

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

1 ***De novo* genome assembly of *Camptotheca acuminata*, a natural source of the anti-cancer**
2 **compound camptothecin**

3 Dongyan Zhao¹, John P. Hamilton¹, Gina M. Pham¹, Emily Crisovan¹, Krystle Wiegert-Rininger¹,
4 Brieanne Vaillancourt¹, Dean DellaPenna², and C. Robin Buell^{1*}

5 ¹Department of Plant Biology, Michigan State University, East Lansing, MI 48824 USA

6 ²Department of Biochemistry & Molecular Biology, Michigan State University, East Lansing, MI
7 48824 USA

8 **Email addresses:** Dongyan Zhao <zhaodon4@msu.edu>, John P. Hamilton <jham@msu.edu>,
9 Gina M. Pham <phamgina@msu.edu>, Emily Crisovan <pankeyem@msu.edu>, Krystle Wiegert-
10 Rininger <wiegertk@msu.edu>, Brieanne Vaillancourt <vaillan6@msu.edu>, Dean Dellapenna
11 <dellapen@msu.edu>, C Robin Buell <buell@msu.edu>

12 *Correspondence should be addressed to: C. Robin Buell, buell@msu.edu

13
14 **Manuscript type:** Data note

15
16 **Note:** Reviewers can access the genome sequence and annotation using the following
17 temporary URL: <http://datadryad.org/review?doi=doi:10.5061/dryad.nc8qr>.

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

18 **Abstract**

19 **Background:** *Camptotheca acuminata* is one of a limited number of species that produce
20 camptothecin, a pentacyclic quinoline alkaloid with anti-cancer activity due to its ability to
21 inhibit DNA topoisomerase. While transcriptome studies have been performed previously with
22 various camptothecin-producing species, no genome sequence for a camptothecin-producing
23 species is available to date.

24 **Findings:** We generated a high quality *de novo* genome assembly for *C. acuminata* representing
25 403,174,860 bp on 1,394 scaffolds with an N50 scaffold size of 1,752 kbp. Quality assessments
26 of the assembly revealed robust representation of the genome sequence including genic
27 regions. Using a novel genome annotation method, we annotated 31,825 genes encoding
28 40,332 gene models. Based on sequence identity and orthology with validated genes from
29 *Catharanthus roseus* as well as Pfam searches, we identified candidate orthologs for genes
30 potentially involved in camptothecin biosynthesis. Extensive gene duplication including tandem
31 duplication was widespread in the *C. acuminata* genome with 3,315 genes belonging to 1,245
32 tandem duplicated gene clusters.

33 **Conclusions:** To our knowledge, this is the first genome sequence for a camptothecin-producing
34 species, and access to the *C. acuminata* genome will permit not only discovery of genes
35 encoding the camptothecin biosynthetic pathway but also reagents that can be used for
36 heterologous expression of camptothecin and camptothecin analogs with novel pharmaceutical
37 applications.

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

38 **Keywords:** *Camptotheca acuminata*, camptothecin, genome assembly, genome annotation,
39 tandem duplications

40
41 **Data Description**

42 **Background information on camptothecin, a key anti-cancer natural product**

43 *Camptotheca acuminata* Decne, also known as the Chinese Happy Tree (Figure 1), is an eudicot
44 asterid Cornales tropical tree species within the Nyssaceae family [1] that also contains *Nyssa*
45 spp (tupelo) and *Davidia involucrate* (dove tree); no genome sequence is available for any
46 member of this family. *C. acuminata* is one of a limited number of plant species that produce
47 camptothecin, a pentacyclic quinoline alkaloid (Figure 2A) with anti-cancer activity due to its
48 ability to inhibit DNA topoisomerase [2]. Due to poor solubility, derivatives such as irinotecan
49 and topotecan, rather than camptothecin are currently in use as approved cancer drugs. The
50 significance of these derivatives as therapeutics is highlighted by the listing of irinotecan on the
51 World Health Organization Model List of Essential Medicines [3]. While transcriptome studies
52 have been performed previously with various camptothecin-producing species including *C.*
53 *acuminata* and *Ophiorrhiza pumila* (e.g., [4-6]), no genome sequence for a camptothecin-
54 producing species is available to date. We report on the assembly and annotation of the *C.*
55 *acuminata* genome, the characterization of genes implicated in camptothecin biosynthesis, and
56 highlight the extent of gene duplication that provides new templates for gene diversification.

57 **RNA isolation, library construction, sequencing, and transcriptome assembly**

1
2
3
4 58 Transcriptome assemblies were constructed using nine developmental RNA-sequencing (RNA-
5
6
7 59 seq) datasets described in a previous study [4] that included immature bark, cotyledons,
8
9
10 60 immature flower, immature fruit, mature fruit, mature leaf, root, upper stem, and lower stem.
11
12 61 Adapters and low-quality nucleotides were removed from the RNA-seq reads using Cutadapt
13
14 62 (v1.8) [7] and contaminating ribosomal RNA reads were removed. Cleaned reads from all nine
15
16
17 63 libraries were assembled using Trinity (v20140717) [8] with a normalization factor of 50x using
18
19
20 64 default parameters. Contaminant transcripts (5,669 total) were identified by searching the *de*
21
22 65 *novo* transcriptome assembly against the National Center for Biotechnology Information (NCBI)
23
24
25 66 non-redundant nucleotide database using BLAST+ (v2.2.30) [9, 10] with an E-value cutoff of 1e-
26
27
28 67 5; transcripts with their best hits being a non-plant sequence were removed from the
29
30 68 transcriptome.

31
32
33
34 69 For additional transcript support for use in a genome-guided transcriptome assembly to
35
36 70 support genome annotation, strand-specific RNA-seq reads were generated by isolating RNA
37
38
39 71 from root tissues and sequencing of Kappa TruSeq Stranded libraries on an Illumina HiSeq 2500
40
41 72 platform generating 150 nt paired-end reads (BioSample ID: SAMN06229771). Root RNA-seq
42
43
44 73 reads were assessed for quality using FASTQC (v0.11.2) [11] using default parameters and
45
46 74 cleaned as described above.

