1 **Single molecule sequencing of THCA synthase reveals copy number variation in**
2 **modern drug-type *Cannabis sativa* L.**

3 Kevin J. McKernan, Yvonne Helbert, Vasisht Tadigotla, Stephen McLaughlin, Jessica
4 Spangler, Lei Zhang, Douglas Smith
5 Courtagen Life Sciences, 12 Gill Street Woburn MA, 01801

6

7 **Summary**

8    • Cannabinoid expression is an important genetically determined feature of
9      *cannabis* that presents clinical and legal implications for patients seeking
10      cannabinoid specific therapies like Cannabidiol (CBD).
11    • Cannabinoid, terpenoid, and flavonoid marker assisted selection can accelerate
12      breeding efforts by offering genetic tools to select for desired traits at an early
13      stage in growth. To this end, multiple models for chemotype inheritance have
14      been described suggesting a complex picture for chemical phenotype
15      determination.
16    • Here we explore the potential role of copy number variation of THCA Synthase
17      using phased single molecule sequencing and demonstrate that copy number
18      and sequence variation of this gene is common and suggests a more nuanced
19      view of chemotype prediction.

20

21 **Introduction**
22 The genetics of cannabis chemotype has been extensively studied due to the highly
23 selected THCA and CBDA phenotypes. A co-dominant model of inheritance of THCA
24 and CBDA synthase (THCAS, CBDAS)(de Meijer *et al.*, 2003) has been demonstrated
25 describing a Bt:Bd allele in linkage to the synthase genes(Onofri *et al.*, 2015; Weiblen *et*
26 *al.*, 2015).  In addition to this Bt:Bd allele, single nucleotide polymorphisms in the FAD
27 binding domain of THCAS have been described that impair the function of the
28 enzyme(Sirikantaramas *et al.*, 2004). Cascini *et al.* demonstrated varying copy number
29 of a highly conserved region in THCAS but did not find this assay to be correlative to
30 THCAS expression(Cascini *et al.*, 2012; Cascini *et al.*, 2013). Complimenting de Meijer
31 and Sirikantaramas work, van Bakel *et al.* reported a gene of unknown function (AAE3)
32 that was differentially replicated in drug versus fiber type cannabis(van Bakel *et al.*,
33 2011).

34

35 Whole genome sequencing using 454 and Illumina sequence data demonstrated
36 excessive polymorphic coverage over the single exon THCAS gene(McKernan,
37 https://archive.org/details/SequencingTheCannabisGenome) ) but allelic phasing of
38 these important genes has not been possible using short read sequencing and
39 conventional genome assembly technologies. Using bacterial cloning, Onofri *et al*.
40 demonstrated a polyploid status for THCA synthase (THCAS) confirming the putative
41 copy number variation observed with whole genome sequencing. Nevertheless, this
42 polyploid status appeared to be a rare event occurring in a single United States bred
43 cultivar and was not correlated with increased THCA expression. Weiblen *et*
44 *al(Weiblen et al., 2015).* and Cascini *et al.(Cascini et al., 2012; Cascini et al., 2013)*
45 further supported polyploidy with quantitative PCR and bacterially cloned THCAS DNA
46 but deep sampling of alleles may have been limited by bacterial cloning methods. Here,

1  we apply over 10,000 long-read single molecule sequences per sample in combination
2  with Illumina sequencing methods, to phase the multiple copies of THCA synthase in
3  over a dozen modern medicinal cultivars(Eid *et al.*, 2009). Allelic phasing of this 1.6kb
4  gene adds further support to additional mechanism of polyploidy in chemical
5  inheritance.
6
7  Multiple primer pairs were designed to amplify THCAS, and the amplicons were
8  sequenced using Pacific Bioscience's circular consensus sequencing (CCS) method and
9  a size-selected Illumina paired-250bp read system (Table1). Given the 40kb reads of
10 the Pacific Biosciences (PacBio) platform, it is ideally suited for phasing polymorphic
11 genes that are longer than the read lengths accessible by Sanger sequencing. Bacterial
12 cloning and primer walking has been used traditionally to address this, but the
13 methods can be cumbersome, expensive and are susceptible to cloning bias(Metzker,
14 2010). Single molecule sequencing enables one to deeply sequence a PCR product
15 without bacterial cloning to eliminate this bias. The higher raw error rate of single
16 molecule sequencing is overcome by ligating the target amplicon into a circle and
17 performing 20X rolling circle sequencing(Chaisson *et al.*, 2015). This CCS procedure
18 significantly reduces the stochastic error associated with single molecule detection
19 while providing long phased haplotypes of each copy of the amplified gene. These
20 phased haplotypes were then verified by paired-end Illumina sequencing to confirm
21 pure phase with homozygous alignments.  Cultivars that vary in CBDA and THCA
22 content were sequenced to assess the diverse haplotypes present in *cannabis* and to
23 resolve putative pseudogenes from intact open reading frames.
24
25 **Results**
26 We demonstrate that multiple copies of THCAS with distinct sequences are present in
27 many common dispensary and coffee house cultivars in circulation in 2015 (Figure 1).
28 The ascertainment of THCAS copy number is dependent on primer selection
29 suggesting local divergence in the 5' and 3' untranslated regions (UTRs).
30
31 Utilizing two sets of primer pairs (within and external to THCAS), we amplified
32 polyploid copies of THCAS and sequenced them using the Pacific Biosciences CCS
33 method. To mitigate PCR error we required a minimum of 10 independent, but
34 identical sequences from different zero mode waveguides (ZMW) to establish a
35 haplogroup. For each CCS, the molecule was read on the forward and reverse strands
36 20-30 times to form a circular consensus that was over 99.95% accurate (Table 2).
37
38 Once PacBio long reads were error corrected using CCS, Illumina Nextera libraries
39 were generated from size-selected 850 bp fragments derived from the PCR products,
40 and were sequenced with 250 base paired-end reads. The resulting read pairs could be
41 accurately mapped to each PacBio haplogroup to produce homozygous alignments for
42 each allele (Figure 2.). These homozygous alignments were helpful in ascertaining
43 whether all haplogroups were accurately represented in PacBio data and if there were
44 any remaining variants ambiguously mapped. Unphased polymorphic data is difficult
45 to translate into amino acid sequence and obscures accurate assessment of functional
46 alleles.

