



A physical and genetic map of *Cannabis sativa* identifies extensive rearrangement at the THC/CBD acid synthase locus

Kaitlin U Laverty, Jake M Stout, Mitchell J Sullivan, et al.

Genome Res. published online November 8, 2018

Access the most recent version at doi:[10.1101/gr.242594.118](https://doi.org/10.1101/gr.242594.118)

P<P	Published online November 8, 2018 in advance of the print journal.
Accepted Manuscript	Peer-reviewed and accepted for publication but not copyedited or typeset; accepted manuscript is likely to differ from the final, published version.
Open Access	Freely available online through the <i>Genome Research</i> Open Access option.
Creative Commons License	This manuscript is Open Access. This article, published in <i>Genome Research</i> , is available under a Creative Commons License (Attribution-NonCommercial 4.0 International license), as described at http://creativecommons.org/licenses/by-nc/4.0/ .
Email Alerting Service	Receive free email alerts when new articles cite this article - sign up in the box at the top right corner of the article or click here .



To subscribe to *Genome Research* go to:
<https://genome.cshlp.org/subscriptions>

Published by Cold Spring Harbor Laboratory Press

1 **A physical and genetic map of *Cannabis sativa* identifies extensive rearrangement at the**
2 **THC/CBD acid synthase locus**

3 Kaitlin U. Lavery¹, Jake M. Stout², Mitchell J. Sullivan³, Hardik Shah^{3,4}, Navdeep Gill⁵, Larry
4 Holbrook⁶, Gintaras Deikus^{3,4}, Robert Sebra^{3,4}, Timothy R. Hughes^{1,7,8,*}, Jonathan E. Page^{5,9,*},
5 and Harm van Bakel^{1,3,4,*}

6 ¹Department of Molecular Genetics, University of Toronto, Toronto, ON, Canada, M5S 1A8

7 ²Department of Biological Sciences, University of Manitoba, Winnipeg, MB, Canada, R3T 2N2

8 ³Department of Genetics and Genomic Sciences, Icahn School of Medicine at Mount Sinai, New
9 York City, NY, USA, 10029

10 ⁴Icahn Institute for Data Science and Genomic Technology, Icahn School of Medicine at Mount
11 Sinai, New York City, NY, USA, 10029

12 ⁵Department of Botany, University of British Columbia, Vancouver, BC, Canada, V6T 1Z4

13 ⁶CanniMed Therapeutics Inc., Saskatoon, SK, Canada, S7K 3J8

14 ⁷Donnelly Centre, University of Toronto, Toronto, ON, Canada, M5S 3E1

15 ⁸Canadian Institute for Advanced Research, MaRS Centre, West Tower, 661 University Avenue,
16 Suite 505, Toronto, ON, Canada, M5G 1M1

17 ⁹Anandia Labs, Vancouver, BC, Canada, V6T 1Z4

18

19 *To whom correspondence should be addressed:

20 t.hughes@utoronto.ca, harm.vanbakel@mssm.edu, jon.page@botany.ubc.ca

21

22 **KEYWORDS**

23 Cannabis genome; cannabinoids; genome assembly; genetic map; tetrahydrocannabinol;
24 cannabidiol, cannabichromene

25 ABSTRACT

26 ***Cannabis sativa* is widely cultivated for medicinal, food, industrial, and recreational use,**
27 **but much remains unknown regarding its genetics, including the molecular determinants**
28 **of cannabinoid content. Here, we describe a combined physical and genetic map derived**
29 **from a cross between the drug-type strain ‘Purple Kush’ and the hemp variety ‘Finola’.**
30 **The map reveals that cannabinoid biosynthesis genes are generally unlinked, but that**
31 **aromatic prenyltransferase (AP) – which produces the substrate for THCA and CBDA**
32 **synthases (THCAS and CBDAS) – is tightly linked to a known marker for total**
33 **cannabinoid content. We further identify the gene encoding CBCA synthase (CBCAS)**
34 **and characterize its catalytic activity, providing insight into how cannabinoid diversity**
35 **arises in cannabis. Strikingly, *THCAS* and *CBDAS* (which determine the drug vs hemp**
36 **chemotype) are contained within large (>250 kb) retrotransposon-rich regions that are**
37 **highly non-homologous between drug- and hemp-type alleles, and are furthermore**
38 **embedded within ~40 Mb of non-recombining repetitive DNA. The chromosome**
39 **structures are similar to those in grains such as wheat, with recombination focused in**
40 **gene-rich, repeat-depleted regions near chromosome ends. The physical and genetic**
41 **map should facilitate further dissection of genetic and molecular mechanisms in this**
42 **commercially and medically important plant.**

43

44 INTRODUCTION

45 Domesticated thousands of years ago (Li 1974), *C. sativa* has been subjected to intensive
46 breeding, resulting in extensive variation in morphology and chemical composition. It is perhaps
47 best known for producing cannabinoids, a unique class of compounds that may function in
48 chemical defense (Pate 1994), but also have pharmaceutical and psychoactive properties. Heat
49 converts the cannabinoid acids (e.g. tetrahydrocannabinolic acid, THCA) to neutral molecules
50 (e.g. (-)-*trans*- Δ^9 -tetrahydrocannabinol, THC) that bind to endocannabinoid receptors found in

51 the vertebrate nervous system. This pharmacological activity leads to analgesic, antiemetic, and
52 appetite-stimulating effects and may alleviate symptoms of neurological disorders including
53 epilepsy (Devinsky et al. 2014) and multiple sclerosis (van Amerongen et al. 2017). There are
54 over 113 known cannabinoids (Elsohly and Slade 2005), but the two most abundant natural
55 derivatives are THC and cannabidiol (CBD). THC is responsible for the well-known
56 psychoactive effects of cannabis consumption, but CBD, while non-intoxicating, also has
57 therapeutic properties, and is specifically being investigated as a treatment for both
58 schizophrenia (Osborne et al. 2017) and Alzheimer's disease (Watt and Karl 2017). Cannabis
59 has traditionally been classified as having a drug ("marijuana") or hemp chemotype based on
60 the relative proportion of THC to CBD, but types grown for psychoactive use produce relatively
61 large amounts of both. Cannabis containing high levels of CBD is increasingly grown for medical
62 use.

63
64 THCA and CBDA are both synthesized from cannabigerolic acid by the related enzymes THCA
65 synthase (THCAS) and CBDA synthase (CBDAS), respectively (Sirikantaramas et al. 2004;
66 Taura et al. 2007). Expression of THCAS and CBDAS appear to be the major factor
67 determining cannabinoid content, but the mechanisms that underlie the expression of these
68 enzymes remain unresolved. Two competing theories are supported by existing data. In one,
69 *CBDAS* and *THCAS* are mutually exclusive alleles (i.e. very different isoforms, as the protein
70 sequences are only 84% identical). Genetic analysis supports this model, with approximately
71 1:2:1 segregation of chemotypes in a cross of drug-type vs hemp (de Meijer et al. 2003). An
72 alternative model is that *THCAS* and *CBDAS* are closely linked (i.e. adjacent on a
73 chromosome), and one or the other is inactivated in drug-type or hemp strains. This model was
74 motivated by the discovery of a *THCAS*-like gene in hemp plants (Kojoma et al. 2006), and is
75 consistent with the possibility that these related genes are derived from an ancient tandem
76 duplication. In addition, physical linkage of genes involved in specialized metabolic pathways

77 has been repeatedly observed in plants, similar to operons in bacterial genomes (Nutzmann and
78 Osbourn 2014); such a cluster was recently described for benzylisoquinoline alkaloid
79 biosynthesis genes in opium poppy (Guo et al. 2018). It is unknown whether genes involved in
80 cannabinoid biosynthesis are clustered, although genetic analyses have previously indicated
81 that at least one locus unlinked to *THCAS/CBDAS* contributes to cannabinoid content (Weiblen
82 et al. 2015).

83
84 The draft genome and transcriptome of *C. sativa* described in 2011 (van Bakel et al. 2011) (for a
85 female plant of the drug-type strain Purple Kush (PK), and resequencing of a plant of the hemp
86 variety 'Finola' (FN)) was unable to discriminate between these models, due to high
87 fragmentation. The *C. sativa* draft genome assembly, done largely with Illumina sequencing,
88 was composed of 136,290 scaffolds, with an N50 of 16.2 Kb. It was subsequently demonstrated
89 that ~70% of the *C. sativa* draft genome is composed of repeat sequence (Pisupati et al. 2018).
90 Measurement of Single Nucleotide Variants (SNVs) in four strains showed rates of
91 heterozygosity ranging from 0.18 - 0.26% and revealed that drug-type and hemp-type strains
92 were well separated by SNVs, the rate of occurrence of SNVs between these types was as high
93 as 0.64% (van Bakel et al. 2011). Cytogenetic analysis has furthermore suggested a high
94 degree of inter- and intracultivar karyotype polymorphisms (i.e. differences in homologous
95 chromosomes that can be observed by microscopy), at least among hemp varieties (Razumova
96 et al. 2016), which may further complicate genome assembly. To address these complications,
97 and to simultaneously leverage the high rate of SNVs between PK and FN, we coupled Pacific
98 Biosciences (PacBio) long-read single-molecule real-time (SMRT) sequencing of PK and FN
99 with Illumina resequencing of 99 F1 progeny between the two, in order to generate a combined
100 genetic and physical map. The combined map provides new insight into the arrangement of the
101 chromosomes and the cannabinoid biosynthetic genes, including discovery of substantial

102 rearrangement and gene duplications at the closely linked THC and CBD acid synthase gene
103 loci.