50 75 **DNA isolation, library construction, and sequencing**

51
52
53 76 The genome size of *C. acuminata* was estimated at 516 Mb using flow cytometry, suitable for
54
55
56 77 *de novo* assembly using the Illumina platform. DNA was extracted from young leaves of *C.*
57
58
59 78 *acuminata* at the vegetative growth stage using CTAB [12]. Multiple Illumina-compatible paired-

1
2
3
4 79 end libraries (Table 1) with insert sizes ranging from 180-609 bp were constructed as described
5
6
7 80 previously [13] and sequenced to 150 nt in paired-end mode on an Illumina HiSeq2000. Mate-
8
9
10 81 pair libraries (Table 1) with size ranges of 1.3-8.9 kb were made using the Nextera Kit (Illumina,
11
12 82 San Diego CA) as per manufacturer's instructions and sequenced to 150 nt in paired-end mode
13
14
15 83 on an Illumina HiSeq2000.

18 84 **Genome assembly**

20
21 85 Paired-end reads (Table 1) were assessed for quality using FASTQC (v0.11.2) [11] using default
22
23
24 86 parameters, cleaned for adapters and low quality sequences using Cutadapt (v1.8) [7] and only
25
26
27 87 reads in pairs with each read ≥ 25 nt were retained for genome assembly. Mate pair libraries
28
29
30 88 (Table 1) were processed using NextClip (v1.3.1) [14] and only reads from Categories A, B, C
31
32 89 were used for the assembly. Using ALLPATHS-LG (v44837) [15] with default parameters, two
33
34
35 90 paired-end read libraries (180 and 268 bp insert libraries) and all five mate pair libraries (Table
36
37 91 1) were used to generate an initial assembly of 403.2 Mb with an N50 contig size of 108 kbp
38
39
40 92 and an N50 scaffold size of 1,752 kbp (Tables 1 and 2). Gaps (5,076) in this initial assembly were
41
42
43 93 filled using SOAP GapCloser (v1.12r6) [16] with four independent paired-end libraries (352, 429,
44
45 94 585, and 609 bp inserts, Table 1); 12,468,362 bp of the estimated 16,471,841 bp of gaps was
46
47
48 95 filled leaving a total of 3,825 gaps (3,772,191 Ns). The assembly was checked for contaminant
49
50
51 96 sequences based on alignments to the NCBI non-redundant nucleotide database using BLASTN
52
53 97 (E-value = $1e-5$) [10]; a single scaffold of 5,156 bp that matched a bacterium sequence with
54
55 98 100% coverage and 100% identity was removed. Subsequently, five scaffolds less than 1 kbp
56
57
58
59
60
61
62
63
64
65

1
2
3
4 99 were removed resulting in the final assembly of 403,174,860 bp comprised of 1,394 scaffolds
5
6
7 100 with an N50 scaffold size of 1,752 kbp (Tables 1 and 2) and 0.9% Ns.
8
9

10 101 Quality assessments revealed a robust high quality assembly with 98% of the paired-end
11
12
13 102 genomic sequencing reads aligning to the assembly, of which, 99.97% aligned concordantly.
14
15 103 With respect to genic representation, 95.3% of RNA-seq-derived transcript assemblies [4] and
16
17
18 104 74,119 of 74,682 (99%) pyrosequencing transcript reads from a separate study [5] aligned to
19
20
21 105 the genome assembly. A total of 93.6% of conserved Embryophyta BUSCO proteins were
22
23 106 present in the assembly as full-length sequences with an additional 2.5% of the Embryophyta
24
25
26 107 proteins fragmented [17].
27
28

29 108 **Genome annotation**

30
31
32

33 109 We used a novel genome annotation method to generate high quality annotation of the *C.*
34
35 110 *acuminata* genome in which we repeat masked the genome, trained an *ab initio* gene finder
36
37
38 111 with a genome-guided transcript assembly, and then refined the gene models using additional
39
40
41 112 genome-guided transcript assembly evidence to generate a high quality gene model set. We
42
43 113 first created a *C. acuminata* specific custom repeat library (CRL) using MITE-Hunter (v2011) [18]
44
45
46 114 and RepeatModeler (v1.0.8) [19]. Protein coding genes were removed from each repeat library
47
48 115 using ProtExcluder.pl (v1.1) [20] and combined into a single CRL, which hard-masked 143.6 Mb
49
50
51 116 (35.6%) of the assembly as repetitive sequence using RepeatMasker (v4.0.6) [21]. Cleaned root
52
53
54 117 RNA-seq reads (Table S1, BioSample ID: SAMN06229771) were aligned to the genome assembly
55
56 118 using TopHat2 (v2.0.13) [22] in strand-specific mode with a minimum intron length of 20 bp and
57
58
59 119 a maximum intron length of 20 kb; the alignments were then used to create a genome-guided
60
61
62
63
64
65

1
2
3
4 120 transcriptome assembly using Trinity (v2.2.0) [23]. The RNA-seq alignments were used to train
5
6
7 121 AUGUSTUS (v3.1) [24] and gene predictions were generated with AUGUSTUS [25] using the
8
9
10 122 hard-masked assembly. Gene model structures were refined by incorporating evidence from
11
12 123 the genome-guided transcriptome assembly using PASA2 (v2.0.2) [26, 27]; with the parameters:
13
14 124 MIN_PERCENT_ALIGNED=90, MIN_AVG_PER_ID=99. After annotation comparison, models that
15
16
17 125 PASA identified as being merged and a subset of candidate camptothecin biosynthetic
18
19
20 126 pathways genes identified as mis-annotated were manually curated. The final high-confidence
21
22 127 gene model set consists of 31,825 genes encoding 40,332 gene models. Functional annotation
23
24
25 128 was assigned using a custom pipeline using WU-BLASTP [28] searches against the *Arabidopsis*
26
27
28 129 *thaliana* annotation (TAIR10; [29]) and Swiss-Prot plant proteins (downloaded on 08-17-2015),
29
30 130 and a search against Pfam (v29) using HMMER (v3.1b2) [30]. This resulted in 34,143 gene
31
32
33 131 models assigned a putative function, 2,011 annotated as conserved hypothetical, and 4,178
34
35 132 annotated as hypothetical.