2

1
2    The THCAS amplicons were initially amplified using the primers described by Onofri *et*
3    *al.(Onofri et al., 2015)* This primer set spans the start and stop codons of the THCAS
4    gene and thus is incapable of providing sequence information for the first and last
5    25bp of the gene. Moving the primers out into the flanking non-coding regions (MGC-
6    2130 external primer set) delivered fewer haplogroups for all cultivars tested,
7    presumably due to the divergence of these regions in the different gene copies (Figure
8    3).
9
10   To ascertain the reproducibility of this approach, Chemdog91 was amplified 3 different
11   times and run on 3 different SMRT cells. One sample was run as a single-plex sample
12   on a SMRT cell to afford high enough sampling to ensure no alleles were missed. While
13   each sample consistently sequenced the active haplogroups to high coverage, a few
14   inactive or unknown haplogroups emerged on only 2 of the 3 SMRT cells. Since one of
15   the intermittent haplogroups has a unique variant (S221F genotype) private to
16   Chemdog91, we can rule out contamination from other gDNA being the source of these
17   intermittent haplogroups. We cannot rule out the possibility of other diverged copies
18   of inactive THCAS stochastically amplifying given the PCR reactions were set up
19   independently. Gradient PCR at lower annealing temperatures with the Onofri primers
20   may unveil additional diverged pseudogenes. A few intermittent haplogroups had
21   depressed consensus accuracy and are detectable with appropriate quality filtering
22   tools. As a result of these conditions we do not believe the amplicon subread coverage
23   numbers derived from the Onofri primer sets can be taken as a direct proxy for
24   absolute copy number in the genome. It is possible that some of the high coverage,
25   active THCAS haplogroups may be more than diploid in copy number yet exist as non-
26   diverged sequence in the genome. With the MGC-2130 data being multiplexed, any
27   inter-strain coverage or inferred copy number differences are more likely the result of
28   normalization of barcoded samples prior to multiplexing. Only intra-strain copy
29   number can be assessed with these data.
30
31   We segregated the haplogroups based on two genbank accession sequences that have
32   been confirmed as encoding active (Q8GTB6.1 also identical to genbank accession
33   number E33090) and inactive (Q33DQ2.1) THCAS(Sirikantaramas *et al.*, 2004; Kojoma
34   *et al.*, 2006). While these two haplogroups differ by 37 amino acids from each other,
35   several other haplogroups diverged by up to 2 and 6 amino acids from these two
36   reference sequences (Figure 4, Table 3). Other haplogroups were identified that
37   contained termination or frame-shifting variants. These had higher similarity to
38   Q33DQ2.1 than Q8GTB6.1, suggesting an inactive heritage. Phylogenetic trees were
39   constructed that clustered the divergent haplogroups around Active, Inactive, and
40   frame-shifted clusters (Figure 4).
41
42   We surveyed CBDA and THCA dominant cultivars seen in Table 1,2 and 3. Cultivars
43   measured to have CBDA-rich chemotypes like WZ_CBD contained 3 haplogroups, one
44   of which was clearly a pseudogene (Figure 1). The other two haplogroups were either
45   identical to, or had 6 amino acid mismatches with the "Inactive THCAS" in genbank
46   (Q33DQ2.1). Since this cultivar has low but measured THCA (0.6% THCA : 20%CBDA),

3

1    it could be assumed that the 6 amino acid divergence from the inactive form
2    represents a very weakly active THCAS gene. The cultivar was positive for the Bd:Bd
3    genotype described by de Miejer, and if this allele governs the binary expression of
4    CBDAS over THCAS then the expression level of THCAS is uncertain and the impact of
5    the 6 amino acid variants on activity is impossible to infer. Some level of chemotype
6    expression exists and suggests a leaky Bd:Bd expression is responsible for the 0.6%
7    THCA. Of note, this haplogroup failed to amplify with the MGC-2130 external primer
8    set (Table 4).
9
10   The 6 amino acid diverged THCAS haplogroup is also found in other THCA dominant
11   cultivars like Chemdog91, however Chemdog91 is a high THCA expressing cultivar and
12   has 5 THCAS copies in total (one with a stop codon, one that is inactive and three that
13   are putatively active). One haplogroup in Chemdog91 is identical in amino acid
14   sequence to the active THCAS sequence (Q8GTB6.1) and the other has a common
15   A250D variant of the active form. Although we did not collect RNA from these samples,
16   Onofri, *et al.* also describe some of these alternative alleles and demonstrates RNA
17   expression(Onofri *et al.*, 2015).
18
19   Adding to the observations of Onofri, *et al.*, we were able to find THCA producing
20   strains that exclusively held the ambiguous genotypes described in their paper as  (1/4
21   , 2/1 , 2/2 , 2/3). The 1/4 genotype described by Onofri, *et al.* is an A250D variant of
22   the active form. The Cultivar "AK-47" presented with 3 haplogroups of which one had a
23   stop codon, one had an inactive allele (Q33D12.1) and one presented with a single
24   active THCAS haplotype with an A250D variant. We can deduce from this that the
25   A250D allele is the only haplogroup capable of synthesis of the high levels of THCA
26   found in AK-47.
27
28   Interestingly, another high THCA strain (Black84) shared this A250D variant but also
29   had the I63L variant described and deemed active by Onofri et al. (Table 3). Only the
30   A250D variant amplified with the MGC-2130 primer set bringing some question to the
31   activity of the I63L variant in our cultivars.  Chemdog91, in addition to its Q8GTB6.1
32   active allele, also has an A250D haplogroup. Both of these haplogroups amplify with
33   the MGC-2130 primer set in Chemdog91 and only the A250D allele amplifies in
34   Black84 suggesting the I63L variant may be inactive (Table 4). Further work is
35   required to characterize the I63L variants functionality.
36
37   The other THCA:CBDA hybrid cultivar "Blueberry Essence" has a single haplogroup
38   that amplified with both primers, suggesting the sample had identical maternal and
39   paternal alleles for THCAS sequence. This sequence has a single amino acid variant
40   from the active form (P333R). Considering the structure of the gene, this proline
41   change may alter the activity of the THCAS and explain the 8.11% THCA, 10.8%CBDA
42   expressed in Blueberry Essence(Shoyama *et al.*, 2012). Figure 6 demonstrates a Hybrid
43   Bd:Bd genotype. The difference in Bd:Bd vs Bt:Bd genotypes relies on measuring the
44   magnitude of the 3rd peak in the electropherogram but in our hands the assay did not
45   produce a clear presence-absence result allowing unambiguous differentiation of these
46   two alleles (Figure 6.).  Interestingly, the B1080/B1192 primers described in Pacifico