104

105 **RESULTS**

106 **A combined genetic and physical map reveals that genes and recombination events are**
107 **concentrated near chromosome ends.**

108 We performed PacBio SMRT sequencing of genomic DNA from the female parent PK and the
109 male parent FN to a depth of ~79x and ~98x, respectively. We used these data to develop an
110 initial set of scaffolds, using the FALCON assembler (Chin et al. 2016), with PK and FN
111 analyzed separately (**Table 1**). The assemblies were further polished with Illumina data using
112 Pilon (Walker et al. 2014) to correct indel errors associated with homopolymer repeats in PacBio
113 data. The FN assembly was more contiguous than the PK assembly (scaffold N50 of 445.6 vs.
114 146 Kbp, respectively) – likely reflecting the increased FN coverage and the use of a more
115 recent sequencing chemistry – and each substantially improved on our original Illumina
116 assembly (van Bakel et al. 2011) (**Supplemental Fig. 1**). *De novo* repeat classification using
117 RepeatModeler (<http://www.repeatmasker.org/RepeatModeler/>) confirmed that the sequence of
118 both assemblies is highly repetitive (~73%) (**Supplemental Fig. 2**), with hundreds of distinct
119 families. The two sets of scaffolds largely mapped to each other 1-to-1 (**Supplemental Fig. 3**),
120 but with differing breakpoints that mostly reflected differences in scaffold boundaries. The total
121 size of the PK and FN assemblies was close to the haploid genome size estimated by flow
122 cytometry (818 Mb and 843 Mb for female and male, respectively (Sakamoto 1998)). Overall,
123 90.3% of 30,074 previously-described PK transcripts (van Bakel et al. 2011) mapped to the PK
124 assembly (82.3% mapping completely within a single scaffold). Each assembly also contained
125 >95% of eudicotyledon single-copy orthologs from OrthoDB, of which >97% were complete
126 (**Supplemental Fig. 4**), indicating that both assemblies represented the vast majority of the
127 cannabis gene space. An ortholog duplication rate >14% and slightly larger than expected

128 assembly sizes suggest that some regions of the diploid genomes were resolved into separate
 129 contigs, which can be an issue for polymorphic species (Shimizu et al. 2017).

130

131 **Table 1. Genome assembly statistics**

Assembly statistic	PK	FN	FN anchored
PacBio sequencing and assembly			
Total PacBio raw reads	9,979,332	10,623,051	N/A
Total PacBio raw read bases (Gbp)	64.62	82.32	N/A
Average PacBio raw read length (bp)	7,179	8,716	N/A
Total assembled bases (Mbp)	892	1,009	784
Scaffold N50 (Kbp)	146.0	445.6	382.9
Number of scaffolds	12,887	5,304	2,952
Largest scaffold (Mbp)	1.41	2.49	2.49
% PK transcriptome in genome ($\geq 50\%$ match)	90.3%	87.3%	78.5%
% PK transcriptome in genome (complete)	82.7%	78.5%	70.4%
% repeat content	73.3%	73.9%	72.2%
Haplotype blocks			
FN haplotype blocks with >10 SNVs	34,197	13,098	10,557*
Number of phased SNVs in haplotype blocks	2,734,893	1,359,019	1,214,845*
% FN scaffolds with ≥ 1 haplotype block	77.2%	86.7%	100%
% sequence in FN haplotype blocks	43.5%	76.1%	78.8%*
FN haplotype block N50 (Kbp)	27.6	92.6	98.7*
Number of blocks used to create genetic map	14,440	4,507	N/A
Number of SNVs used to create genetic map	1,888,187	799,227	N/A
Mean total coverage at SNVs used to create genetic map (parents & F1s)	718 (+/- 350)	660 (+/- 363)	N/A
Illumina sequencing			
Total Illumina raw paired-end reads	105,000,000	162,968,810	N/A
Total filtered Illumina paired-end reads	80,369,366	98,244,687	N/A
Total filtered Illumina reads bases (Gbp)	11.49	28.98	N/A
Coverage of FN FALCON assembly	14.1x	28.5x	N/A

132 *Statistics are based on the subset of FN haplotype blocks that are contained within scaffolds in the
 133 anchored map.

134

135 We reasoned that a genetic map would provide an independent means to link scaffolds, in
 136 addition to being independently useful for genetic analysis. To generate a genetic linkage map,
 137 we employed the SOILOCO pipeline, created by Scaglione *et al.* (Scaglione et al. 2016) to
 138 create a map of the artichoke genome. We applied the pipeline to F1 data from a cross between
 139 a PK female and FN male. SOILOCO requires phasing of the parental scaffolds into blocks in

140 which parental haplotypes can be uniquely identified. It then uses SNVs in the offspring to
141 determine which of the parental haplotypes is inherited for each F1 at each block. The inherited
142 parental haplotypes are called using a Hidden Markov Model, which compensates for
143 uncertainty in genotype calling caused by relatively low coverage typical of resequencing, by
144 taking advantage of the multiple SNVs in each block. Because each of the four parental
145 haplotypes is traced uniquely, recombination frequencies between blocks (and thus between
146 scaffolds) can subsequently be calculated, and the recombination frequencies can be used to
147 place blocks (and scaffolds) into linkage groups. Since the blocks of informative SNVs differ
148 between the parental types, a separate genetic map is created for each parent (in this case, PK
149 and FN). In our implementation, we identified phased haplotype blocks of physically linked
150 unique SNVs in the FN assembly contained within PK or FN PacBio raw reads using HapCUT2
151 (Edge et al. 2017) (**Table 1**), and scored them in 99 F1 progeny using Illumina sequencing
152 (median coverage ~4x). We then ran the SOILOCO pipeline, followed by R/qtl (Broman et al.
153 2003) and MSTmap (Wu et al. 2008) to form linkage groups and order scaffolds within them.

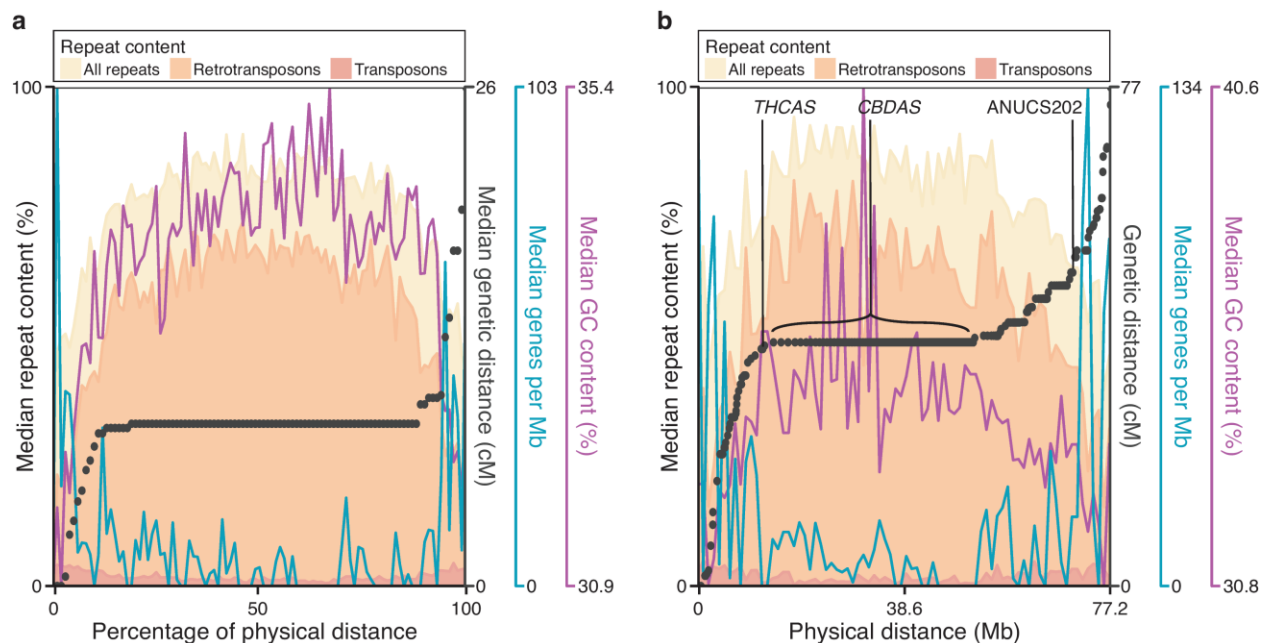
154
155 The blocks formed ten large linkage groups in both PK and FN, which we assume correspond to
156 the established nine autosomes and X/Y (which contain a pseudoautosomal region and
157 recombine (Peil et al. 2003)) and are hereafter referred to as chromosomes. The maps were
158 largely consistent between PK and FN (**Supplemental Fig. 5**), and were therefore merged
159 (MergeMap) (Wu et al. 2011). The merged genetic map is depleted for short scaffolds, repetitive
160 sequence, and scaffolds containing a higher proportion of SNVs with segregation distortion
161 (these SNVs are ignored by SOILOCO). The merged map contains 2,952/5,304 scaffolds,
162 784/1,006 Mb (78%) of the initial sequence, 89% of Eudicotyledon single-copy orthologs and
163 21,168/30,074 of all PK transcripts (70.4%) (**Table 1, Supplemental Fig. 4**).