36
37
38
39 133 *C. acuminata* is insensitive to camptothecin due to mutations within its own DNA
40
41 134 topoisomerase [31] and we identified two topoisomerase genes in our annotated gene set, one
42
43
44 135 of which matches the published *C. acuminata* topoisomerase (99.78% identity, 100% coverage)
45
46
47 136 and includes the two mutations that confer resistance to camptothecin (Figure 2B), one
48
49 137 mutation is specific in *C. acuminata* and the other is present in both *C. acuminata* and two
50
51
52 138 camptothecin-producing *Ophiorrhiza* species. Further quality assessments of our annotation
53
54 139 with 35 nuclear-encoded *C. acuminata* genes available from GenBank revealed an average
55
56
57 140 identity of 99.5% with 100% coverage in our annotated proteome while a single gene encoding
58
59 141 1-deoxy-D-xylulose 5-phosphate reductoisomerase (ABC86579.1) had 88.2% identity with 100%

1
2
3
4 142 coverage that may be attributable to differences in genotypes. One mRNA reported to encode a
5
6
7 143 putative strictosidine beta-D-glucosidase (AES93119.1) was found to have a retained intron that
8
9
10 144 when removed, aligned with 99.3% identity yet reduced coverage (66%) as it was located at the
11
12 145 end of a short scaffold. Collectively, the concordant alignment of whole genome shotgun
13
14 146 sequence reads to the assembly, the high representation of genic regions as assessed by
15
16
17 147 independent transcriptome datasets (RNA-seq and pyrosequencing) as well as the core
18
19
20 148 Embryophyta BUSCO proteins, when coupled with the high quality gene models as revealed
21
22 149 through alignments with cloned *C. acuminata* genes indicate that we have not only generated a
23
24
25 150 high quality genome assembly for *C. acuminata* but also a robust set of annotated gene models.
26
27

28 151 **Gene duplication and orthology analyses**

30
31
32 152 During our annotation efforts, it was readily apparent that there was substantial gene
33
34 153 duplication including tandem gene duplication in the *C. acuminata* genome. Paralogous
35
36
37 154 clustering of the *C. acuminata* proteome revealed 5,768 paralogous groups containing 17,957
38
39
40 155 genes. We identified tandem gene duplications in the *C. acuminata* genome based on if: 1) two
41
42 156 or more *C. acuminata* genes were present within an orthologous/paralogous group; 2) there
43
44
45 157 were no more than 10 genes in between on a single scaffold; and 3) the pairwise gene distance
46
47
48 158 was less than 100 kbp [32]. Under these criteria, a total of 3,315 genes belonging to 1,245
49
50 159 tandem duplicated gene clusters were identified. Gene ontology analysis showed that tandem
51
52
53 160 duplicated genes are significantly enriched in “response to stress” ($p < 0.0001$, χ^2 test) while
54
55 161 under-represented in most other processes, especially “other cellular processes” and “cell
56
57
58 162 organization and biogenesis” ($p < 0.0001$, χ^2 test).
59
60
61
62
63
64
65

1
2
3
4 163 To our knowledge, *C. acuminata* is the first species within the Nyssaceae family with a genome
5
6
7 164 sequence. To better understand the evolutionary relationship of *C. acuminata* with other
8
9
10 165 asterids and angiosperms, we identified orthologous and paralogous groups using our
11
12 166 annotated *C. acuminata* proteome and the proteomes of five other key species (*Arabidopsis*
13
14 167 *thaliana*, *Amborella trichopoda*, *Vitis vinifera*, *Oryza sativa*, and *Catharanthus roseus*) using
15
16
17 168 OrthoFinder (v0.7.1) [33] with default parameters. A total of 12,459 orthologous groups
18
19
20 169 containing at least a single *C. acuminata* protein were identified with 8,521 orthologous groups
21
22 170 common to all six species (Figure 3; Table S2). Interestingly, *C. acuminata* contains the least
23
24
25 171 number of singleton genes (7,177) among the six species, and gene ontology analysis
26
27
28 172 demonstrated that these genes were highly enriched in “transport”, “response to stress”, and
29
30 173 “other cellular and biological processes” ($p < 0.0001$, χ^2 test) while dramatically under-
31
32
33 174 represented in “unknown biological processes” ($p < 0.0001$, χ^2 test), suggesting these genes
34
35 175 may be involved in stress responses and other processes specific to *C. acuminata*.

36 37 38 39 176 **Uses for the *C. acuminata* genome sequence and annotation**

40
41
42 177 Generation of a high-quality genome sequence and annotation dataset for *C. acuminata* will
43
44
45 178 facilitate discovery of genes encoding camptothecin biosynthesis as physical clustering can be
46
47
48 179 combined with sequence similarity and co-expression data to identify candidate genes, an
49
50 180 approach that has been extremely useful in identifying genes in specialized metabolism in a
51
52
53 181 number of plant species (see [34-36]). In *C. acuminata*, geranylgeranyl diphosphate from the 2-
54
55 182 C-methyl-D-erythritol 4-phosphate/1-deoxy-D-xylulose 5-phosphate (MEP) pathway is used to
56
57
58 183 generate secologanic acid via the iridoid pathway and tryptamine from tryptophan
59
60
61
62
63
64
65