1    *et al.* were confirmed to be the B1180/B1192 genotyping primers utilized in Table 2 of
2    Onofri *et al.* (personal communication G. Mandolino). This forward primer targets a
3    region conserved between active and inactive THCAS while the mid portion of the
4    reverse primer rests on a H494P and A495E variant between active and inactive
5    alleles. Although the 3 prime end of the primer is 12 bases upstream of these primed
6    variants, it is possible but unlikely that this primer set differentiates between active
7    and inactive transcripts of THCAS.
8
9    Otto and Sour Tsunami are both CBDA dominant chemotypes but demonstrate 2 and 6
10   haplogroups respectively. Three of the haplogroups in Sour Tsunami have stop codons.
11   The remaining three haplogroups are exact active and inactive haplogroups with an
12   additional 6 amino acid variant of the inactive haplogroup with unknown activity. Otto
13   has 2 haplogroups, one inactive copy that only amplifies with Onofri primers and
14   another that amplifies with both primer sets and represents a single active Q8GTB6.1
15   haplogroup with an A411V variant.
16
17   **Illumina Sequencing of the Bt:Bd alleles**
18   We amplified the Bt:Bd alleles from all of the samples in Table 1 (Figure 6). While 4 of
19   the THCA positive samples demonstrated the 190bp band, 3 did not. A select few of
20   these products were Illumina sequenced to gain a better understanding of the loci the
21   B190:200 primers amplify. Assembling these amplicons produced several contigs per
22   cultivar where the inserts can be mapped to a tandem repeat in CanSat3 Scaffold
23   19079 (Figure 7.). Only 9-10bp of the 3' end of the B190/B200 primers match the
24   reference at appropriate distances implying a Tm sensitive assay that is likely highly
25   sensitive to salts in various DNA isolations. Of note, the referenced methods in
26   Mandolino et al. recommend 38C annealing temperatures(Mandolino, 1999). We
27   confirmed these alignments with cross comparisons to other public whole genome
28   assemblies with of Chemdog91 and LA Confidential(McKernan,
29   https://aws.amazon.com/datasets/the-cannabis-sativa-genome/).  Contig_60162 of
30   LA Confidential has a 100% 159bp alignment of WZ-CBD B190/B200 allele while
31   Chemdog91s contig_96784 has a 148/159bp alignment with many other contigs
32   showing similar homology. Various polymorphisms also exist in the alignments. While
33   this appears to be a repeat, it is also rich in secondary structure notorious for error
34   prone sequencing(Nakamura *et al.*, 2011).
35
36   While these are useful markers for tracking chemotype, its unclear by their sequence
37   alone what function they play. Even though THCAS rests on CanSat3 scaffold 19603
38   and shares little homology to CanSat3 scaffold 19079, one possibility is that the
39   regions are simply in linkage to active and inactive alleles of CBDAS and THCAS and
40   closer markers to the gene may be more predictive in a polymorphic population.
41
42   **Discussion**
43   Even though absolute copy number does not correlate with CBDA dominant or THCA
44   dominant chemotypes, a more nuanced view of the amino acid alterations in the
45   haplogroups reveals an interesting pattern. All CBDA strains contain at least one
46   haplogroup with similarity to the inactive Q33DQ2.1 haplogroup.  The most common

1  second haplogroup in CBDA dominant cultivars is the 6 amino acid diverged form of
2  Q33DQ2.1 (N90F, I266T, L370F, D420G, Y471C, T492S). Additionally, no CBDA
3  dominant strains amplified with perfect Q8GTB6.1 active alleles with the MGC-2130
4  primers. The only two cultivars to amplify with the MGC-2130 primers had variants
5  (Otto-A411V, and Blueberry Essence-P333R). The other three CBDA cultivars failed to
6  amplify with MGC-2130 primers (WZ-CBD, Sour Tsunami, CannaTsu) despite gradient
7  PCR optimization attempts.
8
9  This suggests UTR sequences of various THCAS haplogroups are associated with
10  different chemotypes, and invokes an interest in the linked Bt:Bd allele described by de
11  Meijer *et al*.). Weiblen *et al*. demonstrated this marker is tightly linked to THCAS and
12  CBDAS. Pacifico *et al.* implied these B190/200 markers were imperfect in resolving the
13  Bt:Bd from Bd:Bd alleles but never published the implied and improved B1080/1192
14  primers(Pacifico, 2006). Although, it is possible that the Bt allele is just in phase with
15  the MGC-2130 primer set and predicts active form THCAS while the Bd allele is in
16  phase with functional CBDAS, the current genome assemblies can only anchor partial
17  B190 or B200 primer sequences (used to amplify the Bt:Bd allele) to a genomic contig
18  (scaffold 19079). Sequencing these Bt:Bd products with Illumina sequencing
19  demonstrated high coverage and high homology between the various sized bands
20  suggesting a recently expanded repetitive element that is hard to assemble (Figure 7).
21
22  The UTR based MGC-2130 primers imply the THCAS replicated inactive events did not
23  include 5 or 3' UTRs. Primers that rest 387 & 160 bases back from the start and stop of
24  the gene amplify fewer haplogroups of THCAS presumably due to the genomic
25  divergence flanking the transcribed sequences.  10/13 of the cultivars amplified
26  successfully with the more distant MGC- 2130 primers (WZ-CBD, Sour Tsunami, and
27  CannaTsu failed to amplify). Only one hybrid cultivar demonstrated a single
28  haplogroup with both primers (Blueberry Essence Q8GTB6.1 P333R). The MGC-2130
29  primers appear to predominantly amplify the active forms and could possibly be used
30  to better inform chemotype prediction. Due to the high polymorphism rate of this
31  species and locus, sampling of additional cultivars with well-measured chemotypes is
32  underway.
33
34  These haplogroups help to inform why copy number assays targeting highly conserved
35  regions of THCAS can fail to differentiate active from inactive haplogroups. Cascini *et*
36  *al.* demonstrate copy number variation in THCAS with qPCR but these primer
37  selections would have failed to amplify only one inactive variant haplogroup in this
38  study (Black84_Onofri-lbc1_C0_P0_NumRs92). A comprehensive sequencing strategy
39  can properly identify active from inactive copies.
40
41  **Conclusions**
42  The study of the inheritance of THCAS in modern drug type cultivars must consider
43  polyploidy. Even though many of the replication events in THCAS appear to be
44  pseudogenes or inactive forms, partially active forms are emerging and appear to be
45  associated with more the balanced (Bt:Bd) THCA:CBDA cultivars.  While de Meijer et al.
46  demonstrate (Jurka & Kapitonov, 2001)chemotype prediction from the B190/B200