164

165 **Fig. 1a** plots composite physical vs genetic distance across the chromosomes, with several
166 major trends in the chromosomal sequences also illustrated (**Supplemental Fig. 5** shows
167 similar graphs, and also plots of genetic vs. physical distance as well as a comparison of
168 recombination frequencies, for all individual chromosomes). First, there is a very strong
169 tendency for recombination to occur near chromosome ends, while there are typically large
170 blocks lacking recombination events across the middle of the chromosome. Second, genes are
171 much more frequent near chromosome ends. Because promoters and enhancers are typified by
172 open chromatin, which appears to promote crossovers in diverse species including maize (Liu et
173 al. 2009) and *A. thaliana* (Choi et al. 2013), this arrangement may underlie the observed
174 recombination frequencies. Third, the poorly recombining central parts of chromosomes are not
175 only gene-poor, but also have a higher repeat content, which may be methylated and could
176 suppress recombination (Zamudio et al. 2015). Fourth, assuming that the centromere is located
177 within the non-recombining central segments of the chromosomes, then chromosomes 5, 9, and
178 10 appear to be telocentric (i.e. behave as if they have a single long arm). These may represent
179 the sex chromosome, one end of which is non-homologous and thus non-recombining, and
180 chromosomes 8 and 9 (as determined by cytogenetics (Divashuk et al. 2014)), which harbor 5S
181 rDNA and 45S rDNA on one arm, respectively. The repetitive nature of these regions would be
182 expected to impede both assembly and mapping. Indeed, four of five male-specific markers are
183 found in the Finola assembly, but none were placed on the genetic map, and the 45S and 5S
184 rDNA are not in the assembly (**Supplemental Table 1**).

185
186 Overall, the organization of *C. sativa* genes, repeats, and recombination frequency along
187 chromosomes is similar to what is commonly observed in the grains (e.g. maize, barley, and
188 wheat) (Gore et al. 2009; Liu et al. 2009; Mascher et al. 2017). To our knowledge, such an
189 organization is unusual outside the grains: it has been observed in walnut (Luo et al. 2015), but
190 not thale cress (*Arabidopsis thaliana*) (Meinke et al. 2009), apple (Di Pierro et al. 2016),

191 strawberry (Davik et al. 2015), or mulberry (He et al. 2013), suggesting that this property is rare
 192 among Rosales.



193 **Figure 1. Comparison of physical and genetic distance in *C. sativa* and arrangement of**
 194 **sequence features on chromosomes. (a)** Median values are indicated for all metacentric
 195 linkage groups (chromosomes 5, 9, and 10 are excluded), scaled to the same physical length.
 196 Black points indicate the median increase in genetic distance every 1/100th of the physical
 197 distance. Shaded histograms superimposed show density of repeat sequences. Density of
 198 genes and GC content are also indicated by blue and purple lines. **(b)** Values for Chromosome
 199 6, which contains the *THCAS/CBDAS* loci, here black points are the representative of individual
 200 scaffolds.
 201
 202
 203

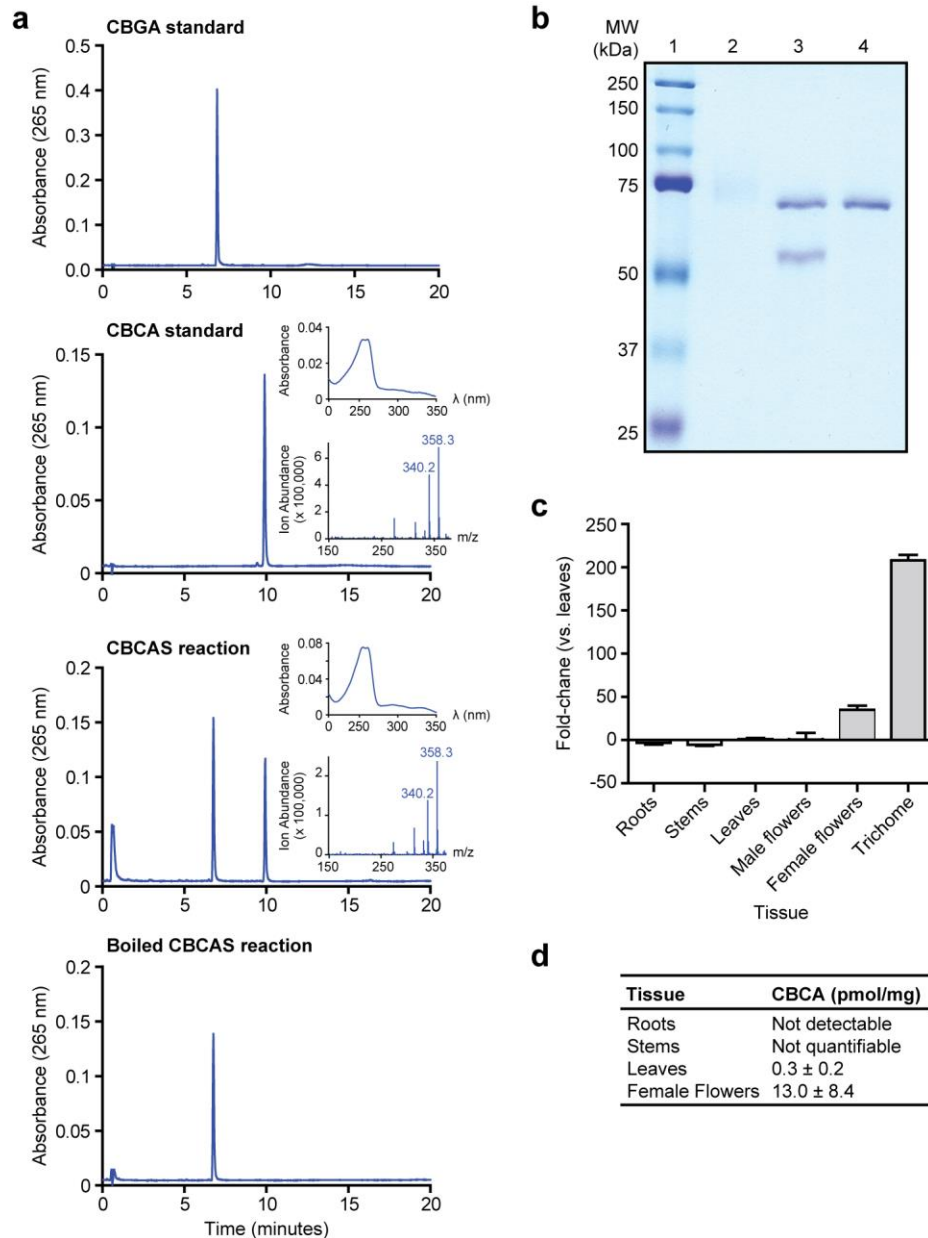
204 Genomic organization of cannabinoid pathway genes.

205 We next examined the positions of genes encoding known cannabinoid biosynthetic enzymes
 206 on the chromosomes. With the exception of the functional copies of *CBDAS* and *THCAS*, which
 207 are considered below, the cannabinoid-related genes are distributed in a mostly random fashion
 208 across the genome (indicated in **Supplemental Fig. 5**). The new map also finds that *C. sativa*
 209 encodes one copy of *AAE1* (hexanoyl-CoA synthetase) and two tandem copies of tetraketide
 210 synthase (“olivetol synthase”). The genome sequences of both PK and FN also contain the
 211 *THCAS*-like gene described by Kojoma (Kojoma et al. 2006) which led to the two-locus
 212 *THCAS/CBDAS* hypothesis. This *THCAS*-like gene is 96% identical to *THCAS* at the nucleotide
 213

214 level, and encodes a protein that is 93% identical to THCAS at the amino acid level. One copy
215 of the *THCAS*-like gene is found in the PK assembly (scaffold 005500F:2986-4620) and two are
216 found in the FN assembly (scaffold 004887F: 13943-15577 and 001793F:69162-70796).

217

218 We examined the possibility that this *THCAS*-like gene encoded cannabichromenic acid (CBCA)
219 synthase (CBCAS), which is found in both drug-type and hemp strains and resembles the THCA
220 and CBDA synthases in its catalytic mechanism (Morimoto et al. 1997). We expressed the
221 predicted open reading frame as a secreted protein in *Pichia pastoris* strain X33. We then
222 added cannabigerolic acid (CBGA) substrate to clarified culture media to test for enzyme
223 activity. The products of this reaction were analyzed by High Performance Liquid
224 Chromatography (HPLC), which revealed a specific signal for CBCA (**Fig. 2a**). Purification of the
225 *Pichia* secreted protein through a series of chromatographic steps yielded a 59 kDa product at
226 the expected size of CBCAS without its secretory signal sequence (calculated to be 58.9 kDa)
227 (**Fig. 2b**). We next determined the kinetic properties of CBCAS after optimizing reaction
228 conditions using the purified protein (**Supplemental Fig. 6**). At the optimal temperature of 40 °C
229 and pH of 5.5, the reaction followed Michaelis Menten reaction kinetics with a K_m of 9.3 ± 2.3
230 μM and a k_{cat} of 0.02 sec^{-1} . These values are similar to those reported for CBCAS purified from
231 cannabis floral tissue ($K_m = 23 \mu\text{M}$, $k_{\text{cat}} = 0.04 \text{ sec}^{-1}$) (Morimoto et al. 1998). Finally, the
232 accumulation of CBCA correlates well with the expression of *CBCAS* in various cannabis
233 tissues, with the highest concentration observed in female floral tissue and minimal amounts in
234 leaf, stem, and root (**Fig. 2c**). Taken together, these data confirm that we identified the gene
235 encoding CBCA synthase.