1
2
3
4 184 decarboxylase are condensed by strictosidinic acid synthase to generate strictosidinic acid that
5
6
7 185 is then converted into camptothecin in the alkaloid pathway via a set of unknown steps [37]
8
9
10 186 (Figure 4A). *Catharanthus roseus*, Madagascar periwinkle, produces vinblastine and vincristine
11
12 187 via the MEP and iridoid pathways for which all genes leading to the biosynthesis of the iridoid
13
14
15 188 secologanin have been characterized [35]. Using sequence identity and coverage with
16
17 189 characterized *C. roseus* genes from the MEP and iridoid pathway (Figure 4A), we were able to
18
19
20 190 identify candidate genes for all steps in the MEP and iridoid pathway in *C. acuminata* (Table 3).
21
22 191 The downstream steps in camptothecin biosynthesis subsequent to formation of strictosidinic
23
24
25 192 acid involve a broad set of enzymes responsible for reduction and oxidation [37] and a total of
26
27
28 193 343 cytochrome P450s (60 paralogous gene clusters and 86 singletons; Table S3) were
29
30 194 identified which can serve as candidates for the later steps in camptothecin biosynthesis.
31
32
33 195 Though not absolute, physical clustering of genes involved in specialized metabolism has been
34
35
36 196 observed in a number of species across a number of classes of specialized metabolites [34, 38].
37
38
39 197 With an N50 scaffold size of 1,752 kbp, we observed several instances of physical clustering of
40
41
42 198 genes with homology to genes involved in monoterpene indole alkaloid biosynthesis which may
43
44 199 produce related compounds in *C. acuminata*. Using characterized genes involved in the
45
46
47 200 biosynthesis of vinblastine and vincristine from *C. roseus* as queries [35] (Figure 4A, Table 3), we
48
49 201 identified a single *C. acuminata* scaffold (907 kbp, 86 genes; Figure 4B) that encoded genes with
50
51
52 202 sequence identity to isopentenyl diphosphate isomerase II within the MEP pathway, 8-
53
54
55 203 hydroxygeraniol oxidoreductase (GOR, three complete and one partial paralogs), 7-
56
57 204 deoxyloganic acid 7-hydroxylase (7DLH) within the iridoid pathway, and a protein with
58
59
60 205 homology to *C. roseus* 16-hydroxy-2,3-dihydro-3-hydroxytabersonine N-methyltransferase
61
62
63
64
65

1
2
3
4 206 (NMT) within the alkaloid pathway suggesting that access to a high contiguity genome assembly
5
6
7 207 may facilitate discovery of genes involved in specialized metabolism in *C. acuminata*. Tandem
8
9
10 208 duplications of genes involved in specialized metabolism have been reported previously [39, 40]
11
12 209 and via divergence either in the coding region or promoter sequence which lead to neo- and
13
14
15 210 sub-functionalization at the enzymatic or expression level, respectively, have been shown to
16
17 211 contribute to the extensive chemical diversity within a species [40, 41].
18
19

20
21 212 The *C. acuminata* genome can also be used to facilitate our understanding of the mechanisms
22
23 213 by which camptothecin production evolved independently in distinct taxa such as *C. acuminata*
24
25
26 214 (*Nyssaceae*) and *O. pumila* (*Rubiaceae*). For example, a comparative analysis of *C. acuminata*
27
28
29 215 and *O. pumila* may be highly informative in not only delineating genes involved in camptothecin
30
31 216 biosynthesis but also in revealing key evolutionary events that led to biosynthesis of this critical
32
33
34 217 natural product across a wide phylogenetic distance. As noted above, camptothecin is
35
36 218 cytotoxic and as a consequence, derivatives of camptothecin are used as anti-cancer drugs.
37
38
39 219 Perhaps most exciting, the ability to decipher the full camptothecin biosynthetic pathway will
40
41 220 yield molecular reagents that can be used to not only synthesize camptothecin in heterologous
42
43
44 221 systems such as yeast, but also produce less toxic analogs with novel pharmaceutical
45
46 222 applications.
47
48

50 223 **Availability of Supporting Information**

51
52

53 224 Raw genomic sequence reads and transcriptome reads derived from root tissues are available
54
55
56 225 in the NCBI Sequence Read Archive under project number PRJNA361128. All other RNA-seq
57
58
59 226 transcriptome reads were from Bioproject PRJNA80029 [4]. The genome assembly and
60
61
62
63
64
65

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

227 annotation are available in the Dryad Digital Repository under doi (to be released upon
228 publication), through the Medicinal Plant Genomics Resource [42] via a genome browser and
229 search and analysis tools, and GigaDB.

230 **Note:** Reviewers can access the genome sequence and annotation using the following
231 temporary URL: <http://datadryad.org/review?doi=doi:10.5061/dryad.nc8qr>.

232

233 **Abbreviations**

234 2-C-methyl-D-erythritol 4-phosphate/1-deoxy-D-xylulose 5-phosphate (MEP), 7-deoxyloganic
235 acid 7-hydroxylase (7DLH), 8-hydroxygeraniol oxidoreductase (GOR), 16-hydroxy-2,3-dihydro-3-
236 hydroxytabersonine N-methyltransferase (NMT), custom repeat library (CRL), National Center
237 for Biotechnology Information (NCBI), RNA-sequencing (RNA-seq)

238 **Competing Interests**

239 The authors have declared that no competing interests exists.

240 **Author Contributions**

241 CRB oversaw the project. DZ performed the genome assembly, assisted in genome annotation
242 and analyzed data. JH annotated the genome and analyzed data. EC, GP, and KWR constructed
243 libraries and analyzed data. BV analyzed data. DDP provided intellectual oversight. DZ, JH, and
244 CRB wrote the manuscript.

245 **Acknowledgements**

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

246 Funding for this work was provided in part by a grant to CRB and DDP from the National
247 Institute of General Medical Sciences (1RC2GM092521) and funds to CRB and DDP from
248 Michigan State University. The funders had no role in study design, data collection and analysis,
249 decision to publish, or preparation of the manuscript.