1    primers for the Bt:Bd allele, its precise location in the genome is difficult to uniquely
2    resolve and while it can inform THCA and CBDA chemotype it has not been able to
3    predict the magnitude of THCA or CBDA expression and some ambiguity on Bt:Bd vs
4    Bd:Bd genotyping has been noticed.  The sequencing of these alleles demonstrates a
5    repetitive and polymorphic locus that appears difficult to design higher throughput
6    qPCR assays for.
7
8    The observation of diverged THCAS genes and copy number variation may provide
9    another tool in refining the chemotype prediction for *Cannabis* in light of previous
10   work.  Novel variants have been identified and phased with other variants of interest
11   in THCAS. In light of both copy number variation, amino acid variation, B190-B200
12   alleles and van Bakels AAE3 observations, it is possible that the highly selected
13   chemotype is experiencing convergent evolution. Of note are the several
14   CWCTTA/TTAGWG "P Instability Factor" (PIF) sequences present in THCAS. While it is
15   unclear if this replication event is transposon driven, the multiple, albeit small,
16   recognition sequence signals in the gene and the repetitive sequence flanking THCAS is
17   of note(Jurka & Kapitonov, 2001; Zhang *et al.*, 2001).  Structural variation in the human
18   genome is often attributed to segmental duplication of ancestral alleles followed by
19   Nonhomologous end joining (NHEJ) or non allelic homologous recombination (NAHR).
20   VNTRs and retrotransposition are also known to drive these events(Kidd *et al.*, 2008).
21   While gene duplication of THCAS is common, the exact mechanism driving it remains
22   to be determined.
23
24   **Methods**
25   Two different amplicons for THCAS were selected. Primers described by Onofri *et al*
26   were utilized but also redesigned as these previously described primers are internal to
27   the gene (5' ATG in the Forward primer is the start codon for the gene) and cannot
28   inform on start codon or N-terminal signal peptide sequence. To multiplex multiple
29   cultivars per SMRT cell, 5 prime dual 16 base pair molecular barcodes were utilized as
30   described by the vendor.
31
32   *THCAS PCR*
33   PCR was performed with 50ul 2x LongAmp (New England Biolabs- #M0287S), 8ul
34   10uM each primer, 42ul DNA+ddH20 (30ng total). Plant DNA was extracted with
35   SenSATIVAx according to manufacturers instructions (Medicinal Genomics part
36   #420001). DNA is eluted with 50ul ddH20.
37
38   THCAS was amplified from genomic 30ng of genomic DNA using primers Onofri-F
39   (Forward;ATGAATTGCTCAGCATTTTCCTT) and Onofri-R (Reverse;
40   ATGATGATGCGGTGGAAGA). Product was cycled using 98C 30 s followed by 28 cycles
41   of 98C 10 s, 48C 30 s, 72C 90 s followed by a 72C for 5 min.
42
43   THCAS was also amplified from genomic 30ng of genomic DNA using primers MGC-
44   2130-F (Forward; ATGGGAACCATAATAAACTATAAAAGTCATT) and MGC-2130-R
45   (Reverse; TATGATTGTCTACACAGTTCTAATGTAGATTTTC). Product was cycled using