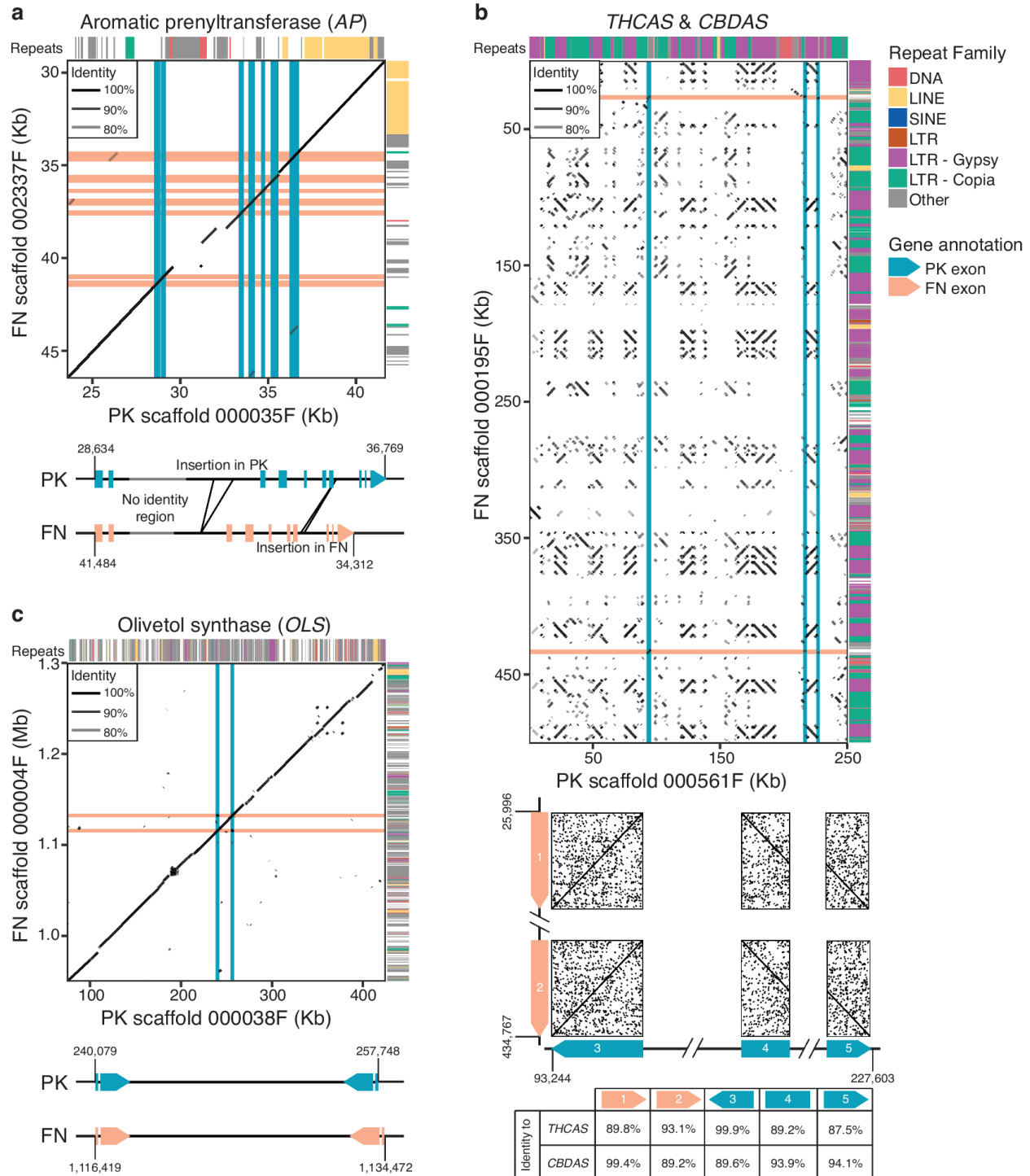


236

237 **Figure 2. Characterization of CBCAS activity and expression. (a)** HPLC analysis of CBCAS
 238 activity detected in *Pichia pastoris* cell cultures. Chromatograms of the CBGA substrate and
 239 CBCA standards are shown together with chromatograms of the enzyme reaction in media
 240 sampled from *Pichia* expressing CBCA in the presence of CBGA substrate, before and after
 241 boiling at 95 °C for 10 min. Insets correspond to the UV-absorbance spectrum (top) and the
 242 mass spectrum derived from a single quadrupole mass spectrometer (bottom) of the compound
 243 that eluted at 10 min. **(b)** SDS-PAGE analysis of CBCAS expressed in *Pichia pastoris* and
 244 purified by protein chromatography. Lane 1: Protein ladder. Lane 2: Concentrated protein
 245 fraction exhibiting CBCAS activity. The high molecular weight smear is glycosylated CBCAS.
 246 Lane 3: Same fraction as lane 2, treated with EndoHf (MW = 70 kDa). Lane 4: EndoHf only. **(c)**
 247 qRT-PCR analysis of *CBCAS* expression in cannabis tissues. cDNA derived from cannabis

248 tissues was used as a template for PCR reactions using *CBCAS*-specific primers and *EF1 α* as a
249 reference gene. Differential expression of *CBCAS* is depicted as fold-change between tissue
250 types compared to leaves. Trichome tissue consisted of isolated trichome secretory cells. **(d)**
251 Quantification of CBCA content of the developing seedlings by HPLC (bottom).
252

253 A previous study (Weiblen et al. 2015) used QTL analysis in *C. sativa* to associate 121 genetic
254 markers with total cannabinoid content and THCA/CBDA ratio. Outside of *THCAS/CBDAS*, this
255 study identified only one locus displaying a strong association with total cannabinoid content, at
256 a distance of ~1.2 cM between the trait and the marker. In our genetic map, this locus (marker
257 ANUCS501) is linked to aromatic prenyltransferase (AP), which catalyzes the production of
258 CBGA, the substrate of *THCAS*, *CBDAS* and *CBCAS*, with a similar recombination frequency
259 (2.1 cM in PK; 4cM in FN). This observation suggests that either polymorphisms or differential
260 regulation of AP contributes to cannabinoid production, presumably by controlling substrate
261 concentration for *THCAS* and *CBDAS*. PK has >5-fold higher transcript levels of AP than FN
262 (van Bakel et al. 2011), with no difference in copy number, suggesting that AP enzyme levels
263 may be higher in drug-type plants partly due to differences in transcript levels. In addition to
264 polymorphisms, there are multiple large (>100 bp) indels in and around the AP locus (including
265 two within introns), which correspond mainly to LTRs, LINEs, and simple-repeat-like insertions,
266 which could conceivably alter regulation of transcription or splicing (**Fig. 3a**).



267
268
269
270
271
272
273
274
275

Figure 3. Comparison of scaffolds between PK and FN assemblies. Alignment of scaffolds from PK and FN FALCON assemblies containing key cannabinoid biosynthesis enzymes are shown. Locations of exons are indicated by pink and blue lines for FN and PK, respectively. Repeat classes given are from RepeatModeler. Individual repeat types indicated were identified by manual analysis. Features of genes are further described and compared beneath alignments. **(a)** Aromatic prenyltransferase (*AP*). **(b)** *THCAS* and *CBDAS*. **(c)** Olivetol synthase (*OLS*, or tetraketide synthase).

276
277 **Extensive rearrangement of the cannabinoid synthase locus underlies chemotype**
278 **differences between Purple Kush and Finola.**

279 Finally, we examined *THCAS* and *CBDAS* in the PK and FN genomes. The PK assembly
280 contains only a single copy of *THCAS*, and no exact copies of *CBDAS*: none have greater than
281 95% identity to *CBDAS* at the nucleotide level. Similarly, the FN assembly contains only a single
282 functional copy of *CBDAS*, while no *THCAS* gene is detected. These observations are
283 confirmed by raw sequencing reads; no reads from FN map to *THCAS*, and no reads from PK
284 map to *CBDAS*. Both genomes include the aforementioned *CBCAS*. This supports claims made
285 using the draft genome and transcriptome (van Bakel et al. 2011). As expected from established
286 segregation patterns, *THCAS* and *CBDAS* map to roughly the same region on Chromosome 6
287 (**Fig. 1b**). Remarkably, however, the scaffolds that contain *THCAS* (in PK) and *CBDAS* (in FN)
288 are dramatically different from each other, and neither has a clear counterpart in the other
289 genome. The scaffold containing *THCAS* in PK does, however, contain a pseudogenic copy of
290 *CBDAS*, with ~94% identity to the known *CBDAS* sequence. The gene is likely non-functional
291 as it has a Gypsy element insertion at its center. Assuming these loci share common ancestry,
292 there has clearly been extensive rearrangement since their divergence. The scaffold containing
293 *THCAS* is ~250 Kb, and that containing *CBDAS* is ~750 Kb, but the dotplot shown in **Fig. 3b**
294 illustrates almost complete lack of similarity over this span, with the exception of a large number
295 of LTR-class retroelements. The extreme rearrangement clearly shows that these two genes do
296 not have a simple isogenic relationship; **Fig. 3a and 3c** illustrate more typical patterns of
297 sequence similarity between PK and FN. Intriguingly, the scaffolds containing *CBDAS* and
298 *THCAS* are both located within a much larger repeat-rich and gene-poor region of ~40 Mb, in
299 the central section of Chromosome 6, encompassing 152 scaffolds with no recombination in
300 either parent observed among the 99 F1s (**Fig. 1b**). The scaffold containing *THCAS*, however,
301 was separated from this region in a single recombination event among the 99 crosses, thus

302 placing it at one end of this region, and indicating that the *THCAS* and *CBDAS* scaffolds are at
303 separate loci. We suggest that this repeat-rich segment of the chromosome may have hosted a
304 series of tandem duplications and rearrangements amplifying an ancestral gene, leading to the
305 present chromosomal organization; there is also a pseudogene with 89-93% identity to each of
306 *THCAS*, *CBDAS*, and *CBCAS* in this region. We note that this observation represents a
307 modification of both previous models of *CBDAS* and *THCAS* arrangement: they are not isoforms
308 at an otherwise equivalent locus, and no equivalent of *THCAS* (deactivated or not) is found in
309 hemp.

310

311 **DISCUSSION**

312 The combined sequence/genetic map presented here is consistent with the known *C. sativa*
313 karyotype and genome size, contains the vast majority of known transcripts, and largely
314 correlates between PK and FN. To completely finish the sequence, it will most likely be
315 necessary to further improve the resolution of the genetic map and/or leverage hybrid
316 scaffolding technologies; e.g. by incorporating single-molecule genomic maps (Pendleton et al.
317 2015) or Hi-C data that provides >1Mb phasing information (Kronenberg et al. 2018). Another
318 future goal will be to identify and fully assemble the X/Y chromosomes. There are numerous
319 scaffolds in both PK and FN with no obvious counterpart in the other genome, which could
320 represent distinctive components of the sex chromosomes, and which were not captured in our
321 genetic map.