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

250 **Figure Legends**

251 **Figure 1. *Camptotheca acuminata* Decne, the Chinese Happy Tree, is a member in the**
252 **Nyssaceae family that produces the anticancer compound camptothecin.**

253 **Figure 2. Genome aspects of *Camptotheca acuminata*. (A) Structure of camptothecin. (B) Key**
254 **amino acid mutations (red rectangles) in DNA topoisomerase I in camptothecin-producing**
255 **and non-producing species and their phylogenetic relationship.**

256 **Figure 3. Venn diagram showing orthologous and paralogous groups between *Camptotheca***
257 ***acuminata*, *Amborella trichopoda*, *Oryza sativa*, *Arabidopsis thaliana*, *Vitis vinifera*, and**
258 ***Catharanthus roseus*.**

259 **Figure 4. Key portions of the proposed camptothecin biosynthetic pathway and an example of**
260 **physical clustering of candidate genes in *Camptotheca acuminata*. (A) The methylerythritol**
261 **phosphate (MEP) pathway (green), iridoid pathway (blue), and condensation of secologanic**
262 **acid with tryptamine via strictosidinic acid synthase (STRAS) to form strictosidinic acid prior**
263 **to downstream dehydration, reduction, and oxidation steps yielding camptothecin. DXS, 1-**
264 **deoxy-D-xylulose 5-phosphate synthase 2; DXR, 1-deoxy-D-xylulose-5-phosphate**
265 **reductoisomerase; CMS, 4-diphosphocytidyl-methylerythritol 2-phosphate synthase; CMK, 4-**
266 **diphosphocytidyl-2-C-methyl-D-erythritol kinase; MCS, 2C-methyl-D-erythritol 2,4-**
267 **cyclodiphosphate synthase; HDS, GCPE protein; HDR, 1-hydroxy-2-methyl-butenyl 4-**
268 **diphosphate reductase; IPI, plastid isopentenyl pyrophosphate, dimethylallyl pyrophosphate**
269 **isomerase; GPPS, geranyl pyrophosphate synthase; GES, plastid geraniol synthase; G8H,**
270 **geraniol 8-hydroxylase; GOR, 8-hydroxygeraniol oxidoreductase; CYC1, iridoid cyclase 1; 7-DLS,**

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

271 7-deoxyloganetic acid synthase; 7-DLGT, 7-deoxyloganetic acid glucosyltransferase; 7-DLH, 7-deoxyloganic acid hydroxylase; SLAS, secologanic acid synthase; TDC, tryptophan decarboxylase.

273 **(B) Physical clustering of homologs of genes involved in the methylerythritol phosphate, iridoid, and alkaloid biosynthetic pathways of *Catharanthus roseus* on scaffold 151 of *C.***

274 ***acuminata*.** GOR: 8-hydroxygeraniol oxidoreductase; NMT: 16-hydroxy-2,3-dihydro-3-hydroxytabersonine N-methyltransferase; 7DLH: 7-deoxyloganic acid 7-hydroxylase; IPP2: isopentenyl diphosphate isomerase II. Gene IDs are below the arrows.

278

279 **Table 1. Input libraries and sequences for *de novo* assembly of the *Camptotheca acuminata***
 280 **genome.**

BioProject ID	BioSample ID	Fragment size (bp)	No. of cleaned read pairs	Use
Paired end				
PRJNA361128	SAMN06220985	180	96,955,546	ALLPATHS-LG assembly
PRJNA361128	SAMN06220986	268	89,381,055	ALLPATHS-LG assembly
PRJNA361128	SAMN06220987	352	61,207,691	GapCloser
PRJNA361128	SAMN06220988	429	50,688,562	GapCloser
PRJNA361128	SAMN06220989	585	21,856,610	GapCloser
PRJNA361128	SAMN06220990	609	22,217,954	GapCloser
Mate pair				
PRJNA361128	SAMN06220991	8,111	9,923,643	ALLPATHS-LG assembly
PRJNA361128	SAMN06220992	7,911	7,652,519	ALLPATHS-LG assembly
PRJNA361128	SAMN06220993	1,377	12,800,554	ALLPATHS-LG assembly
PRJNA361128	SAMN06220994	3,179	13,138,503	ALLPATHS-LG assembly
PRJNA361128	SAMN06220995	8,879	13,599,241	ALLPATHS-LG assembly

All libraries were sequenced in paired end mode generating 150 nt reads.

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

282 **Table 2. Metrics of the final assembly of *Camptotheca acuminata* genome.**

Metric	Value
Total scaffold length (bp)	403,174,860
Total no. of scaffolds (bp)	1,394
Maximum scaffold length (bp)	8,423,530
Minimum scaffold length (bp)	1,002
N50 scaffold size (bp)	1,751,747
N50 contig size (bp)	107,594
No. Ns	3,772,191 (0.9%)
No. gaps	3,825

283

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

284 **Table 3. Identification of candidate camptothecin biosynthetic pathway genes in the *Camptotheca acuminata* genome as revealed**
 285 **by sequence identity and coverage with characterized genes from the 2-C-methyl-D-erythritol 4-phosphate/1-deoxy-D-xylulose 5-**
 286 **phosphate and iridoid biosynthetic pathways from *Catharanthus roseus*.**

Description	Abbreviation	Protein	Camptotheca Gene ID	% coverage	% identity
MEP					
1-deoxy-D-xylulose 5-phosphate synthase 2	DXS	ABI35993.1	Cac_g024944.t1	98	77.60
1-deoxy-D-xylulose-5-phosphate reductoisomerase	DXR	AAF65154.1	Cac_g016318.t1	100	88.82
4-diphosphocytidyl-methylerythritol 2-phosphate synthase	CMS	ACI16377.1	Cac_g018722.t1	88	77.82
4-diphosphocytidyl-2-C-methyl-D-erythritol kinase	CMK	ABI35992.1	Cac_g021688.t1	99	76.17
2C-methyl-D-erythritol 2,4-cyclodiphosphate synthase	MCS	AAF65155.1	Cac_g008169.t1	100	73.77
GCPE protein	HDS	AAO24774.1	Cac_g022763.t1	100	88.65
1-hydroxy-2-methyl-butenyl 4-diphosphate reductase	HDR	ABI30631.1	Cac_g014659.t1	100	83.77
plastid isopentenyl pyrophosphate:dimethylallyl pyrophosphate isomerase	IPI	ABW98669.1	Cac_g008847.t1	76	91.06
geranyl pyrophosphate synthase	GPPS	ACC77966.1	Cac_g026508.t1	51	76.50
Iridoid					
geraniol 8-hydroxylase	G8H	CAC80883.1	Cac_g017987.t1	95	76.71
8-hydroxygeraniol oxidoreductase	GOR	AHK60836.1	Cac_g027560.t1	100	71.69
iridoid synthase	ISY	AFW98981.1	Cac_g006027.t1	100	65.65
iridoid oxidase	IO	AHK60833.1	Cac_g032709.t1	97	78.44
UDP-glucose iridoid glucosyltransferase	7DLGT	BAO01109.1	Cac_g008744.t1	100	77.11
7-deoxyloganic acid 7-hydroxylase	7DLH	AGX93062.1	Cac_g012663.t1	96	69.58
loganic acid methyltransferase	LAMT	ABW38009.1	Cac_g005179.t1	95	53.91
secologanin synthase	SLS	AAA33106.1	Cac_g012666.t1	99	64.94

Note: Only the top hit from the BLAST search is presented.