1    94C 30s followed by 28 cycles of 98C 10s, 58C 30s, 65C 150s followed by a 65C for 5
2    min.
3
4    PCR products were purified with equal volume of SenSATIVAx (100ul) followed by two
5    100ul 70% EtOH washes.  Samples were eluted in 25ul ddH20 and quantified with
6    Qubit and Agilent HS chips, normalized and pooled for Pacific Biosciences circular
7    consensus sequencing. DNA barcodes recommended by the manufacturers were used
8    to multiplex sequencing on each SMRT cell. Sequencing was performed by University
9    of Florida service center.
10
11    **Nextera Libraries for Illumina phase confirmation.**
12    To generate sized shotgun libraries 9ul of THCAS PCR product (1.7ng/ul) was added to
13    20ul Nextera Buffer with 0.5ul Nextera Enzyme and 0.5ul ddH20. Nextera reaction was
14    driven 7 minutes at 51C. Reaction was quickly purified with 40ul of SenSATIVAx. Tip
15    mixes 10 times and placed on a magnet plate (Medicinal Genomics part#420202) for
16    10 minutes, removed supernatent and washed with 100ul 70% EtOH twice. After 5
17    minutes of benchtop room temp drying of residual EtOH the samples were eluted with
18    25ul ddH20.
19
20    **Post Nextera end healing or NT-PCR**.
21    To end repair the transposition event a nick translation event was performed prior to
22    initial denaturization (NT-PCR). 31ul of LongAmp 2X master mix (New England
23    Biolabs- #M0287S) was added to 6ul primer+indicies (2ul Index 1, 2ul Index 2, 2ul
24    ILMN primers). 25ul Purified Nextera library DNA (62ul total reaction volume) was
25    used for NT-PCR.
26
27    Reactions were cycled with 72C for 3 minutes, 98C for 30 seconds, 98C for 10 seconds,
28    63C for 30s, 72C for 1minute, Goto step 3 for 14cycles, 10C-hold. After NT-PCR, the
29    sample was purified with 60ul SenSATIVAx added to the 62ul NT-PCR reaction.
30    SenSATIVAx was tip mixed 10 times and placed on a magnet plate for 10 minutes. After
31    separation and supernatant removal the sample was washed twice with 100ul of 70%
32    EtOH. The beads were dried on a magnet plate for 10 minutes and eluted in 25ul
33    ddH20
34
35    **SAGE size separation**
36    25ul sample was added to 10ul R2 Marker (Sage Sciences) and loaded into a Sage
37    Science Blue Pippin, 1.5% Agarose Dye-Free 40ul EM. Software was set to size cut at
38    600-800bp in size. 40-60ul of elution was captured in the output port. Optionally QC
39    1ul on an Agilent 2100 Bioanalyzer High Sensitivity chips. Loaded library according to
40    Illumina V2 instructions protocol with 2x250bp reads.
41
42    **Bt:Bd Allele Amplification and sequencing**
43    Primers described by de Meijer *et al.* were used named B190F(Forward; B190F
44    TGCTCTGCCCAAAGTATCAA) and B200R (Reverse; CCACTCACCACTCCACCTTT).
45
46    Using 5ng input DNA was added to 25uL Q5 polymerase 2X Master Mix (NEB #).

8

1  1.25uL B190F Primer (10uM) with 1.25uL B200R Primer (10uM) were mixed with
2  22.5uL Water. Total volume was brought to 50ul with water + DNA (5ng total).
3
4  *Genomic amplification*
5  PCR reactions were heated to 98C for 30 seconds and 40 cycles of 98C for 10s, 56C for
6  30s, 72C for 30s. A final 72C for 5 minutes was performed before 8C hold. PCR
7  products were purified with 75ul SenSATIVAx. After two 100ul 70% EtOH washes, the
8  beads were dried and eluted in 25ul ddH20. Yield was measured with a Qubit and an
9  Agilent HS chip.
10
11  *PCR product Cloning*
12  DNA libraries were constructed with 250ng DNA using NEB's NEBNext Quick ligation
13  module (NEB # E6056S). End Repair used 3ul of Enzyme Mix, 6.5ul of Reagent Mix,
14  55.5ul of DNA + ddH20. After End Repair, Ligation was performed directly with 15ul of
15  Blunt End TA Mix, 2.5ul of Ilumina Adaptor (10uM) and 1ul of Ligation enhancer
16  (assumed to be 20% PEG 6000). After 15-minute ligation at 25C, 3ul of USER enzyme
17  was added to digest the hairpin adaptors and prepare for PCR. The USER enzyme was
18  tipmixed and incubated at 37C for 20 minutes. After USER digestion, 86.5ul of
19  SenSATIVAx was added and mixed. The samples were placed on a magnet for 15
20  minutes until the beads cleared and the supernatent could be removed. Beads were
21  washed twice with 150ul of 70% EtOH. Beads were left for 10 minute to air dry and
22  then eluted in 25ul of 10mM Tris-HCl.
23
24  *Amplification of libraries*
25  Samples adapted with NEB adaptors were amplified with 25ul of 2X Q5 Polymerase
26  (NEB). 1ul of 25uM i7 Primer and 1ul Universal primer were added. 23ul of Eluted
27  DNA was added to make a 50ul total reaction. 98C for 30s was used to hot start the
28  enzyme, followed by 6 cycles of 98C for 10s, 65C for 30s and 72C for 30s. A final 72C
29  for 5 minutes was performed before 8C hold. PCR products were then purified using
30  equal volumes (50ul) of SenSATIVAx. Samples were run on Qubit and Agilent to gauge
31  their quality and further size selected on a SAGE Blue Pippin 1.2% Agarose dye free
32  system to remove off target amplicons. (100-400bp). Quantified libraries were run on
33  a MiSeq with V2 chemistry using 2x250bp reads.
34
35  *Analysis*
36  Circular Consensus was achieved using Pacific Biosciences tools for de-multiplexing
37  symmetric barcodes with the RS Long Amplicon Analysis.1 using SMRTanalysis 2.3.0
38  workflows (http://www.pacb.com/devnet/). Minimum size was selected as 1500 and
39  2300 for Onofri and MGC-2130 amplicons respectively with a maximum subread of
40  700. Fasta files were aligned with a CLCbio workstation and confirmed using BLAST to
41  active and inactive accessions numbers (Q8GTB6.1 and Q33DQ2.1). Phylogenetic trees
42  were graphed with CLC bio with E33090 set as the root.
43
44  Illumina reads were mapped with BWA-MEM version 0.7.12-r1044 to the E33090 with
45  600-650bp inserts selected. These reads were subsequently mapped to each PacBio
46  derived haplogroup for each respective cultivar. Reads mapped to PacBio haplogroup

1    references were displayed in the Integrated Genome Viewer or IGV(Robinson *et al.*,
2    2011). Genbank accession numbers for this manuscript: SRP064442
3
4    **Author Contributions**
5    KJM- Experimental design, barcoded PCR, Nextera, Sequencing and Manuscript
6    drafting
7    YH- Barcoded PCR, Nextera, Sequencing, B190, B200 tests
8    VT- THCAS primer design and PacBio software installment.
9    SM- Illumina alignments to phased PacBio data and Figure generation
10   JS- Experimental optimizations and method development
11   LZ- Experimental optimizations and method development
12   DS- Manuscript review and drafting.
13
14   **Acknowledgments**

