTGGCGAGTGGTTATATCTGAGGCTCTCGATCGAAGCTCGCTGGAGCAATTT	240
T5 α H3	CTCTTCGGTCTAAAGCCGACTCCTCCCTAAACTCTCTCTGGAAATTAGGATCCCTTTTCATTGGCGAGTGGTTATATCTGAGGCTCTCGATCGAAGCTCGCTGGAGCAATTT	240
Consensus	c t c t c g t t c t a a a g c c a c t c c t c c c t t a a c t t c c t c c t g g g a a a t t a g g a t c c c t t t c a t t g g c g a g t c g t t a t a t t c c t g a g g g c t c t t c g a t c g a a c t c g c t g g a g c a a t t t	
T5 α H1	TTTGACGAGAGAGTGAAGCAATTCGGCTCTGTGTCAAGACCTCCTTAATGGGCATCCGACAGTACTCTGCGGCCCGGGGAAACCGGCTATTCTGTCCAACGAGGAGAAAGCTG	360
T5 α H2	TTTGACGAGAGAGTGAAGCAATTCGGCTCTGTGTCAAGACCTCCTTAATGGGCATCCGACAGTACTCTGCGGCCCGGGGAAACCGGCTATTCTGTCCAACGAGGAGAAAGCTG	360
T5 α H3	TTTGACGAGAGAGTGAAGCAATTCGGCTCTGTGTCAAGACCTCCTTAATGGGCATCCGACAGTACTCTGCGGCCCGGGGAAACCGGCTATTCTGTCCAACGAGGAGAAAGCTG	360
Consensus	t t t g a c g a g a g a g t g a a g c a a t t c g g t c t g t t c a a g a c c t c c t t a a t t g g g c a t c c c a c a g t a g t a c t c t g c g g c c t g o g g g a a c c g g c t t a t t c t g t c c a a c g a g g a a g c t g	
T5 α H1	GTGCAGATGTGCTGGCCGCTCAATTTATGAAGCTCATGGGAGAGAAATTCGTTGCCACGAGGAGGGTGAAGACCATATAGTTATGGCTCTGCTCTTGCAGGTTTTTTCGGCCCTGGT	480
T5 α H2	GTGCAGATGTGCTGGCCGCTCAATTTATGAAGCTCATGGGAGAGAAATTCGTTGCCACGAGGAGGGTGAAGACCATATAGTTATGGCTCTGCTCTTGCAGGTTTTTTCGGCCCTGGT	480
T5 α H3	GTGCAGATGTGCTGGCCGCTCAATTTATGAAGCTCATGGGAGAGAAATTCGTTGCCACGAGGAGGGTGAAGACCATATAGTTATGGCTCTGCTCTTGCAGGTTTTTTCGGCCCTGGT	480
Consensus	g t g c a g a t g t c g t g g c c c g c t c a a t t t a t g a a g c t c a t g g a g a g a a t t c c g t t g c c a c c a g g a g g g t g a a g a c c a t a t a g t t a t g c g c t c t g c t c t g c a g g t t t t t c g g c c t g g t	
T5 α H1	CGCGTGCAGAGTACATTTGGTAAAATGAATACAGAGATCCAGATTCATATCAACGAGCAATGGAAGGGAAGAAGTGAAGTGAATGTACTTCCTTTGTTAAGAGAGCTCGTCTCAACATC	600
T5 α H2	CGCGTGCAGAGTACATTTGGTAAAATGAATACAGAGATCCAGATTCATATCAACGAGCAATGGAAGGGAAGAAGTGAAGTGAATGTACTTCCTTTGTTAAGAGAGCTCGTCTCAACATC	600
T5 α H3	CGCGTGCAGAGTACATTTGGTAAAATGAATACAGAGATCCAGATTCATATCAACGAGCAATGGAAGGGAAGAAGTGAAGTGAATGTACTTCCTTTGTTAAGAGAGCTCGTCTCAACATC	600
Consensus	g c g c t g c a g a g t a c a t t t g g t a a a t g a a t a c a g a g a t c c a g a t t c a t a t c a a c g a g c a a t t g a a g g g a a a a g a g t g a a g t g a a t g t a c t t c c t t t g t t a a g a g a c t c g t c t c a a c a t c	
T5 α H1	TCGGCCATCTTTTTCACACATATATGATAAAGCAGAAACAGGATCGTCTGCATAAAGCTTTTGAAACTATTCTGCTCGAAGTTTTGCTCTTCCGATTGACTTGCOCGGATTTGGTTTTG	720
T5 α H2	TCGGCCATCTTTTTCACACATATATGATAAAGCAGAAACAGGATCGTCTGCATAAAGCTTTTGAAACTATTCTGCTCGAAGTTTTGCTCTTCCGATTGACTTGCOCGGATTTGGTTTTG	720
T5 α H3	TCGGCCATCTTTTTCACACATATATGATAAAGCAGAAACAGGATCGTCTGCATAAAGCTTTTGAAACTATTCTGCTCGAAGTTTTGCTCTTCCGATTGACTTGCOCGGATTTGGTTTTG	720
Consensus	t c g g c c a t c t t t t t c a c a c a t a t a t g a t a a g c a g a a c a g g a t c g t c t g c a t a a g c t t t g g a a c t a t t c t g t c g g a a g t t t t g c t c t c c g a t t g a c t t g c c c g g a t t g g t t t c	
T5 α H1	CATAGAGCACTCCAGGGAACGGGCACTCAACAAAATTTATGGTGTCTTTAAATAAAAGAGAAAAGAGAGTGTGCACTGTGATCGGCAACAGGCACTCAGGATCTGCTCTCTGTTTTG	840
T5 α H2	CATAGAGCACTCCAGGGAACGGGCACTCAACAAAATTTATGGTGTCTTTAAATAAAAGAGAAAAGAGAGTGTGCACTGTGATCGGCAACAGGCACTCAGGATCTGCTCTCTGTTTTG	840
T5 α H3	CATAGAGCACTCCAGGGAACGGGCACTCAACAAAATTTATGGTGTCTTTAAATAAAAGAGAAAAGAGAGTGTGCACTGTGATCGGCAACAGGCACTCAGGATCTGCTCTCTGTTTTG	840
Consensus	c a t a g a g c a c t c c a g g g a c g g g c a a g c t c a a c a a a t t a t g c t g t c t t t a a t t a a a a g a g a a a a a g a g a t c t g a c t c t g g a t c g g c a c a g c a c t c a g g a t c t g c t c t g t t t t g	