1
2
3
4 287 **References**
5
6
7

- 8 288 1. Angiosperm Phylogeny Group III. An update of the Angiosperm Phylogeny Group classification
9 289 for the orders and families of flowering plants: APG III. *Botanical Journal of the Linnean Society*.
10 290 2009;161:105-21.
11 291 2. Lorence A, Nessler, C.L. Molecules of Interest: Camptothecin, over four decades of surprising
12 292 findings. *Phytochemistry*. 2004; 65:2735–49.
13 293 3. World Health Organization: 19th WHO Model List of Essential Medicines.
14 294 http://www.who.int/medicines/publications/essentialmedicines/EML2015_8-May-15.pdf.
15 295 Accessed 26 March 2017.
16 296 4. Gongora-Castillo E, Childs KL, Fedewa G, Hamilton JP, Liscombe DK, Magallanes-Lundback M, et
17 297 al. Development of transcriptomic resources for interrogating the biosynthesis of monoterpene
18 298 indole alkaloids in medicinal plant species. *PLoS One*. 2012;7 12:e52506.
19 299 doi:10.1371/journal.pone.0052506.
20 300 5. Sun Y, Luo H, Li Y, Sun C, Song J, Niu Y, et al. Pyrosequencing of the *Camptotheca acuminata*
21 301 transcriptome reveals putative genes involved in camptothecin biosynthesis and transport. *BMC*
22 302 *Genomics*. 2011;12:533. doi:10.1186/1471-2164-12-533.
23 303 6. Yamazaki M, Mochida K, Asano T, Nakabayashi R, Chiba M, Udomson N, et al. Coupling deep
24 304 transcriptome analysis with untargeted metabolic profiling in *Ophiorrhiza pumila* to further the
25 305 understanding of the biosynthesis of the anti-cancer alkaloid camptothecin and anthraquinones.
26 306 *Plant Cell Physiol*. 2013;54 5:686-96. doi:10.1093/pcp/pct040.
27 307 7. Martin M. Cutadapt removes adapter sequences from high-throughput sequencing reads.
28 308 *EMBnetjournal*. 2011;17 1 doi:<http://dx.doi.org/10.14806/ej.17.1.200>.
29 309 8. Haas BJ, Papanicolaou A, Yassour M, Grabherr M, Blood PD, Bowden J, et al. De novo transcript
30 310 sequence reconstruction from RNA-seq using the Trinity platform for reference generation and
31 311 analysis. *Nat Protoc*. 2013;8 8:1494-512. doi:10.1038/nprot.2013.084.
32 312 9. Altschul SF, Madden TL, Schaffer AA, Zhang J, Zhang Z, Miller W, et al. Gapped BLAST and PSI-
33 313 BLAST: a new generation of protein database search programs. *Nucleic Acids Res*. 1997;25
34 314 17:3389-402.
35 315 10. Camacho C, Coulouris G, Avagyan V, Ma N, Papadopoulos J, Bealer K, et al. BLAST+: architecture
36 316 and applications. *BMC Bioinformatics*. 2009;10:421. doi:10.1186/1471-2105-10-421.
37 317 11. FastQC. <http://www.bioinformatics.babraham.ac.uk/projects/fastqc/>. Accessed 26 March 2017.
38 318 12. Saghai-Marroof MA, Soliman KM, Jorgensen RA and Allard RW. Ribosomal DNA spacer-length
39 319 polymorphisms in barley - Mendelian inheritance, chromosomal location, and population-
40 320 dynamics. *PRoc Natl Acad USA*. 1984;81 24:8014-8. doi:Doi 10.1073/Pnas.81.24.8014.
41 321 13. Hardigan MA, Crisovan E, Hamilton JP, Kim J, Laimbeer P, Leisner CP, et al. Genome Reduction
42 322 Uncovers a Large Dispensable Genome and Adaptive Role for Copy Number Variation in
43 323 Asexually Propagated *Solanum tuberosum*. *Plant Cell*. 2016;28 2:388-405.
44 324 doi:10.1105/tpc.15.00538.
45 325 14. Leggett RM, Clavijo BJ, Clissold L, Clark MD and Caccamo M. NextClip: an analysis and read
46 326 preparation tool for Nextera Long Mate Pair libraries. *Bioinformatics*. 2014;30 4:566-8.
47 327 doi:10.1093/bioinformatics/btt702.
48 328 15. Gnerre S, Maccallum I, Przybylski D, Ribeiro FJ, Burton JN, Walker BJ, et al. High-quality draft
49 329 assemblies of mammalian genomes from massively parallel sequence data. *Proc Natl Acad Sci U*
50 330 *S A*. 2011;108 4:1513-8. doi:1017351108 [pii]10.1073/pnas.1017351108.
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4 331 16. Luo R, Liu B, Xie Y, Li Z, Huang W, Yuan J, et al. SOAPdenovo2: an empirically improved memory-
5 332 efficient short-read de novo assembler. *Gigascience*. 2012;1 1:18. doi:10.1186/2047-217X-1-18.
6 333 17. Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV and Zdobnov EM. BUSCO: assessing
7 334 genome assembly and annotation completeness with single-copy orthologs. *Bioinformatics*.
8 335 2015;31 19:3210-2. doi:10.1093/bioinformatics/btv351.
9 336 18. Han Y and Wessler SR. MITE-Hunter: a program for discovering miniature inverted-repeat
10 337 transposable elements from genomic sequences. *Nucleic Acids Res*. 2010;38 22:e199.
11 338 doi:10.1093/nar/gkq862.
12 339 19. Repeat Modeler. <http://www.repeatmasker.org/>. Accessed 26 March 2017.
13 340 20. Campbell MS, Law M, Holt C, Stein JC, Moghe GD, Hufnagel DE, et al. MAKER-P: a tool kit for the
14 341 rapid creation, management, and quality control of plant genome annotations. *Plant Physiol*.
15 342 2014;164 2:513-24. doi:10.1104/pp.113.230144.
16 343 21. RepeatMasker. <http://www.repeatmasker.org/>. Accessed 26 March 2017.
17 344 22. Kim D, Pertea G, Trapnell C, Pimentel H, Kelley R and Salzberg SL. TopHat2: accurate alignment
18 345 of transcriptomes in the presence of insertions, deletions and gene fusions. *Genome Biol*.
19 346 2013;14 4:R36. doi:10.1186/gb-2013-14-4-r36.
20 347 23. Grabherr MG, Haas BJ, Yassour M, Levin JZ, Thompson DA, Amit I, et al. Full-length
21 348 transcriptome assembly from RNA-Seq data without a reference genome. *Nature Biotechnology*.
22 349 2011;29 7:644-52. doi:10.1038/nbt.1883.
23 350 24. Stanke M, Schoffmann O, Morgenstern B and Waack S. Gene prediction in eukaryotes with a
24 351 generalized hidden Markov model that uses hints from external sources. *BMC Bioinformatics*.
25 352 2006;7:62.
26 353 25. Stanke M and Morgenstern B. AUGUSTUS: a web server for gene prediction in eukaryotes that
27 354 allows user-defined constraints. *Nucleic Acids Res*. 2005;33 Web Server issue:W465-7.
28 355 doi:10.1093/nar/gki458.
29 356 26. PASA2. <http://pasapipeline.github.io/>. Accessed 26 March 2017.
30 357 27. Haas BJ, Delcher AL, Mount SM, Wortman JR, Smith RK, Jr., Hannick LI, et al. Improving the
31 358 Arabidopsis genome annotation using maximal transcript alignment assemblies. *Nucleic Acids*
32 359 *Res*. 2003;31 19:5654-66.
33 360 28. Altschul SF and Gish W. Local alignment statistics. *Methods Enzymol*. 1996;266:460-80.
34 361 29. The Arabidopsis Information Resource. Arabidopsis.org. Accessed 26 March 2017.
35 362 30. Eddy SR. Accelerated Profile HMM Searches. *PLoS Comput Biol*. 2011;7 10:e1002195.
36 363 doi:10.1371/journal.pcbi.1002195.
37 364 31. Sirikantaramas S, Yamazaki M and Saito K. Mutations in topoisomerase I as a self-resistance
38 365 mechanism coevolved with the production of the anticancer alkaloid camptothecin in plants.
39 366 *Proc Natl Acad Sci U S A*. 2008;105 18:6782-6. doi:0801038105 [pii]10.1073/pnas.0801038105.
40 367 32. Hanada K, Zou C, Lehti-Shiu MD, Shinozaki K and Shiu SH. Importance of lineage-specific
41 368 expansion of plant tandem duplicates in the adaptive response to environmental stimuli. *Plant*
42 369 *Physiol*. 2008;148 2:993-1003.
43 370 33. Emms DM and Kelly S. OrthoFinder: solving fundamental biases in whole genome comparisons
44 371 dramatically improves orthogroup inference accuracy. *Genome Biol*. 2015;16:157.
45 372 doi:10.1186/s13059-015-0721-2.
46 373 34. Nutzmann HW and Osbourn A. Gene clustering in plant specialized metabolism. *Curr Opin*
47 374 *Biotechnol*. 2014;26:91-9. doi:10.1016/j.copbio.2013.10.009.
48 375 35. Kellner F, Kim J, Clavijo BJ, Hamilton JP, Childs KL, Vaillancourt B, et al. Genome-guided
49 376 investigation of plant natural product biosynthesis. *Plant J*. 2015;82 4:680-92.
50 377 doi:10.1111/tpj.12827.