17
18

19   **Cascini F, Passerotti S, Boschi I. 2013.** Analysis of THCA synthase gene expression in
20           cannabis: a preliminary study by real-time quantitative PCR. *Forensic Sci Int*
21           **231**(1-3): 208-212.
22   **Cascini F, Passerotti S, Martello S. 2012.** A real-time PCR assay for the relative
23           quantification of the tetrahydrocannabinolic acid (THCA) synthase gene in
24           herbal Cannabis samples. *Forensic Sci Int* **217**(1-3): 134-138.
25   **Chaisson MJ, Huddleston J, Dennis MY, Sudmant PH, Malig M, Hormozdiari F,**
26           **Antonacci F, Surti U, Sandstrom R, Boitano M, Landolin JM,**
27           **Stamatoyannopoulos JA, Hunkapiller MW, Korlach J, Eichler EE. 2015.**
28           Resolving the complexity of the human genome using single-molecule
29           sequencing. *Nature* **517**(7536): 608-611.
30   **de Meijer EP, Bagatta M, Carboni A, Crucitti P, Moliterni VM, Ranalli P, Mandolino**
31           **G. 2003.** The inheritance of chemical phenotype in Cannabis sativa L. *Genetics*
32           **163**(1): 335-346.
33   **Eid J, Fehr A, Gray J, Luong K, Lyle J, Otto G, Peluso P, Rank D, Baybayan P,**
34           **Bettman B, Bibillo A, Bjornson K, Chaudhuri B, Christians F, Cicero R, Clark**
35           **S, Dalal R, Dewinter A, Dixon J, Foquet M, Gaertner A, Hardenbol P, Heiner**
36           **C, Hester K, Holden D, Kearns G, Kong X, Kuse R, Lacroix Y, Lin S, Lundquist**
37           **P, Ma C, Marks P, Maxham M, Murphy D, Park I, Pham T, Phillips M, Roy J,**
38           **Sebra R, Shen G, Sorenson J, Tomaney A, Travers K, Trulson M, Vieceli J,**
39           **Wegener J, Wu D, Yang A, Zaccarin D, Zhao P, Zhong F, Korlach J, Turner S.**
40           **2009.** Real-time DNA sequencing from single polymerase molecules. *Science*
41           **323**(5910): 133-138.
42   **Jurka J, Kapitonov VV. 2001.** PIFs meet Tourists and Harbingers: a superfamily
43           reunion. *Proc Natl Acad Sci U S A* **98**(22): 12315-12316.
44   **Kidd JM, Cooper GM, Donahue WF, Hayden HS, Sampas N, Graves T, Hansen N,**
45           **Teague B, Alkan C, Antonacci F, Haugen E, Zerr T, Yamada NA, Tsang P,**
46           **Newman TL, Tuzun E, Cheng Z, Ebling HM, Tusneem N, David R, Gillett W,**

1  **Phelps KA, Weaver M, Saranga D, Brand A, Tao W, Gustafson E, McKernan**
2  **K, Chen L, Malig M, Smith JD, Korn JM, McCarroll SA, Altshuler DA, Peiffer**
3  **DA, Dorschner M, Stamatoyannopoulos J, Schwartz D, Nickerson DA,**
4  **Mullikin JC, Wilson RK, Bruhn L, Olson MV, Kaul R, Smith DR, Eichler EE.**
5  **2008.** Mapping and sequencing of structural variation from eight human
6  genomes. *Nature* **453**(7191): 56-64.
7  **Kojoma M, Seki H, Yoshida S, Muranaka T. 2006.** DNA polymorphisms in the
8  tetrahydrocannabinolic acid (THCA) synthase gene in "drug-type" and "fiber-
9  type" Cannabis sativa L. *Forensic Sci Int* **159**(2-3): 132-140.
10 **Mandolino C. 1999.** Identification of DNA markers linked to the male sex in dioecious
11 hemp
12 (Cannabis sativa L.). *Theor Appl Genet (1999) 98: 86Ð92.*
13 **McKernan KJ. https://archive.org/details/SequencingTheCannabisGenome)**
14 Sequencing The Cannabis Genome.
15 **McKernan KJ. https://aws.amazon.com/datasets/the-cannabis-sativa-genome/.**
16 The Cannabis Sativa Genome.
17 **Metzker ML. 2010.** Sequencing technologies - the next generation. *Nat Rev Genet*
18 **11**(1): 31-46.
19 **Nakamura K, Oshima T, Morimoto T, Ikeda S, Yoshikawa H, Shiwa Y, Ishikawa S,**
20 **Linak MC, Hirai A, Takahashi H, Altaf-Ul-Amin M, Ogasawara N, Kanaya S.**
21 **2011.** Sequence-specific error profile of Illumina sequencers. *Nucleic Acids Res*
22 **39**(13): e90.
23 **Onofri C, de Meijer EP, Mandolino G. 2015.** Sequence heterogeneity of cannabidiolic-
24 and tetrahydrocannabinolic acid-synthase in Cannabis sativa L. and its
25 relationship with chemical phenotype. *Phytochemistry* **116**: 57-68.
26 **Pacifico D. 2006.** Genetics and marker-assisted selection of the chemotype
27 in Cannabis sativa L. *Molecular Breeding* **17**: 257–268.
28 **Robinson JT, Thorvaldsdottir H, Winckler W, Guttman M, Lander ES, Getz G,**
29 **Mesirov JP. 2011.** Integrative genomics viewer. *Nat Biotechnol* **29**(1): 24-26.
30 **Shoyama Y, Tamada T, Kurihara K, Takeuchi A, Taura F, Arai S, Blaber M,**
31 **Shoyama Y, Morimoto S, Kuroki R. 2012.** Structure and function of 1-
32 tetrahydrocannabinolic acid (THCA) synthase, the enzyme controlling the
33 psychoactivity of Cannabis sativa. *J Mol Biol* **423**(1): 96-105.
34 **Sirikantaramas S, Morimoto S, Shoyama Y, Ishikawa Y, Wada Y, Shoyama Y, Taura**
35 **F. 2004.** The gene controlling marijuana psychoactivity: molecular cloning and
36 heterologous expression of Delta1-tetrahydrocannabinolic acid synthase from
37 Cannabis sativa L. *J Biol Chem* **279**(38): 39767-39774.
38 **van Bakel H, Stout JM, Cote AG, Tallon CM, Sharpe AG, Hughes TR, Page JE. 2011.**
39 The draft genome and transcriptome of Cannabis sativa. *Genome Biol* **12**(10):
40 R102.
41 **Weiblen GD, Wenger JP, Craft KJ, ElSohly MA, Mehmedic Z, Treiber EL, Marks MD.**
42 **2015.** Gene duplication and divergence affecting drug content in Cannabis
43 sativa. *New Phytol*.
44 **Zhang X, Feschotte C, Zhang Q, Jiang N, Eggleston WB, Wessler SR. 2001.** P
45 instability factor: an active maize transposon system associated with the