T5 α H1	CTCACTTTCAAGAGTGAACAAGGGAACCTCACTCACCATGACGAGATCTGCACAACCTTTCTCTCTGCTCCATGCCCTCCTATGACACCACCACCTTCGCCAATGGCTTTGATTTTCAAG	960
T5 α H2	CTCACTTTCAAGAGTGAACAAGGGAACCTCACTCACCATGACGAGATCTGCACAACCTTTCTCTCTGCTCCATGCCCTCCTATGACACCACCACCTTCGCCAATGGCTTTGATTTTCAAG	960
T5 α H3	CTCACTTTCAAGAGTGAACAAGGGAACCTCACTCACCATGACGAGATCTGCACAACCTTTCTCTCTGCTCCATGCCCTCCTATGACACCACCACCTTCGCCAATGGCTTTGATTTTCAAG	960
Consensus	c t c a c t t t c a g a g a t g a c a a a g g g a c t c c a c t c a c c a a t g a c g a g a t a c t g c a c a a c t t t c t c t c t g c t c c a t g c c t a t g a c a c c a c c a c t t c g c c a a t g g c t t g a t t t c a a g	
T5 α H1	CTCTTGTCTTCCAATCCAGATGCTATCAAAAAGTAGTTCAGAGCAATTCGAGATAGTTTCGAAACAAGAGGAGGCGAAGAAATCACATGGAAGGATCTCAAGCCATGAAATACACA	1080
T5 α H2	CTCTTGTCTTCCAATCCAGATGCTATCAAAAAGTAGTTCAGAGCAATTCGAGATAGTTTCGAAACAAGAGGAGGCGAAGAAATCACATGGAAGGATCTCAAGCCATGAAATACACA	1080
T5 α H3	CTCTTGTCTTCCAATCCAGATGCTATCAAAAAGTAGTTCAGAGCAATTCGAGATAGTTTCGAAACAAGAGGAGGCGAAGAAATCACATGGAAGGATCTCAAGCCATGAAATACACA	1080
Consensus	c t c t t g t c t t c c a a t c c a a a t g c t a t c a a a a g t a g t t c a a g a c a a t t g g a g a t a c t t c c a a c a a g a g g a g g g c a a g a a a t c a c a t g g a a g a t c t c a a g c c a t g a a a t a c a c a	
T5 α H1	TGGCAAGTAGCTCAGGAAACGCTGGCAGTGTTCCTCCAGTTTTGGAACATTCGCAAGGCCACTCACTGACATTCAGTATGATGGTTACACAATTCGAAAGGGTGGAAAGCTGTGTGG	1200
T5 α H2	TGGCAAGTAGCTCAGGAAACGCTGGCAGTGTTCCTCCAGTTTTGGAACATTCGCAAGGCCACTCACTGACATTCAGTATGATGGTTACACAATTCGAAAGGGTGGAAAGCTGTGTGG	1200
T5 α H3	TGGCAAGTAGCTCAGGAAACGCTGGCAGTGTTCCTCCAGTTTTGGAACATTCGCAAGGCCACTCACTGACATTCAGTATGATGGTTACACAATTCGAAAGGGTGGAAAGCTGTGTGG	1200
Consensus	t g g c a a g t a g c t c a g g a a a c g c t g c g g a t g t t t c c t c c a g t t t t g g a a c a t t c g c a a g g c a c t c a c t g a c a t t c a g t a t g a t g g t t a c a c a a t t c g a a a g g g t g g a a g c t g t t g t g g	
T5 α H1	ACAACCTACAGTACACATCCGAAGGACTTGTATTTCAATGAACCAAGAAAATTCATGCCCTCAAGATTCGATCAGGAAGGAAGCATGTAGCTCCTACACATTTTTGCCCTTCGGTGGG	1320
T5 α H2	ACAACCTACAGTACACATCCGAAGGACTTGTATTTCAATGAACCAAGAAAATTCATGCCCTCAAGATTCGATCAGGAAGGAAGCATGTAGCTCCTACACATTTTTGCCCTTCGGTGGG	1320
T5 α H3	ACAACCTACAGTACACATCCGAAGGACTTGTATTTCAATGAACCAAGAAAATTCATGCCCTCAAGATTCGATCAGGAAGGAAGCATGTAGCTCCTACACATTTTTGCCCTTCGGTGGG	1320
Consensus	a c a a c t t a c a g t a c a c a t c c g a a g g a c t t g t a t t t c a a t g a a c c a g a g a a a t t c a t g c c t t c a a g a t t c g a t c a g g a a g g a a g c a t g t a g c t c c t t a c a c a t t t t t t g c c c t t c g g t g g a	
T5 α H1	GGCAGCGGTGATGTGGGATGGGACTTTCAAAGATGGAGATTTTACTGTTGCTTCATCATTTTGTCAAAACTTTTAGCAGCTACACCCAGTTGATCCCGACGAAAATAATCAGGG	1440
T5 α H2	GGCAGCGGTGATGTGGGATGGGACTTTCAAAGATGGAGATTTTACTGTTGCTTCATCATTTTGTCAAAACTTTTAGCAGCTACACCCAGTTGATCCCGACGAAAATAATCAGGG	1440
T5 α H3	GGCAGCGGTGATGTGGGATGGGACTTTCAAAGATGGAGATTTTACTGTTGCTTCATCATTTTGTCAAAACTTTTAGCAGCTACACCCAGTTGATCCCGACGAAAATAATCAGGG	1440
Consensus	g g c a g c g g t c a t g t g g g a t g g g a c t t t t c a a a g a t g g a g a t t t t a c t g t t c g t t c a t c a t t t t g t t a a a a c t t t t a g c a g c t a c a c c c a g t t g a t c c c g a c g a a a a a a t a t c a g g g	
T5 α H1	GATCCACTCCCTCCTCTTCTCTCAAAGGATTTTCCATTAAACTGTTTCCGAGACATA	1499
T5 α H2	GATCCACTCCCTCCTCTTCTCTCAAAGGATTTTCCATTAAACTGTTTCCGAGACATA	1499
T5 α H3	GATCCACTCCCTCCTCTTCTCTCAAAGGATTTTCCATTAAACTGTTTCCGAGACATA	1499
Consensus	g a t c c a c t c c c t c c t c t c c t t c c a a g g a t t t t c c a t t a a a c t g t t c c g a g a c a t a	