322

323 The identification of *CBCAS* allows for a number of potential applications. Cannabichromene
324 (CBC) is a weaker agonist of the cannabinoid CB1 and CB2 receptors compared to THC and
325 CBD. However, unlike THC, both CBD and CBC have been shown to decrease nociception by
326 both blocking the activity of ankyrin-type transient receptor potential channels that play roles in
327 the perception of pain-inducing signals, and by inhibiting the re-uptake of endocannabinoids

328 such as anandamide (Maione et al. 2011). Furthermore, CBC operates as a gastrointestinal
329 anti-inflammatory agent in mice, and protects adult neuronal stem progenitor cells *in vitro* (Izzo
330 et al. 2012; Shinjyo and Di Marzo 2013). It therefore may be useful to breed medical cannabis
331 strains with higher quantities of CBCA to treat specific ailments such as inflammatory bowel
332 disease and Crohn's disease. Finally, the high degree of sequence similarity between *CBCAS*,
333 *THCAS*, and *CBDAS* and the presence of multiple pseudogenes suggest that gene duplication
334 and divergence has been the key driver of cannabinoid end-product diversification in cannabis.
335 Comparative sequence analysis of the enzymes will help ascertain which amino acids are
336 important in catalysis, and may lead to the rational design of cannabinoid biosynthetic enzymes
337 that produce novel cannabinoids not observed in nature.

338

339 Our identification of *CBCAS* also clarifies a puzzling finding of Kojoma et al (2006), who used
340 PCR to amplify a *THCAS*-like gene from "fibre-type" (hemp) cannabis that contained no THCA.
341 Based on the sequence of the gene that we show has *CBCAS* activity, the THCA-like gene
342 amplified by Kojoma et al (2006) is *CBCAS*. This result makes sense, since non-drug/hemp
343 forms of cannabis also contain CBCA.

344

345 Cannabis and cannabinoids are increasingly employed in medicine, and recently have been
346 legalized for recreational use in many jurisdictions. The new map should facilitate vastly
347 improved genetic analysis, including QTL mapping, which will accelerate crop improvement
348 efforts. Drug prohibition has restricted access to cannabis by plant breeders and researchers,
349 and as a result it has received less attention than other crops. Cannabis suffers from insect
350 pests, widespread fungal diseases and has a number of agronomic issues such as flowering
351 time requirements that make it difficult to grow in some environments. In addition, breeding of
352 cannabis types with specific cannabinoid and terpene profiles is desirable for the development
353 of new varieties for medical and recreational use. The fact that a strong and interpretable result

354 was obtained by re-examining a previously described marker correlating with total cannabinoid
355 content (Weiblen et al. 2015) clearly shows the potential of this approach as it applies to
356 cannabinoid metabolism. Due to the relatively high rate of polymorphism in cannabis, it should
357 be possible to employ resequencing (e.g. low-coverage short-read Illumina protocols) on either
358 crosses or at a population level to associate variants or variation with traits and genes, using the
359 genetic map.

360

361 **METHODS**

362

363 **Plant cultivation and genomic DNA isolation**

364 A female PK plant, produced through multiple vegetative propagation generations from the
365 original source plant used to produce the draft *Cannabis sativa* genome (van Bakel et al. 2011),
366 was pollinated by a male FN plant in an indoor growth chamber. Seeds produced from this
367 cross were germinated under standard conditions and grown to seedling stage. Genomic DNA
368 was isolated from young leaves using a GenElute genomic miniprep kit (Sigma). The secure
369 facilities used for plant growing were licensed by Health Canada.

370

371 **PacBio SMRT sequencing of the PK and FN genomes.**

372 Genomic DNA (gDNA) library preparation and sequencing was performed according to the
373 manufacturer's instructions and reflects the P6-C4 sequencing enzyme and chemistry,
374 respectively. PK and FN gDNA was first re-purified using a 0.8X AMPure XP purification step
375 (0.80X AMPure beads added, by volume, to each DNA sample dissolved in 200 μ L EB,
376 vortexed for 10 minutes at 2,000 rpm, followed by two washes with 70% alcohol and finally
377 diluted in EB), to remove small fragments and/or biological contaminant. The purified DNA
378 sample was taken into DNA damage and end-repair steps. Briefly, the DNA fragments were
379 repaired using DNA Damage Repair solution (1X DNA Damage Repeat Buffer, 1X NAD⁺, 1 mM
380 ATP high, 0.1 mM dNTP, and 1X DNA Damage Repeat Mix) with a volume of 21.1 μ L and
381 incubated at 37°C for 20 minutes. DNA ends were repaired next by adding 1X End Repair Mix
382 to the solution, which was incubated at 25°C for 5 minutes, followed by the second 0.45X
383 Ampure XP purification step. Next, 0.75 μ M of Blunt Adapter was added to the DNA, followed by
384 1X template Prep Buffer, 0.05 mM ATP low and 0.75 U/ μ L T4 ligase to ligate (final volume of
385 47.5 μ L) the SMRTbell adapters to the DNA fragments. This solution was incubated at 25°C
386 overnight, followed by a 65°C 10-minute ligase denaturation step. After ligation, the library was

387 treated with an exonuclease cocktail to remove un-ligated DNA fragments using a solution of
 388 1.81 U/ μ L Exo III 18 and 0.18 U/ μ L Exo VII, then incubated at 37°C for 1 hour. Two additional
 389 0.80X Ampure XP purifications steps were performed to remove < 1000 bp molecular weight
 390 DNA and organic contaminant.

391 Size-selection was confirmed using the Agilent Bioanalyzer and the mass was quantified using
 392 a Qubit assay, before proceeding with primer annealing and DNA sequencing. For PK, 100pM
 393 of SMRTbell libraries were mag bead loaded and sequenced with a combination of P5/C3 and
 394 P6/C4 chemistry on a PacBio RSII machine with 6-hour movies. For FN, 3pM of SMRTbell
 395 libraries were diffusion-loaded and sequenced on a Sequel machine with v2 chemistry and 10-
 396 hour movies.

397 **Falcon assembly and Illumina polishing**

398 FALCON (Chin et al. 2016) was used to generate genome assemblies for PK (v0.4.0) and FN
 399 (v1.8.6). Briefly, raw subread data was filtered to remove the shortest reads to an approximate
 400 coverage of 70x for each genome, leaving 8,003,220 (80.2%) of subreads for PK and 6,646,226
 401 (62.6%) of subreads for FN, or approximately 58 Gbps for each. Preassembled reads (i.e. error-
 402 corrected reads) were then created with a length cutoff of $\geq 6,000$ bp for PK and $\geq 7,000$ bp for
 403 FN, resulting in 2,239,051 and 5,323,023 preassembled reads respectively. The PK and FN
 404 genomes were then assembled using preassembled reads with a minimum length of 9 Kbp or 7
 405 Kbp, respectively. Additional relevant assembly parameter settings for FN were as follows:

```
406 pa_HPCdaligner_option      : -B128 -t16 -e0.8 -M24 -l1200 -k18 -h256 -w8 -s100 -T12
407 ovlp_HPCdaligner_option   : -B128 -M24 -k24 -h600 -e.92 -l1800 -s100 -T12
408 falcon_sense_option        : --output_multi --min_cov_aln 4 --min_idt 0.70 --min_cov 4
409                           : --max_n_read 200
410 falcon_sense_skip_contained : False
411 overlap_filtering_setting  : --max_diff 120 --max_cov 120 --min_cov 4
```

412
 413 Similar assembly parameters were used for PK, except that min_cov was set to 3.

414

415 Each FALCON assembly was corrected with paired-end Illumina reads using Pilon version 1.22
416 (Walker et al. 2014) after mapping available Illumina sequencing data (van Bakel et al. 2011) to
417 the Falcon assembled genomes using BWA-MEM (Li 2013) (version 0.7.8) with an average of
418 96x (PK) and 23x (FN) coverage. Correction was performed with the “diploid” flag and the
419 “bases” flag set to correct only indels and snps. A total of 1,511,828 insertions and 228,876
420 deletions were corrected in the FN assembly, and 1,807,453 insertions and 283,918 deletions
421 were corrected in the PK assembly.

422

423 **Repeat content analysis**

424 Repeats in the FN and PK genomes were predicted *de novo* and classified using
425 RepeatModeler (v1.0.11; <http://www.repeatmasker.org/RepeatModeler/>). RepeatModeler was
426 applied to each assembly with the ‘ncbi’ engine (RMBlas v2.2.28) provided with
427 RepeatModeler. Other prerequisite components installed with the RepeatModeler package
428 included RECON v1.0.8 and RepeatScout v1.0.5 (Price et al. 2005), Tandem Repeat Finder
429 v4.0.4 (Benson 1999), and Repbase-derived RepeatMasker libraries
430 (<http://www.girinst.org/server/RepBase/>) from January 2017. The *de novo* repeat classification
431 provided by RepeatModeler was filtered to remove families with a >1kb BLAT (Kent 2002)
432 alignment to PK transcripts. The final filtered RepeatModeler output was then used as input for
433 RepeatMasker (<http://www.repeatmasker.org>) to produce a masked version of the assembly
434 and obtain the genomic positions of annotated repeats.

435

436 **Assessment of genome assembly completeness**

437 The completeness of each genome assembly was assessed using BUSCO v3.0 (Simao et al.
438 2015) and the set of eudicotyledons single-copy orthologs from OrthoDB v10, with default
439 arguments in the provided virtual machine instance.