1    amplification of Tourist-like MITEs and a new superfamily of transposases. *Proc*
2    *Natl Acad Sci U S A* **98**(22): 12572-12577.
3
4
5
6
7
8
9
10
11

| Cultivar | Source | Chemotype | B190-B200 |
|---|---|---|---|
| Black84_THCA | Uncle Ikes, Seattle, WA | THCA | Bt:Bd |
| BlueBerry Essence_CBDA_THCA | Cannabis City, Seattle,WA | 8.11%THCA:10.8%CBDA | Bd:Bd |
| Sour Tsunami_CBDA | Werc Shop, Seattle,WA | 11.3% CBDA, 0.5%THCA | Bd:Bd |
| BlueDream_THCA | MCR labs, MA | THCA | Bd:Bd |
| Otto_CBDA | Centennial Seeds, CO | 11.3% CBDA, 0.5%THCA | Bd:Bd |
| WZ-CBD_CBDA | Rocky Mountain Remedies, CO | 20% CBDA, 05%THCA | Bd:Bd |
| WIFI_THCA | MCR labs, MA | THCA | Bt:Bt |
| ChemDog91_THCA | Greenhouse Seeds, NL | THCA | Bt:Bt |
| CannaTsu_CBDA | Salmon Creek, CA | 11.4% CBDA:0.4%THCA | Bd:Bd |
| C4 X CannaTsu | Salmon Creek, CA | Unknown | Bd:Bd |
| C4_THCA | Salmon Creek, CA | THCA | Bt:Bt |
| AK-47_THCA | Salmon Creek, CA | THCA | Bd:Bd |
| Alaskan Ice_THCA | MCR labs, MA | THCA | Bt:Bt |
| Chemdog91_THCA repeat | Greenhouse Seeds, NL | THCA | Bt:Bt |
| Chemdog91_SinglePlex | Greenhouse Seeds, NL | THCA | Bt:Bt |

**Table1**. THCAS was amplified with two primer sets (Onofri and MGC-2130) for each Cultivar and single molecule sequenced with Pacific BioSciences SMRT cell circular consensus sequencing. B190, B200 alleles were measured according to Mandolino *et al.* Red highlighted samples represent B190-B200 Genotype-Chemotype discordant samples.

**Samples Run on 3 different chips**

| Singleplex ChemDog91 | SubRead Coverage (#ZMWs) | Chemotype | PredictedAccuracy | ConsensusConverged | Activity | Reference_AAVariant |
|---|---|---|---|---|---|---|
| Chemdog91_1plex_Onofri_lbc0_C0_P2_NumReads320_Q33DQ2.1_S221F | 320 | THCA | 0.9999369 | TRUE | Unknown | Q33DQ2.1_F221S |
| Chemdog91_1plex_Onofri_lbc0_C1_P1_NumReads312_Q8GTB6.1_A250D | 312 | THCA | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1_A250D |
| Chemdog91_1plex_Onofri_lbc0_C1_P0_NumReads188_Q8GTB6.1 | 188 | THCA | 0.9999369 | TRUE | Active | Q8GTB6.1 |
| Chemdog91_1plex_Onofri_lbc0_C0_P0_NumReads93_Q33DQ2.1_6AA | 93 | THCA | 0.9996669 | TRUE | Unknown | Q33DQ2.1_N90F,I266T,L370F,D420G,Y471C,T492S |
| Chemdog91_1plex_Onofri_lbc0_C0_P1_NumReads87_Q33DQ2.1_STOP_438 | | | 0.9936923 | TRUE | Stop/Frame-Shift | Q33DQ2.1_M134I_STOP438 |

| FR09904480 | SubRead Coverage (#ZMWs) | Chemotype | PredictedAccuracy | ConsensusConverged | Activity | Reference_AAVariant |
|---|---|---|---|---|---|---|
| ChemDog91_lbc10_NumReads318_Q33DQ2.1 | 318 | THCA | 0.9999369 | TRUE | Inactive | Q33DQ2.1 |
| ChemDog91_lbc10_NumReads307_Q8GTB6.1_A250D | 307 | THCA | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1-A250D |
| ChemDog91_lbc10_NumReads193_Q8GTB6.1 | 193 | THCA | 0.999931 | TRUE | Active | Q8GTB6.1 |
| ChemDog91_lbc10_NumReads95_Q33DQ2.1_M134I_STOP438 | 95 | THCA | 0.992139 | TRUE | Stop/Frame-Shift | Q33DQ2.1_M134I_STOP438 |
| ChemDog91_lbc10_NumReads87_Q33DQ2.1_6AA | 87 | THCA | 0.9921376 | TRUE | Unknown | Q33DQ2.1_N90F,I266T,L370F,D420G,Y471C,T492S |