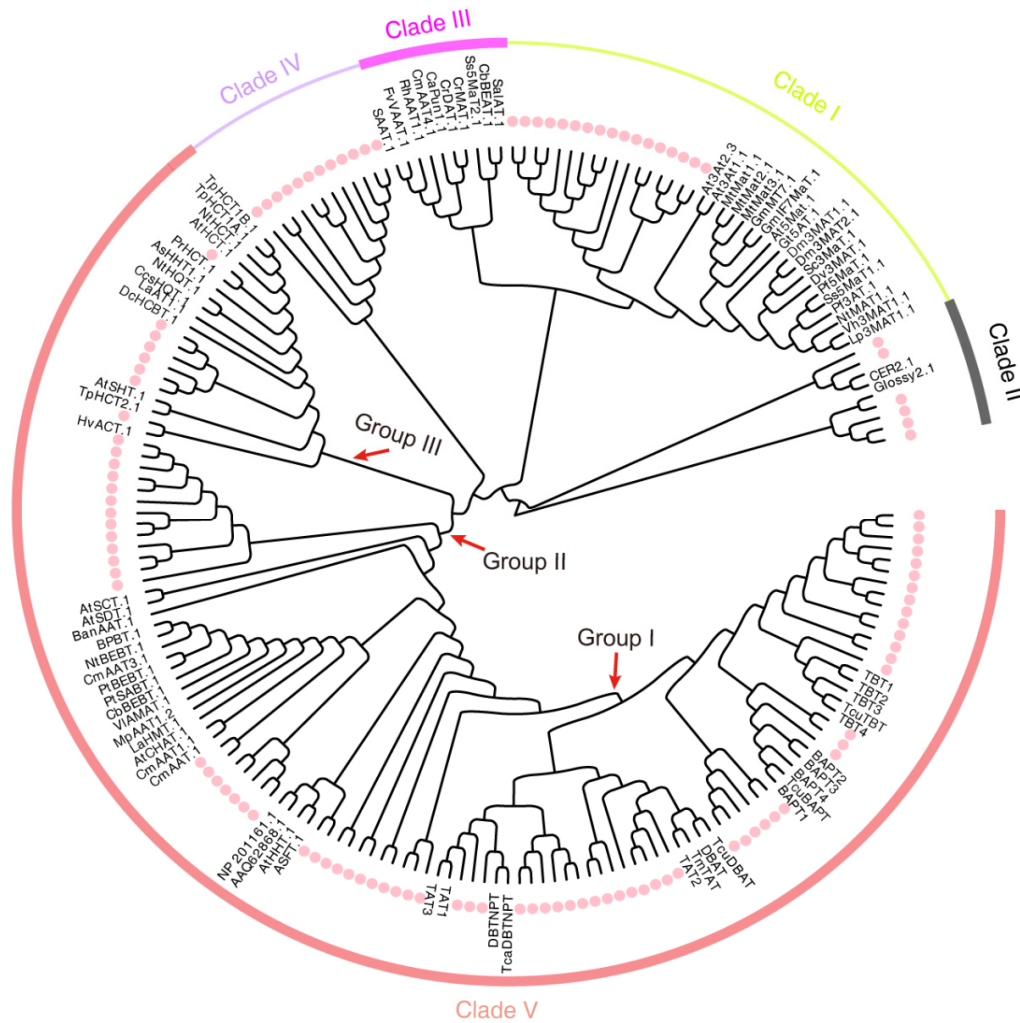
211  
212  
213  
214  
215  
216  
217  
218  
219