440

441 **Comparison of PK and FN scaffolds**

442 PK and FN assemblies were aligned using LASTZ (Harris 2007) version 1.04.00 with the –
443 ungapped and –notransitions options, and a step of 20. Alignments with an identity of $\leq 95\%$
444 and a length of $\leq 2000\text{bp}$ were removed. To produce a dotplot, FN contigs were initially
445 ordered by size along the Y axis. Next, PK contigs were ordered and orientated on the X axis by
446 the position of their best hit on the Y axis. FN contigs were then reordered on the Y axis
447 according to their best hit to the newly ordered contigs on the X axis. This process was repeated
448 until the order of contigs on the X axis, and the order of contigs on the Y axis converged.

449

450 **Illumina sequencing of the FN and F1 individuals.**

451 Dual-indexed libraries were prepared using the Nextera DNA Library Preparation Kit (Illumina),
452 pooled equimolar, and sequenced on the HiSeq2500 platform, yielding 529.9 Gbp total. FN was
453 sequenced independently on the NextSeq500 platform, yielding 49.9 Gbp.

454

455 **Building the genetic map**

456 *Quality filtering.* Barcode and adaptor sequences were filtered from all FN and F1 Illumina PE
457 reads. FN reads were further filtered using sickle with the flags -q 20 -l 125. [Joshi NA, Fass JN.
458 (2011). Sickle: A sliding-window, adaptive, quality-based trimming tool for FastQ files (Version
459 1.33)]. PK Illumina 2x100 PE reads from the 2011 draft genome were also filtered using sickle,
460 with the flags -q 20 -l 90.

461

462 *Variant calling.* BWA-MEM (Li 2013) was used to map Illumina paired-end reads for FN, PK,
463 and the F1s to the PK FALCON assembly, after which Picard
464 [<http://broadinstitute.github.io/picard/>] was used for sorting, duplicate marking, and indexing the
465 alignments. To call variants for the F1s, we used the mpileup function from bcftools (Li et al.
466 2009) over all of the F1 individuals and both parents to overcome spots of lower coverage in the

467 F1s. Variants were also called individually for each parent using the GATK HaplotypeCaller
468 (McKenna et al. 2010) to be used as input for haplotype phasing.

469
470 *Phasing the parental haplotypes.* Haplotypes for the parents were phased using HapCUT2
471 (Edge et al. 2017), using the `-pacbio 1` argument to improve accuracy with PacBio reads and
472 the `-ea 1` argument to calculate switch quality scores. As input, parental SNPs called by
473 HaplotypeCaller on the Illumina data were provided in conjunction with PacBio raw reads. This
474 was done to both increase the length of the resulting haplotype blocks and boost confidence in
475 the phasing by requiring agreement between the two sets of data. To further increase
476 confidence we only used SNPs with a quality score > 25 and read coverage between 6 and 46,
477 and that were >5 bases away from an indel. Haplotype blocks were then split if the switch
478 quality score was less than 30. Finally, only blocks with > 10 SNPs were retained to use as
479 input for SOILoCo.

480
481 *Genotyping the F1s.* The SOILoCo method (Scaglione et al. 2016) was used to genotype the
482 F1s at each haplotype block, using the output of HapCUT2 and the variants called by mpileup.
483 Required values and divergence from the default parameters are as follows. For `vcf2strings.pl`
484 minor allele frequencies 1 and 2 (`--MAF-1` and `--MAF-2`) were set to 0.25 and 0.75 respectively.
485 This step allows the removal of any markers that may display segregation distortion (8.5% of
486 markers show some degree of segregation distortion; scaffolds that do not get incorporated into
487 the genetic map have an average of 20% of markers displaying significant distortion). When
488 running `gt-hmm.pl` the minimum number of variant calls in a haplotype block (`--min-string`) was
489 set to 6, the probability of a crossing over event (`--switch-prob`) was set to $1E10^{-6}$, and the
490 probability of having reads containing both alleles at a heterozygous site (`--HCALL-prob`) was
491 set to 0.15. Lastly, population type (`--pop`) for `calls2csvr.pl` is set to cross pollinated (CP). This
492 process is run separately for each parent, with the two respective sets of haplotype blocks.

493
494 The scaffold containing *CBDAS* and the scaffold from the PK FALCON assembly containing
495 *THCAS* were genotyped separately. As both scaffolds do not have a counterpart in the other
496 parental assembly, genotypes were extracted from variant loci that meet the following criteria: an
497 allele frequency of 0.5 in the parent harboring the scaffold, no coverage in the opposing parent,
498 an allele frequency of 0.5 in the F1s, and at which all F1s are homozygous. The scaffold
499 containing *THCAS* is the only scaffold from the PK FALCON assembly that was placed in the
500 genetic map.

501
502 *Forming linkage groups.* R/qtl (Broman et al. 2003) was employed to divide haplotype blocks for
503 each parent separately across linkage groups using the `formLinkageGroups` function with
504 maximum recombination frequency (`max.rf`) set to 0.05 and minimum LOD (`min.LOD`) set to 15.
505 The resulting linkage groups were compared against one another to identify any pairs of linkage
506 groups with a mean recombination frequency of > 0.8 between the haplotype blocks they
507 contain, in which case the `switchAlleles` function was used to swap the alleles for all the
508 haplotype blocks in the smaller linkage group, and `formLinkageGroups` was called again.
509 Afterwards, R/qtl functions `checkAlleles`, `switchAlleles`, and `formLinkageGroups` were run in
510 succession two more times to further identify and fix haplotype blocks with swapped alleles. All
511 linkage groups with >100 haplotype blocks were passed to the ordering step. For PK, there
512 were 11 linkage groups with > 100 haplotype blocks, however two of them just missed the cutoff
513 for being joined together and were therefore combined. Further support for combining these
514 linkage groups came from a comparison with the FN map, in which the scaffolds held in these
515 two PK linkage groups were held in a single FN linkage group.

516
517 *Ordering scaffolds.* Haplotype blocks were ordered within each linkage group using MSTmap
518 (Wu et al. 2008) with the Kosambi distance function. Three rounds of ordering were done with a

519 smoothing step in between carried out using the Perl implementation of the SMOOTH correction
520 algorithm (van Os et al. 2005) that is provided with the SOILoCo pipeline using an error
521 threshold of 0.85. Correspondence between the two parental sets of linkage groups was
522 determined based on similarity in the sets of scaffolds belonging to each linkage group. To
523 handle ambiguity in scaffold placement, if the haplotype blocks for any given scaffold were
524 distributed over more than one linkage group within or between parental maps, a census was
525 taken to determine the correct linkage group and haplotype blocks that did not agree with the
526 majority were removed. If fewer than half of the haplotype blocks were in agreement, all
527 haplotype blocks for that scaffold were removed and the scaffold was not placed in either
528 parental map. Finally, for each scaffold within each map, a distribution of the genetic positions
529 (in cM) for all haplotype blocks belonging to the scaffold was established and any outlier blocks
530 were removed. After removal of ambiguous haplotype blocks and scaffolds, a final round of
531 ordering was carried out for each parental map.

532

533 *Merging the genetic maps.* To translate each parental map from haplotype blocks to scaffolds in
534 order that they could be merged, scaffold placements were determined by averaging the
535 locations of the haplotype blocks belonging to each scaffold. The genetic maps for PK and FN
536 were then merged using MergeMap (Wu et al. 2011) with the weight of the FN genetic map set
537 to 2 and the weight of the PK genetic map set to 1 because it was based off the FN FALCON
538 assembly.

539

540 **Gene cloning**

541 *CBDAS* was amplified from DNA isolated from Finola leaves using gene-specific primers
542 (Forward: 5' – CTGCAGGAATGAAGTACTCAACATTCTCCTTTTGG – 3', Reverse: 5' –
543 AAGCTTTCATGGTACCCCATGATGATGCCGTGGAAGAG – 3'). PCR products were cloned
544 into pCR8/GW/TOPO (Invitrogen), excised as PstI/KpnI fragments, and cloned into pPICz-alpha

545 B (Invitrogen). The expression vectors were then transformed into *Pichia pastoris* strain X-33
546 (Invitrogen) by electroporation. Positive recombinants were selected for by plating transformed
547 cells on YPD plates supplemented with 25 $\mu\text{g mL}^{-1}$ phleomycin (Invivogen). To screen for
548 activity, colonies were used to inoculate 5 mL BMG cultures, which were grown for two days at
549 37 °C with shaking. The cells were then pelleted by centrifugation, resuspended in 5 mL BMM
550 media and grown for four days at 20 °C with shaking with the addition of 1% methanol daily.
551 Enzyme activity was tested by directly adding CBGA to clarified culture media, incubating
552 overnight at 37 °C, and then analyzing products by HPLC as previously described (Stout et al.
553 2012).

554

555 **Quantitative PCR**

556 RNA extraction, cDNA generation, and qRT-PCR conditions were identical to those previously
557 reported (Stout et al. 2012). *CBCAS* primers (Forward: 5' – CGGATGTACTGTTATGCTCCAA –
558 3'; Reverse: 5' – CATTCTCCATTAAAATAAGAAAGACAA – 3') were designed from alignments
559 of *THCAS-like* genes identified in the cannabis genome to ensure their selectivity. Primers were
560 tested using cloned *THCAS*, *CBDAS* and *CBCAS* as templates. Any primer set that amplified a
561 non-target cDNA was discarded. Primer efficiencies were extrapolated from raw amplification
562 data using LinRegPCR (Ruijter et al. 2009).