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

378 36. Itkin M, Heinig U, Tzfadia O, Bhide AJ, Shinde B, Cardenas PD, et al. Biosynthesis of
379 antinutritional alkaloids in solanaceous crops is mediated by clustered genes. *Science*. 2013;341
380 6142:175-9. doi:10.1126/science.1240230.

381 37. Sadre R, Magallanes-Lundback M, Pradhan S, Salim V, Mesberg A, Jones AD, et al. Metabolite
382 Diversity in Alkaloid Biosynthesis: A Multilane (Diastereomer) Highway for Camptothecin
383 Synthesis in *Camptotheca acuminata*. *Plant Cell*. 2016;28 8:1926-44. doi:10.1105/tpc.16.00193.

384 38. DellaPenna D and O'Connor SE. Plant science. Plant gene clusters and opiates. *Science*. 2012;336
385 6089:1648-9. doi:10.1126/science.1225473.

386 39. Chae L, Kim T, Nilo-Poyanco R and Rhee SY. Genomic signatures of specialized metabolism in
387 plants. *Science*. 2014;344 6183:510-3. doi:10.1126/science.1252076.

388 40. Kliebenstein DJ. A role for gene duplication and natural variation of gene expression in the
389 evolution of metabolism. *PLoS One*. 2008;3 3:e1838. doi:10.1371/journal.pone.0001838.

390 41. Kliebenstein DJ, Lambrix VM, Reichelt M, Gershenzon J and Mitchell-Olds T. Gene duplication in
391 the diversification of secondary metabolism: tandem 2-oxoglutarate-dependent dioxygenases
392 control glucosinolate biosynthesis in *Arabidopsis*. *Plant Cell*. 2001;13 3:681-93.

393 42. The Medicinal Plant Genomics Resource. <http://medicinalplantgenomics.msu.edu/>. Accessed 26
394 March 2017.