**Additional Supplementary Fig.3 Nucleotide sequence alignment of T5αHs. T5αH1, 2, and 3 refer to *T. chinensis* var. *mairei* taxa-4(5),11(12)-diene-5α-hydroxylase 1, 2, and 3, respectively.**

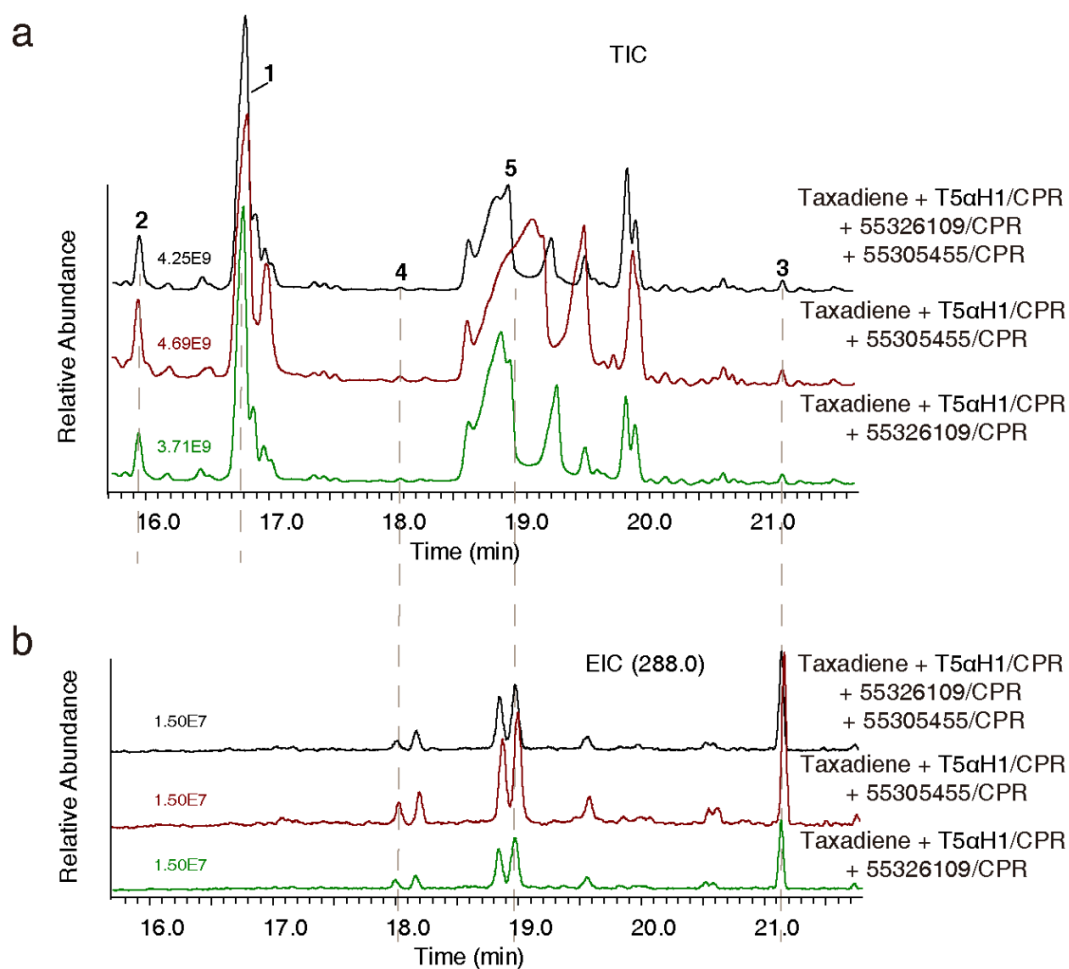


220  
 221 **Additional Supplementary Fig.4 GC-MS analysis of taxanes (peaks 1, 2, 3, 4, and 5) and**  
 222 **protein expression of the three CYP725As (T5 $\alpha$ H1, 9G55326109, and 9G55305455) in yeast**  
 223 **to complement Fig. 3. a**, GC-MS analysis of the crude taxadiene produced by *E. coli* T2<sup>25</sup> and  
 224 the purified taxadiene by thin-layer chromatography. Peaks **1** and **2** are taxa-4(5),11(12)-diene  
 225 and taxa-4(20),11(12)-diene, respectively, produced by TSs; **b**, GC-MS analysis of preparative **4**  
 226 and **5**. OCT (**4**) was oxidized from **1** by hydrogen peroxide, and **4** spontaneously formed iso-  
 227 OCT (**5**) at 37°C overnight. Peaks **3**, **4** and **5** represent the compounds taxa-4(20),11(12)-diene-  
 228 5 $\alpha$ -ol, OCT and iso-OCT, respectively, generated by T5 $\alpha$ Hs. Chromatograms at the extracted  
 229 ions of  $m/z^+$  288. The MS/MS fragmentation of peaks **1** (c), **2** (d), **3** (e), **4** (f) and **5** (g). **h**,