563

564 **Recombinant CBCAS enzyme expression and purification**

565 The culture with the highest CBCAS activity was selected for scaled up production. 1 mL of the
566 initial culture was used to inoculate two 40 mL BMG cultures, which were grown for two days at
567 37 °C. These cultures were then used to initiate two 400 mL modified BMM cultures that were
568 buffered with 10 mM HEPES (pH 7) and were supplemented with riboflavin at 20 mg L^{-1} . These
569 cultures were grown at 20 °C with shaking at 100 RPM for five days, with methanol was added

570 to 1% by volume each day. The cultures were then clarified by centrifugation, and the resulting
571 media was filtered and passed over two Bio-scale Mini CHT hydroxyapatite cartridges (Biorad)
572 at a flow rate of 1.5 mL min⁻¹ at 4 °C. The cartridges were then attached in series to an AKTA
573 FPLC system* (GE Healthcare) and eluted with a 75 mL linear gradient from 5 mM sodium
574 phosphate pH 7 to 500 mM sodium phosphate pH 7. Active fractions were pooled, concentrated
575 with a 30 kDa cutoff Centricon filter (Millipore), and buffer exchanged into 20 mM citrate pH 4.7
576 using a PD10 column (GE Healthcare). The resulting fraction was then injected onto a MonoS
577 5/50 cation exchange column (GE Healthcare) and eluted with a 40 mL linear gradient of 20 mM
578 citrate pH 4.7 to 20 mM citrate pH 4.7 + 500 mM NaCl. Active fractions were pooled,
579 concentrated with a 30 kDa cutoff Centricon filter and injected onto a Hiload 26/60 Superdex
580 200 size exclusion column (GE Healthcare). Proteins were eluted with a single column volume
581 of 20 mM citrate pH 5.0 + 150 mM NaCl. Throughout the purification, 1/10th volume of each
582 fraction was retained for analysis to judge purity. Protein was isolated from each fraction using
583 15 µL of StrataClean resin (Stratagene) and analyzed by SDS PAGE.

584

585 **Enzyme assays and HPLC quantification of reaction products**

586 To test for CBCAS enzyme activity during the protein purification, 150 µL of protein fraction was
587 mixed with 50 µL of 500 µM sodium citrate buffer pH 5.0 and 20 µmoles of CBGA and incubated
588 overnight at 37 °C. The reactions were then extracted twice with ethyl acetate, and the organic
589 fractions were pooled and dried in a Speedvac concentrator. The products were then
590 resuspended in 16 µL 50% methanol, of which 10 µL were analyzed by HPLC as previously
591 described (Stout et al. 2012). Reactions for enzyme kinetic analyses were comprised of 1 µg of
592 purified CBCAS, 100 mM sodium citrate pH 5.0, and 100 mM NaCl. These reactions were
593 performed under Michaelis-Menten conditions at 40 °C for one hour. Reaction product
594 extraction and analyses were the same as above.

595

596 **DATA ACCESS**

597 The PacBio sequence read data generated for genome assembly, the Illumina sequencing data
598 for the FN and F1 individuals, and the PK and FN genome assemblies have been deposited
599 publicly with links to BioProject accession number PRJNA73819 in the NCBI BioProject
600 database (<https://www.ncbi.nlm.nih.gov/bioproject/>).

601

602 **ACKNOWLEDGMENTS**

603 This work was supported by grants from the Canadian Institutes of Health Research (Operating
604 Grant MOP-126070 to TRH, JEP, and HvB, and Foundation Grant FDN-148403 to TRH). TRH
605 is a Scholar of the Canadian Institute for Advanced Research (CIFAR). HvB was supported in
606 part by R01 AI119145. This work was supported in part through the computational resources
607 and staff expertise provided by Scientific Computing at the Icahn School of Medicine at Mount
608 Sinai. We are grateful to the Donnelly Sequencing Centre (DSC) for assistance with Illumina
609 sequencing.

610

611 **DISCLOSURE DECLARATION**

612 Jonathan E. Page is the Chief Executive Officer and shareholder of a for-profit cannabis science
613 company, Anandia Laboratories Inc, based on Vancouver, Canada. In this position he receives
614 a salary. Anandia Laboratories Inc performs analytical testing for licensed cannabis producers in
615 Canada as well as working to develop new cannabis cultivars through breeding.

616

617 Larry Holbrook was the Chief Research Officer at CanniMed Therapeutics Inc until May 2018.
618 CanniMed is a for-profit company based in Saskatoon, Canada that produces medical cannabis
619 products for authorized patients. In this position he received compensation in the form of a

620 salary and stock options. In June 2018 he moved to CB3 Life Sciences where he is Chief
621 Scientific Officer.

622

623 Jonathan E. Page and Jake M. Stout have filed Patent WO2015196275 on the nucleotide
624 sequence encoding the enzyme CBCAS based reagents, and methods for producing
625 cannabinoids and/or altering cannabinoid production.

626 **REFERENCES**

- 627
- 628 Benson G. 1999. Tandem repeats finder: a program to analyze DNA sequences. *Nucleic acids*
629 *research* **27**: 573-580.
- 630 Broman KW, Wu H, Sen S, Churchill GA. 2003. R/qtl: QTL mapping in experimental crosses.
631 *Bioinformatics* **19**: 889-890.
- 632 Chin CS, Peluso P, Sedlazeck FJ, Nattestad M, Concepcion GT, Clum A, Dunn C, O'Malley R,
633 Figueroa-Balderas R, Morales-Cruz A et al. 2016. Phased diploid genome assembly with
634 single-molecule real-time sequencing. *Nature methods* **13**: 1050-1054.
- 635 Choi K, Zhao X, Kelly KA, Venn O, Higgins JD, Yelina NE, Hardcastle TJ, Ziolkowski PA,
636 Copenhaver GP, Franklin FC et al. 2013. Arabidopsis meiotic crossover hot spots
637 overlap with H2A.Z nucleosomes at gene promoters. *Nature genetics* **45**: 1327-1336.
- 638 Davik J, Sargent DJ, Brurberg MB, Lien S, Kent M, Alsheikh M. 2015. A ddRAD Based Linkage
639 Map of the Cultivated Strawberry, *Fragaria xananassa*. *PloS one* **10**: e0137746.
- 640 de Meijer EP, Bagatta M, Carboni A, Crucitti P, Moliterni VM, Ranalli P, Mandolino G. 2003. The
641 inheritance of chemical phenotype in *Cannabis sativa* L. *Genetics* **163**: 335-346.
- 642 Devinsky O, Cilio MR, Cross H, Fernandez-Ruiz J, French J, Hill C, Katz R, Di Marzo V, Jutras-
643 Aswad D, Notcutt WG et al. 2014. Cannabidiol: pharmacology and potential therapeutic
644 role in epilepsy and other neuropsychiatric disorders. *Epilepsia* **55**: 791-802.
- 645 Di Pierro EA, Gianfranceschi L, Di Guardo M, Koehorst-van Putten HJ, Kruisselbrink JW, Longhi
646 S, Troggo M, Bianco L, Muranty H, Pagliarani G et al. 2016. A high-density, multi-
647 parental SNP genetic map on apple validates a new mapping approach for outcrossing
648 species. *Horticulture research* **3**: 16057.
- 649 Divashuk MG, Alexandrov OS, Razumova OV, Kirov IV, Karlov GI. 2014. Molecular cytogenetic
650 characterization of the dioecious *Cannabis sativa* with an XY chromosome sex
651 determination system. *PLoS One* **9**: e85118.

- 652 Edge P, Bafna V, Bansal V. 2017. HapCUT2: robust and accurate haplotype assembly for
653 diverse sequencing technologies. *Genome research* **27**: 801-812.
- 654 Elsohly MA, Slade D. 2005. Chemical constituents of marijuana: the complex mixture of natural
655 cannabinoids. *Life sciences* **78**: 539-548.
- 656 Gore MA, Chia JM, Elshire RJ, Sun Q, Ersoz ES, Hurwitz BL, Peiffer JA, McMullen MD, Grills
657 GS, Ross-Ibarra J et al. 2009. A first-generation haplotype map of maize. *Science* **326**:
658 1115-1117.
- 659 Guo L, Winzer T, Yang X, Li Y, Ning Z, He Z, Teodor R, Lu Y, Bowser TA, Graham IA et al.
660 2018. The opium poppy genome and morphinan production. *Science* **362**: 343-347.
- 661 Harris RS. 2007. Improved pairwise alignment of genomic DNA. Ph.D. Thesis, The
662 Pennsylvania State University.
- 663 He N, Zhang C, Qi X, Zhao S, Tao Y, Yang G, Lee TH, Wang X, Cai Q, Li D et al. 2013. Draft
664 genome sequence of the mulberry tree *Morus notabilis*. *Nature communications* **4**: 2445.
- 665 Izzo AA, Capasso R, Aviello G, Borrelli F, Romano B, Piscitelli F, Gallo L, Capasso F, Orlando
666 P, Di Marzo V. 2012. Inhibitory effect of cannabichromene, a major non-psychoactive
667 cannabinoid extracted from *Cannabis sativa*, on inflammation-induced hypermotility in
668 mice. *Br J Pharmacol* **166**: 1444-1460.
- 669 Kent WJ. 2002. BLAT--the BLAST-like alignment tool. *Genome research* **12**: 656-664.
- 670 Kojoma M, Seki H, Yoshida S, Muranaka T. 2006. DNA polymorphisms in the
671 tetrahydrocannabinolic acid (THCA) synthase gene in "drug-type" and "fiber-type"
672 *Cannabis sativa* L. *Forensic science international* **159**: 132-140.
- 673 Kronenberg ZN, Hall RJ, Hiendleder S, Smith TPL, Sullivan ST, Williams JL, Kingan SB. 2018.
674 FALCON-Phase: Integrating PacBio and Hi-C data for phased diploid genomes. *bioRxiv*
675 doi:10.1101/327064.
- 676 Li H-L. 1974. An Archaeological and Historical Account of Cannabis in China. *Economic Botany*
677 **28**: 437-448.