395

396

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

397 **Additional files**

398 **Supplemental tables:**

399 **Table S1. RNA-sequencing libraries used in this study.**

BioProject ID	BioSample ID	Tissue	No. cleaned	
			reads	Estimated bases
PRJNA80029	SAMN00255206	mature leaf	90,862,580	5,451,754,800
PRJNA80029	SAMN00255207	immature bark	84,537,958	5,072,277,480
PRJNA80029	SAMN00255208	root	88,940,668	5,336,440,080
PRJNA80029	SAMN00255215	young flower	71,435,806	4,286,148,360
PRJNA80029	SAMN00255216	immature fruit	84,250,338	5,055,020,280
PRJNA80029	SAMN00255217	mature fruit	47,811,342	2,868,680,520
PRJNA80029	SAMN00255222	cotyledons	74,037,722	4,442,263,320
PRJNA80029	SAMN00255223	upper stem	76,105,786	4,566,347,160
PRJNA80029	SAMN00255224	lower stem	72,680,940	4,360,856,400
PRJNA361128	SAMN06229771	root	55,435,804	7,224,198,331
Total			771,909,254	49,309,244,481

400

401 **Table S2. Orthologous groups of genes from *Camptotheca acuminata* and five other plant**
402 **species.**

403 This is available as a separate XLS file

404

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

405 **Table S3. P450 paralogous genes in *Camptotheca acuminata*.**

406 This is available as a separate XLS file

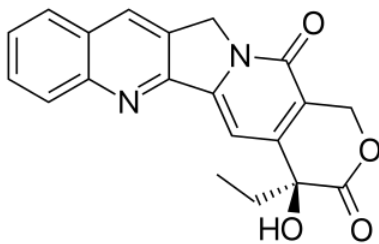
407 **Table S4. Expression abundance matrix (fragments per kbp exon model per million mapped**
408 **reads) from different tissues of *Camptotheca acuminata*.**

409 This is available as a separate XLS file

410



Figure 2



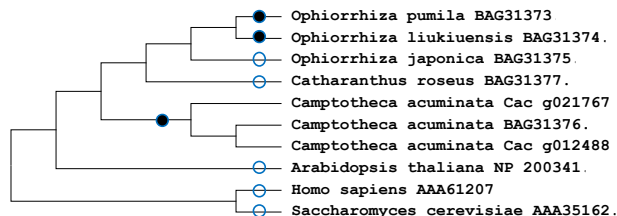
B

Homo sapiens AAA61207	F RG R GN H PK M GM L KRR I MP E DI I NC S DK A K V SP P P-G H K W KE V R H DN K V T W L SV T EN I Q S -S 423
Camptotheca acuminata BAG31376	F RG R GE H PK M GM L K K C I R F SD I T I N I G K DA P IE C PI P GES W KE I R H DN T V T W L AF W ND P IK P RE 556
Camptotheca acuminata Cac_g012488	F RG R GE H PK M GM L K K C I R F SD I T I N I G K DA P IE C PI P GES W KE I R H DN T V T W L AF W ND P IK P RE 555
Camptotheca acuminata Cac_g021767	F RG R GE H PK M GM L K L I R SP D I T I N I G K D A P IE C PI P GES W KE I R H DN T V T W L AF W ND P IN P RE 560
Ophiorrhiza pumila BAG31373	F RG R GE H PK V GM L K K R I R P RD I T I N I G K DA P IE C PI P GER W KE V R ND NT V W L AY W ND P V N L K E 587
Ophiorrhiza liukuensis BAG31374	F RG R GE H PK M GM L K K R I R P RD I T I N I G K DA P IE C PI P GER W KE V R ND NT V W L AF W ND P IN Q K E 587
Ophiorrhiza japonica BAG31375	F RG R GE H PK M GM L K K R I R P RD I T I N I G K DA P IE C PI P GER W KE V R ND NT V W L AF W ID P IN Q K E 588
Catharanthus roseus BAG31377	F RG R GE H PK M GM L K K R I R P CD I T I N I G K DA P IE C VP P GER W KE V R H DN T V T W L AF W ND P IN P K E 570
Arabidopsis thaliana NP_200341	F RG R GE H PK M GM L K K R I HP E IT L N I G K DA P IE C PI A GER W KE V R H DN T V T W L AF W AD P IN P K E 575
Saccharomyces cerevisiae AAA35162	F X G R G A H P T G K L K R R V N P E D I V L N L S K D A P V P A P E -G H K W GE I R H DN T V Q W L A M W R EN F N S -S 355

Direct/indirect camptothecin binding

Homo sapiens AAA61207	E SK K K A V Q R L E E Q L M K L E V Q AT D R E N K Q I A L G T S K I N Y L D P R I T V A W C K 734
Camptotheca acuminata BAG31376	E AL E R K I G Q T NA K I E K M ER D K E T K E G L K T I A L G T S K I S Y L D P R I T V A W C K 864
Camptotheca acuminata Cac_g012488	E AL E R K I G Q T NA K I E K M ER D K E T K E G L K T I A L G T S K I S Y L D P R I T V A W C K 863
Camptotheca acuminata Cac_g021767	E AL G R K I A Q T SA K I E K M ER D K A T K E G L K T V A L S T S K I S Y L D P R I T V A W C K 868
Ophiorrhiza pumila BAG31373	E AL E R K I A Q T NA K I E K M ER D K K T K E D L K A V A L S T S K I S Y L D P R I T V A W C K 896
Ophiorrhiza liukuensis BAG31374	E S L E R K I A Q T N A K I E K M ER D K K T K E D L K A V A L S T S K I S Y L D P R I T V A W C K 896
Ophiorrhiza japonica BAG31375	E AL E R K MA I NA K I E K M ER D K E T K E D L K T V A L G T S K I N Y L D P R I T V A W C K 897
Catharanthus roseus BAG31377	E S L E K K I A Q T N A K I E K M ER D K E T K E D L K T V A L G T S K I N Y L D P R I T V A W C K 880
Arabidopsis thaliana NP_200341	N A W E K I A Q Q S A K I E K M E R D M H T K E D L K T V A L G T S K I N Y L D P R I T V A W C K 883
Saccharomyces cerevisiae AAA35162	E K I A Q V E K L Q R I Q T S S I Q L K D K E EN S Q V S L G T S K I N Y I D P R L S V P F C K 738

Direct/indirect camptothecin binding



camptothecin

● present

○ absent

Figure 3

[Click here to download Figure Fig_3.pptx](#)