230 Immunoblotting assays of three enzymes (T5 $\alpha$ H1, 9G55326109, and 9G55305455) in the  
 231 WAT11 strain. The mouse monoclonal antibody for the FLAG tag and sheep polyclonal antibody  
 232 conjugated with horseradish peroxidase (Beyotime, China) were used as primary and secondary  
 233 antibodies, respectively. The immunoblotting assays were repeated 3 times.  
 234



235  
 236 **Additional Supplementary Fig.5 Phylogenetic analysis of Taxus BAHD acyltransferase**  
 237 **proteins.** The tree is generated using the neighbor-joining method and 1000 bootstraps based on  
 238 putative amino acid full-length transferase sequences with MEGA 7.0 software and EvolView.  
 239 The pink-red circle denotes transferase genes in *T. chinensis* var. *mairiei*. The subgroups were  
 240 designated based on a previous report<sup>30</sup>. All defined paclitaxel pathway BAHD acyltransferase  
 241 genes are included in Group 1 of Clade V.  
 242



243

244

245 **Additional Supplementary Fig.6 Enzymatic analysis of CYP725A 55326109 and 55305455**

246 **with taxane as the substrate. a**, Total ion chromatograms (TICs) of the *in vitro* assays. Peaks 1-

247 5 are the substrates taxa-4(5),11(12)-diene, taxa-4(20),11(12)-diene, taxa-4(20),11(12)-dien-5α-

248 ol, OCT and iso-OCT, respectively. The taxadiene substrates (peaks **1** and **2**) were exogenously

249 added, while the oxidized taxadienes (peaks **3**, **4**, and **5**) were from the T5αH/CPR catalytic

250 products. **b**, Chromatograms of the extracted ions at 288  $m/z^+$  for the oxidized taxadienes (peaks

251 **3**, **4**, and **5**). There were no new peaks produced by the reaction mixtures, including T5αH/CPR

252 and 55326109/CPR, T5αH/CPR and 55305455/CPR, and T5αH/CPR, and 55326109/CPR and

253 54359314/CPR.

254

255 **3. Supplementary Tables**

256

257 **Supplementary Table 1 The Taxus genome assembly.**

258 **A. Sequencing data used to assemble the Taxus genome.**

<b>Material</b>	<b>Data Type</b>	<b>Size (Gb)</b>
Endosperm	Illumina (WGS)	693.73
Endosperm	Pacbio (CLR)	318.05
Leaf	Pacbio (HiFi)	26.06
Endosperm	Hi-C	1135.15
Total		2172.99

259 WGS, Whole-Genome Sequencing; CLR, Continuous Long Read; HiFi, High fidelity reads.

260

261 **B. The statistics of the PacBio (CLR) assembly.**

	<b>Contig</b>	<b>Scaffold</b>
Total	10,232,190,857 bp	10,237,952,243 bp
Count	15,880	3,813
Average	644,344.51 bp	2,685,012.39 bp
Max	21,889,323 bp	1,051,374,208 bp
N25	4,384,376 bp	957,407,293 bp
N50	2,436,244 bp	903,737,476 bp
N75	1,124,632 bp	769,507,007 bp
N90	426,395 bp	659,955,358 bp
N95	216,560 bp	411,570,604 bp

262

263

264

265

266

267

268

269



270 **Supplementary Table 2 Statistics of Hi-C sequencing and mapping.**

**A. Raw data aligned to the genome.**

<b>Sample</b>	<b>Hi-C data</b>
Clean Paired-end Reads	3,783,831,838
Unmapped Paired-end Reads	231,750,895
Unmapped Paired-end Reads Rate (%)	6.12
Paired-end Reads with Singleton	691,566,783
Paired-end Reads with Singleton Rate (%)	18.28
Multi Mapped Paired-end Reads	791,574,538
Multi Mapped Paired-end Reads Ratio (%)	20.92
Unique Mapped Paired-end Reads	2,068,939,622
Unique Mapped Paired-end Reads Ratio (%)	54.68

**B. Valid data of reads aligned to the genome.**

<b>Sample</b>	<b>Hi-C data</b>
Unique Mapped Paired-end Reads	2,068,939,622
Dangling End Paired-end Reads	47,562,609
Dangling End Paired-end Reads Rate (%)	2.30
Self Circle Paired-end Reads	5,581,545
Self Circle Paired-end Reads Rate (%)	0.27
Dumped Paired-end Reads	10,525,787
Dumped Paired-end Reads Rate (%)	0.51
Interaction Paired-end Reads	2,005,269,681
Interaction Paired-end Reads Rate (%)	96.92
Valid Paired-end Reads	2,001,492,791
Valid Paired-end Reads Rate (%)	96.74

272

273 **Supplementary Table 3 The lengths of pseudochromosomes generated in the Hi-C**