- 678 Li H. 2013. Aligning sequence reads, clone sequences and assembly contigs with BWA-MEM.
679 *ArXiv e-prints*.
- 680 Li H, Handsaker B, Wysoker A, Fennell T, Ruan J, Homer N, Marth G, Abecasis G, Durbin R,
681 Genome Project Data Processing S. 2009. The Sequence Alignment/Map format and
682 SAMtools. *Bioinformatics* **25**: 2078-2079.
- 683 Liu S, Yeh CT, Ji T, Ying K, Wu H, Tang HM, Fu Y, Nettleton D, Schnable PS. 2009. Mu
684 transposon insertion sites and meiotic recombination events co-localize with epigenetic
685 marks for open chromatin across the maize genome. *PLoS genetics* **5**: e1000733.
- 686 Luo MC, You FM, Li P, Wang JR, Zhu T, Dandekar AM, Leslie CA, Aradhya M, McGuire PE,
687 Dvorak J. 2015. Synteny analysis in Rosids with a walnut physical map reveals slow
688 genome evolution in long-lived woody perennials. *BMC genomics* **16**: 707.
- 689 Maione S, Piscitelli F, Gatta L, Vita D, De Petrocellis L, Palazzo E, de Novellis V, Di Marzo V.
690 2011. Non-psychoactive cannabinoids modulate the descending pathway of
691 antinociception in anaesthetized rats through several mechanisms of action. *Br J*
692 *Pharmacol* **162**: 584-596.
- 693 Mascher M, Gundlach H, Himmelbach A, Beier S, Twardziok SO, Wicker T, Radchuk V, Dockter
694 C, Hedley PE, Russell J et al. 2017. A chromosome conformation capture ordered
695 sequence of the barley genome. *Nature* **544**: 427-433.
- 696 McKenna A, Hanna M, Banks E, Sivachenko A, Cibulskis K, Kernysky A, Garimella K, Altshuler
697 D, Gabriel S, Daly M et al. 2010. The Genome Analysis Toolkit: a MapReduce
698 framework for analyzing next-generation DNA sequencing data. *Genome research* **20**:
699 1297-1303.
- 700 Meinke D, Sweeney C, Muralla R. 2009. Integrating the genetic and physical maps of
701 *Arabidopsis thaliana*: identification of mapped alleles of cloned essential (EMB) genes.
702 *PloS one* **4**: e7386.

- 703 Morimoto S, Komatsu K, Taura F, Shoyama Y. 1997. Enzymological Evidence for
704 Cannabichromenic Acid Biosynthesis. *Journal of Natural Products* **60**: 854-857.
- 705 Morimoto S, Komatsu K, Taura F, Shoyama Y. 1998. Purification and characterization of
706 cannabichromenic acid synthase from *Cannabis sativa*. *Phytochemistry* **49**: 1525-1529.
- 707 Nutzmans HW, Osbourn A. 2014. Gene clustering in plant specialized metabolism. *Current*
708 *opinion in biotechnology* **26**: 91-99.
- 709 Osborne AL, Solowij N, Babic I, Huang XF, Weston-Green K. 2017. Improved Social Interaction,
710 Recognition and Working Memory with Cannabidiol Treatment in a Prenatal Infection
711 (poly I:C) Rat Model. *Neuropsychopharmacology : official publication of the American*
712 *College of Neuropsychopharmacology* **42**: 1447-1457.
- 713 Pate D. 1994. Chemical ecology of *Cannabis*. *Journal of the International Hemp Association* **2**:
714 32-37.
- 715 Peil A, Flachowsky H, Schumann E, Weber WE. 2003. Sex-linked AFLP markers indicate a
716 pseudoautosomal region in hemp (*Cannabis sativa* L.). *TAG Theoretical and applied*
717 *genetics Theoretische und angewandte Genetik* **107**: 102-109.
- 718 Pendleton M, Sebra R, Pang AW, Ummat A, Franzen O, Rausch T, Stutz AM, Stedman W,
719 Anantharaman T, Hastie A et al. 2015. Assembly and diploid architecture of an individual
720 human genome via single-molecule technologies. *Nature methods* **12**: 780-786.
- 721 Pisupati R, Vergara D, Kane NC. 2018. Diversity and evolution of the repetitive genomic content
722 in *Cannabis sativa*. *BMC genomics* **19**: 156.
- 723 Price AL, Jones NC, Pevzner PA. 2005. De novo identification of repeat families in large
724 genomes. *Bioinformatics* **21 Suppl 1**: i351-358.
- 725 Razumova OV, Alexandrov OS, Divashuk MG, Sukhorada TI, Karlov GI. 2016. Molecular
726 cytogenetic analysis of monoecious hemp (*Cannabis sativa* L.) cultivars reveals its
727 karyotype variations and sex chromosomes constitution. *Protoplasma* **253**: 895-901.

- 728 Ruijter JM, Ramakers C, Hoogaars WM, Karlen Y, Bakker O, van den Hoff MJ, Moorman AF.
729 2009. Amplification efficiency: linking baseline and bias in the analysis of quantitative
730 PCR data. *Nucleic acids research* **37**: e45.
- 731 Sakamoto K, Akiyama, Y., Fukui, K., Kamada, H., Satoh, S. 1998. Characterization; Genome
732 Sizes and Morphology of Sex Chromosomes in Hemp (*Cannabis Sativa* L.). *Cytologia* **63**:
733 459-464.
- 734 Scaglione D, Reyes-Chin-Wo S, Acquadro A, Froenicke L, Portis E, Beitel C, Tirone M, Mauro
735 R, Lo Monaco A, Mauromicale G et al. 2016. The genome sequence of the outbreeding
736 globe artichoke constructed de novo incorporating a phase-aware low-pass sequencing
737 strategy of F1 progeny. *Scientific reports* **6**: 19427.
- 738 Shimizu T, Tanizawa Y, Mochizuki T, Nagasaki H, Yoshioka T, Toyoda A, Fujiyama A,
739 Kaminuma E, Nakamura Y. 2017. Draft Sequencing of the Heterozygous Diploid
740 Genome of Satsuma (*Citrus unshiu* Marc.) Using a Hybrid Assembly Approach. *Front*
741 *Genet* **8**: 180.
- 742 Shinjyo N, Di Marzo V. 2013. The effect of cannabichromene on adult neural stem/progenitor
743 cells. *Neurochem Int* **63**: 432-437.
- 744 Simao FA, Waterhouse RM, Ioannidis P, Kriventseva EV, Zdobnov EM. 2015. BUSCO:
745 assessing genome assembly and annotation completeness with single-copy orthologs.
746 *Bioinformatics* **31**: 3210-3212.
- 747 Sirikantaramas S, Morimoto S, Shoyama Y, Ishikawa Y, Wada Y, Shoyama Y, Taura F. 2004.
748 The gene controlling marijuana psychoactivity: molecular cloning and heterologous
749 expression of Delta1-tetrahydrocannabinolic acid synthase from *Cannabis sativa* L. *The*
750 *Journal of biological chemistry* **279**: 39767-39774.
- 751 Stout JM, Boubakir Z, Ambrose SJ, Purves RW, Page JE. 2012. The hexanoyl-CoA precursor
752 for cannabinoid biosynthesis is formed by an acyl-activating enzyme in *Cannabis sativa*
753 trichomes. *Plant J* **71**: 353-365.

- 754 Taura F, Sirikantaramas S, Shoyama Y, Yoshikai K, Shoyama Y, Morimoto S. 2007.
755 Cannabidiolic-acid synthase, the chemotype-determining enzyme in the fiber-type
756 Cannabis sativa. *FEBS letters* **581**: 2929-2934.
- 757 van Amerongen G, Kanhai K, Baakman AC, Heuberger J, Klaassen E, Beumer TL, Strijers RL,
758 Killestein J, van Gerven J, Cohen A et al. 2017. Effects on Spasticity and Neuropathic
759 Pain of an Oral Formulation of Delta9-Tetrahydrocannabinol in Patients With
760 Progressive Multiple Sclerosis. *Clinical therapeutics* doi:10.1016/j.clinthera.2017.01.016.
- 761 van Bakel H, Stout JM, Cote AG, Tallon CM, Sharpe AG, Hughes TR, Page JE. 2011. The draft
762 genome and transcriptome of Cannabis sativa. *Genome Biol* **12**: R102.
- 763 van Os H, Stam P, Visser RG, van Eck HJ. 2005. SMOOTH: a statistical method for successful
764 removal of genotyping errors from high-density genetic linkage data. *TAG Theoretical
765 and applied genetics Theoretische und angewandte Genetik* **112**: 187-194.
- 766 Walker BJ, Abeel T, Shea T, Priest M, Abouelliel A, Sakthikumar S, Cuomo CA, Zeng Q,
767 Wortman J, Young SK et al. 2014. Pilon: an integrated tool for comprehensive microbial
768 variant detection and genome assembly improvement. *PloS one* **9**: e112963.
- 769 Watt G, Karl T. 2017. In vivo Evidence for Therapeutic Properties of Cannabidiol (CBD) for
770 Alzheimer's Disease. *Frontiers in pharmacology* **8**: 20.
- 771 Weiblen GD, Wenger JP, Craft KJ, ElSohly MA, Mehmedic Z, Treiber EL, Marks MD. 2015.
772 Gene duplication and divergence affecting drug content in Cannabis sativa. *The New
773 phytologist* **208**: 1241-1250.
- 774 Wu Y, Bhat PR, Close TJ, Lonardi S. 2008. Efficient and accurate construction of genetic
775 linkage maps from the minimum spanning tree of a graph. *PLoS genetics* **4**: e1000212.
- 776 Wu Y, Close TJ, Lonardi S. 2011. Accurate construction of consensus genetic maps via integer
777 linear programming. *IEEE/ACM transactions on computational biology and
778 bioinformatics* **8**: 381-394.

779 Zamudio N, Barau J, Teissandier A, Walter M, Borsos M, Servant N, Bourc'his D. 2015. DNA
780 methylation restrains transposons from adopting a chromatin signature permissive for
781 meiotic recombination. *Genes & development* **29**: 1256-1270.
782