274 **assembly.**

<b>Pseudochromosome</b>	<b>Length (bp)</b>
Chr01	1,052,277,571
Chr02	995,865,175
Chr03	958,316,394
Chr04	950,298,798
Chr05	904,642,436
Chr06	909,132,773
Chr07	786,775,274
Chr08	774,176,834
Chr09	770,396,028
Chr10	681,618,719
Chr11	660,890,931
Chr12	412,525,553
Total anchored length	9,856,916,486
Unanchored length	381,035,757
Total length	10,237,952,243

275

276

277

278

279

280

281

282

283

284

285

286

287 **Supplementary Table 4 Assembly and annotation statistics of the Taxus genome.**

<b>Parameter</b>	<b>Value</b>
Genome size	10.23 Gb
GC content	36.78%
Contig number	15,879
N50 length (contig)	2.44 Mb
N90 length (contig)	426.40 kb
Longest contig	21,889,323 bp
Gene number	42,746
Average gene length	13,260 bp
Chromosomes	12
Exons per gene	3.5
Total anchored	9,856,916,486 bp
Repeat	76.09%
Anchored rate	96.28%

288

289

290 **Supplementary Table 6 BUSCO analysis.**291 **A. BUSCO analysis of genome annotation (embryophyte\_odb10 (2020-9-10)).**

	<b>Number</b>	<b>Rate (%)</b>
Complete BUSCOs (C)	1,052	65.2
Complete and single-copy BUSCOs (S)	980	60.7
Complete and duplicated BUSCOs (D)	72	4.5
Fragmented BUSCOs (F)	265	16.4
Missing BUSCOs (M)	297	18.4
Total BUSCO groups searched	1,614	100.0

292 **B. BUSCO analysis of the gymnosperm species with genome data.**

Species	Genome Length (Gb)	Contig N50 (kb)	C (%)	S (%)	D (%)	F (%)	M (%)
<i>Pinus lambertiana</i> <sup>31</sup>	31.0	246.6	73.0	67.9	5.1	7.1	19.9
<i>Pinus taeda</i> <sup>32</sup>	20.6	25.4	40.7	35.9	4.8	18.6	40.7
<i>Picea abies</i> <sup>33</sup>	19.6	4.9	27.4	24.0	3.4	25.8	46.8
<i>Abies alba</i> <sup>34</sup>	18.16	14.1	15.4	12.5	2.9	16.9	67.7
<i>Pseudotsuga menziesii</i> <sup>35</sup>	16.6	44.1	67.8	62.5	5.3	11.0	21.2
<i>Ginkgo biloba</i> <sup>36</sup>	10.61	48.2	69.4	61.3	8.1	10.3	20.3
<i>Taxus chinensis</i> var. <i>mairei</i>	9.9	2435.93	65.2	60.7	4.5	16.4	18.4
<i>Sequoiadendron giganteum</i> <sup>37</sup>	8.13	347.95	50.7	47.5	3.2	14.9	34.4
<i>Gnetum montanum</i> <sup>38</sup>	4.07	25.02	83.0	79.3	3.7	4.1	12.9

293 Abbreviation: C, Complete BUSCOs; S, Complete and single-copy BUSCOs; D, Complete and  
 294 duplicated BUSCOs; F, Fragmented BUSCOs; M, Missing BUSCOs.

295

296

297 **Supplementary Table 7 Functional annotation of Taxus genes.**

<b>Database</b>	<b>Number of annotated genes</b>	<b>Percentage (%)</b>
GO	13,694	32.04
InterProScan	33,210	77.69
NR	32,802	76.74
Swiss-Port	28,093	65.72
TAIR	31,126	72.82
Total	36,518	85.43

298

299 Abbreviation: GO, Gene Ontology Database; NR, Non-Redundant Proteins; Swiss-Port, Swiss-  
300 Port Database. TAIR, The Arabidopsis Information Resource.

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

318 **Supplementary Table 8 Statistics of repeated elements in the *Taxus* genome.**

<b>Classification</b>	<b>Number</b>	<b>Length (bp)</b>	<b>Percent of repeats (%)</b>	<b>Percent of genome (%)</b>
Class I: Retroelement	2,712,885	4,219,248,359	54.16	41.21
LTR Retrotransposon	2,571,938	4,080,706,163	52.38	39.86
Ty1/Copia	473,313	504,043,663	6.47	4.92
Ty3/Gypsy	2,082,710	3,561,506,787	45.72	34.79
Other	15,915	15,155,713	0.19	0.15
Non-LTR Retrotransposon	140,947	138,542,196	1.78	1.35
LINE	140,871	138,531,871	1.78	1.35
SINE	76	10,325	0.00	0.00
Class II: DNA elements	299,098	283,535,871	3.64	2.77
CMC	7,686	1,266,320	0.02	0.01
hAT	81,774	69,108,347	0.89	0.68
PIF/Harbinger	219	124,131	0.00	0.00
MuLE-MuDR	135,519	158,064,978	2.03	1.54
Sola	19,436	25,932,250	0.33	0.25
Helitron	36,248	26,622,352	0.34	0.26
Enspm	17,812	2,384,418	0.03	0.02
Mariner	404	33,075	0.00	0.00
Small RNA	502	73,190	0.00	0.00
Simple repeats	53,662	26,564,977	0.34	0.26
Unknown	5,208,975	3,122,338,811	40.08	30.50
Total repeat fraction	8,416,069	7,790,303,404	100.00	76.09

320 **Supplementary Table 27 Gene cluster identified through PlantiSMASH in the *Taxus***

321 **genome.**

<b>Type</b>	<b>Number</b>	<b>Chromosome</b>
putative	9	1, 3, 9, 10
terpene	7	1, 3, 6, 8, 9
saccharide	13	2, 6, 8, 9, 10, 12
alkaloid	1	1
saccharide-polyketide	1	2
terpene-alkaloid	1	4
lignan-terpene	1	5
saccharide-terpene	1	10

322

323 **4. Supplementary References 1-38**

324

- 325 1 Li, Y. *et al.* Induction of half-sib embryonic callus and production of taxiod compounds from  
326 *Taxus chinensis* var. *mairei*. *Int. J. Agric. Biol.* **21**, 719-725, doi:10.17957/IJAB/15.0949  
327 (2019).
- 328 2 Marçais, G. & Kingsford, C. A fast, lock-free approach for efficient parallel counting of  
329 occurrences of k-mers. *Bioinformatics* **27**, 764-770, doi:10.1093/bioinformatics/btr011  
330 (2011).
- 331 3 Li, H. Minimap2: pairwise alignment for nucleotide sequences. *Bioinformatics* **34**, 3094-  
332 3100, doi:10.1093/bioinformatics/bty191 (2018).
- 333 4 Danecek, P. *et al.* Twelve years of SAMtools and BCFtools. *Gigascience* **10**, 1-4,  
334 doi:10.1093/gigascience/giab008 (2021).
- 335 5 Xu, Z. & Wang, H. LTR\_FINDER: an efficient tool for the prediction of full-length LTR  
336 retrotransposons. *Nucleic Acids Res.* **35**, 265-268, doi:10.1093/nar/gkm286 (2007).
- 337 6 Ellinghaus, D., Kurtz, S. & Willhoeft, U. LTRharvest, an efficient and flexible software for  
338 de novo detection of LTR retrotransposons. *BMC Bioinform.* **9**, 18, doi:10.1186/1471-2105-  
339 9-18 (2008).
- 340 7 Ou, S. & Jiang, N. LTR\_retriever: a highly accurate and sensitive program for identification  
341 of long terminal repeat retrotransposons. *Plant Physiol.* **176**, 1410-1422,  
342 doi:10.1104/pp.17.01310 (2018).
- 343 8 Potter, S. C. *et al.* HMMER web server: 2018 update. *Nucleic Acids Res.* **46**, 200-204,  
344 doi:10.1093/nar/gky448 (2018).
- 345 9 Edgar, R. C. MUSCLE: multiple sequence alignment with high accuracy and high  
346 throughput. *Nucleic Acids Res.* **32**, 1792-1797, doi:10.1093/nar/gkh340 (2004).
- 347 10 Arkin, A. P. FastTree 2 - approximately maximum-likelihood trees for large alignments.  
348 *PLoS One* **5**, e9490, doi:10.1371/journal.pone.0009490 (2010).
- 349 11 Keller, I., Bensasson, D. & Nichols, R. A. Transition-transversion bias is not universal: a  
350 counter example from grasshopper pseudogenes. *PLoS Genet.* **3**, e22,  
351 doi:10.1371/journal.pgen.0030022 (2007).
- 352 12 Rice, J., Longden, I. & Bleasby, A. EMBOSS: the European molecular biology open software  
353 suite. *Trends Genet.* **16**, 276-277, doi: 10.1016/s0168-9525(00)02024-2 (2000).



- 354 13 De La Torre, A. R., Li, Z., Van de Peer, Y. & Ingvarsson, P. K. Contrasting rates of molecular  
355 evolution and patterns of selection among gymnosperms and flowering plants. *Mol. Biol.*  
356 *Evol.*, 1363-1377, doi:10.1093/molbev/msx069 (2017).
- 357 14 Kim, D., Langmead, B. & Salzberg, S. L. HISAT: a fast spliced aligner with low memory  
358 requirements. *Nat. Methods* **12**, 357-360, doi:10.1038/nmeth.3317 (2015).
- 359 15 Pertea, M., Kim, D., Pertea, G. M., Leek, J. T. & Salzberg, S. L. Transcript-level expression  
360 analysis of RNA-seq experiments with HISAT, StringTie and Ballgown. *Nat. Protoc.* **11**,  
361 1650-1667, doi:10.1038/nprot.2016.095 (2016).
- 362 16 Robinson, M. D., McCarthy, D. J. & Smyth, G. K. edgeR: a Bioconductor package for  
363 differential expression analysis of digital gene expression data. *Bioinformatics* **26**, 139-140,  
364 doi:10.1093/bioinformatics/btp616 (2009).
- 365 17 Wegrzyn, J. L., Lee, J. M., Tearse, B. R. & Neale, D. B. TreeGenes: a forest tree genome  
366 database. *Int. J. Plant Genomics* **2008**, 412875, doi:10.1155/2008/412875 (2008).
- 367 18 Altschul, S. F., Gish, W., Miller, W., Myers, E. W. & Lipman, D. J. Basic local alignment  
368 search tool. *J. Mol. Biol.* **215**, 403-410, doi:10.1016/S0022-2836(05)80360-2 (1990).
- 369 19 Li, L., Stoeckert, C. J., Jr. & Roos, D. S. OrthoMCL: identification of ortholog groups for  
370 eukaryotic genomes. *Genome Res.* **13**, 2178-2189, doi:10.1101/gr.1224503 (2003).
- 371 20 Edgar, R. MUSCLE: multiple sequence alignment with high accuracy and high throughput.  
372 *Nucleic Acids. Res.* **32**, 1792-1797, doi:10.1093/nar/gkh340 (2004).
- 373 21 Castresana, J. Selection of conserved blocks from multiple alignments for their use in  
374 phylogenetic analysis. *Mol. Biol. Evol.* **17**, 540-552,  
375 doi:10.1093/oxfordjournals.molbev.a026334 (2000).
- 376 22 Stamatakis, A. RAxML version 8: a tool for phylogenetic analysis and post-analysis of large  
377 phylogenies. *Bioinformatics* **30**, 1312-1313, doi:10.1093/bioinformatics/btu033 (2014).
- 378 23 Sanderson, M. J. r8s: inferring absolute rates of molecular evolution and divergence times in  
379 the absence of a molecular clock. *Bioinformatics* **19**, 301-302,  
380 doi:10.1093/bioinformatics/19.2.301 (2003).
- 381 24 De Bie, T., Cristianini, N., Demuth, J. P. & Hahn, M. W. CAFE: a computational tool for the  
382 study of gene family evolution. *Bioinformatics* **22**, 1269-1271,  
383 doi:10.1093/bioinformatics/btl097 (2006).

- 384 25 Yu, G., Wang, L. G., Han, Y. & He, Q. Y. clusterProfiler: an R package for comparing  
385 biological themes among gene clusters. *OMICS* **16**, 284-287, doi:10.1089/omi.2011.0118  
386 (2012).
- 387 26 Urban, P., Mignotte, C., Kazmaier, M., Delorme, F. & Pompon, D. Cloning, yeast expression,  
388 and characterization of the coupling of two distantly related *Arabidopsis thaliana* NADPH-  
389 Cytochrome P450 reductases with P450 CYP73A5. *J. Biol. Chem.* **272**, 19176-19186,  
390 doi:10.1074/jbc.272.31.19176 (1997).
- 391 27 Pompon, D., Louerat, B., Bronine, A. & Urban, P. in *Methods in Enzymology* Vol. **272** (eds  
392 Eric F. Johnson & Michael R. Waterman) 51-64 (Academic Press, 1996).
- 393 28 Biggs, B. W. *et al.* Orthogonal assays clarify the oxidative biochemistry of taxol P450  
394 CYP725A4. *ACS Chem. Biol.* **11**, 1445-1451, doi:10.1021/acscchembio.5b00968 (2016).
- 395 29 Bian, G. *et al.* Production of taxadiene by engineering of mevalonate pathway in *Escherichia*  
396 *coli* and endophytic fungus *Alternaria alternata* TPF6. *Biotechnol. J.* **12**,  
397 doi:10.1002/biot.201600697 (2017).
- 398 30 Kuang, X., Sun, S., Wei, J., Li, Y. & Sun, C. Iso-Seq analysis of the *Taxus cuspidata*  
399 transcriptome reveals the complexity of Taxol biosynthesis. *BMC Plant Biol.* **19**, 210,  
400 doi:10.1186/s12870-019-1809-8 (2019).
- 401 31. Stevens, K. A. *et al.* Sequence of the sugar pine megagenome. *Genetics* **204**, 1613-1626,  
402 doi:10.1534/genetics.116.193227 (2016).
- 403 32. Zimin, A. V. *et al.* An improved assembly of the loblolly pine mega-genome using long-read  
404 single-molecule sequencing. *Gigascience* **6**, 1-4, doi: 10.1093/gigascience/giw016 (2017).
- 405 33. Nystedt, B. *et al.* The Norway spruce genome sequence and conifer genome evolution.  
406 *Nature* **497**, 579-584, doi:10.1038/nature12211 (2013).
- 407 34. Mosca, E. *et al.* A reference genome sequence for the european silver fir (*Abies alba* Mill.): A  
408 community-generated genomic resource. *G3* **9**, 2039-2049, doi: 10.1534/g3.119.400083  
409 (2019).
- 410 35. Neale, D. B. *et al.* The Douglas-Fir genome sequence reveals specialization of the  
411 photosynthetic apparatus in Pinaceae. *G3* **7**, 3157-3167, doi:10.1534/g3.117.300078 (2017).
- 412 36 Guan, R. *et al.* Draft genome of the living fossil *Ginkgo biloba*. *Gigascience* **5**, 49,  
413 doi:10.1186/s13742-016-0154-1 (2016).

- 414 37 Scott, A. D. et al. A reference genome sequence for giant sequoia. *G3* **10**, 3907-3919,  
415 doi:10.1534/g3.120.401612 (2020).
- 416 38 Wan, T. et al. A genome for gnetophytes and early evolution of seed plants. *Nat. Plants* **4**, 82-  
417 89, doi:10.1038/s41477-017-0097-2 (2018).