| FR09904616 | SubRead Coverage (#ZMWs) | Chemotype | PredictedAccuracy | ConsensusConverged | Activity | Reference_AAVariant |
|---|---|---|---|---|---|---|
| ChemDog91-lbc7_C0_P1_NumReads299_Q8GTB6.1_A250D | 299 | THCA | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1_A250D |
| ChemDog91-lbc7_C0_P0_NumReads201_Q8GTB6.1 | 201 | THCA | 0.9999369 | TRUE | Active | Q8GTB6.1 |
| ChemDog91-lbc7_C1_P2_NumReads103_FS | 103 | THCA | 0.9963401 | TRUE | Stop/Frame-Shift | Frame Shift |
| ChemDog91-lbc7_C1_P0_NumReads103_Q33DQ2.1_S221F | 103 | THCA | 0.9981365 | FALSE | | Q33DQ2.1_F221S |
| ChemDog91-lbc7_C1_P1_NumReads75_Q33DQ2.1_6AA | 75 | THCA | 0.9911302 | TRUE | Unknown | N90F,I266T,L370F,D420G,Y471C,T492S |
| ChemDog91-lbc7_C1_P3_NumReads67_Q33DQ2.1_M134I_STOP438 | 67 | THCA | 0.9999369 | TRUE | Stop/Frame-Shift | Q33DQ2.1_M134I_STOP438 |

| Active (Q8GTB6.1) |
|---|
| Inactive (Q33DQ2.1) |
| Unknown |
| Stop FrameShift |

**Table2.** Triplicate Sample analysis. ChemDog91 was amplified 3 distinct times and sequenced on 3 different SMRT cells to assess reproducibility and sampling. Active forms consistently amplify and sequence across all amplification products. A few private inactive or unknown activity haplogroups amplify inconsistently (Q33DQ2.1- Red and Q33DQ2.1_F221S- Green).

| Sample | Coverage | Chemotype | Location | Notes | Transcript_AAVariant |
|---|---|---|---|---|---|
| Black84 lbc1_C0_P1_NumReads408_Q33DQ2.1 | 408 | THCA | Seattle | | Q33DQ2.1 |
| Black84 lbc1_C1_P1_NumReads333_Q8GTB6.1_A250D | 333 | THCA | | AK-47 implies this is active | Q8GTB6.1_A250D |
| Black84 lbc1_C1_P0_NumReads167_Q8GTB6.1_I63L | 167 | THCA | | 2010.24a.35/T-Onofri has an I | Q8GTB6.1_I63L |
| Black84 lbc1_C0_P0_NumReads92_STOP | 92 | | | 2010.24a.35/T | Stop_Frameshift |
| Blue Essence_lbc2_C0_P0_NumReads500_Q8GTB6.1 | >500 | THCA:CBDA | Seattle | | Q8GTB6.1_P333R |
| SourTsunami_lbc5_C1_P0_NumReads500_Q8GTB6.1 | >500 | 11.3% CBDA, 0.5%THCA | Seattle | | Q8GTB6.1 |
| SourTsunami_lbc5_C0_P2_NumReads218_Q33DQ2.1 | 218 | 11.3% CBDA, 0.5%THCA | | | Q33DQ2.1 |
| SourTsunami_lbc5_C0_P0_NumReads117_STOP | 117 | 11.3% CBDA, 0.5%THCA | | | Stop_Frameshift |
| SourTsunami_lbc5_C0_P1_NumReads103_Q33DQ2.1_6AA | 103 | 11.3% CBDA, 0.5%THCA | | | Q33DQ2.1_F90V,I266T,L370F,G420D,Y472C,T491S |
| SourTsunami_lbc5_C0_P3_NumReads62_STOP | 62 | 11.3% CBDA, 0.5%THCA | | | Stop_Frameshift |
| SourTsunami_lbc5_C2_P0_NumReads23_STOP | 23 | 11.3% CBDA, 0.5%THCA | | | Stop_Frameshift |
| BlueDream_lbc6_C1_P0_NumReads500_Q8GTB6.1 | 500 | THCA | Mass | | Q8GTB6.1 |
| BlueDream_lbc6_214_C0_P3_NumReads214_Q33DQ2.1 | 24 | THCA | | | Q33DQ2.1 |
| BlueDream_lbc6_C0_P2_NumReads102_STOP | 102 | THCA | | | Stop_Frameshift |
| BlueDream_lbc6_C0_P0_NumReads92_STOP | 92 | THCA | | | Stop_Frameshift |
| BlueDream_lbc6_C0_P1_NumReads92b_Q33DQ2.1_6AA | 92 | THCA | | | Q33DQ2.1_V90F,I266T,L370F,G420D,Y472C,T491S |
| Otto_lbc7_C0_P0_NumReads500_Q33DQ2.1 | 500 | CBDA | Centennial Seeds | | Q33DQ2.1 reference |
| Otto_lbc7_C1_P1_NumReads500b_Q8GTB6.1_A411V | 500 | CBDA | | Active or not? | Q8GTB6.1_A411V |
| WZ_CBD_lbc8_C0_P2_NumReads335_Q33DQ2.1 | 335 | CBDA | Seattle | | Q33DQ2.1 |
| WZ_CBD_lbc8_C0_P1_NumReads86_STOP | 86 | CDBA | | | Stop_Frameshift |
| WZ_CBD_lbc8_C0_P1_NumReads79_Q33DQ2.1_6AA | 79 | CDBA | | | Q33DQ2.1_F90V,I266T,L370F,G420D,Y472C,T491S |
| WIFI_lbc9_C0_P1_NumReads290_Q8GTB6.1 | 290 | THCA | Mass male | | Q8GTB6.1 |
| WIFI_lbc9_C0_P0_NumReads210_Q8GTB6.1_V125L | 210 | THCA | | | Q8GTB6.1_V125L |
| WIFI_lbc9_C1_P1_NumReads205_Q33DQ2.1 | 205 | THCA | | | Q33DQ2.1 |
| WIFI_lbc9_C1_P0_NumReads119_Q33DQ2.1_N22S_P81S | 119 | THCA | | | Q33DQ2.1_N22S,P81S |
| Canna-Tsu-lbc1_C0_P0_NumReads359_Q8GTB6.1 | 359 | CBDA | Salmon Creek, CA | | Q8GTB6.1_A250D, |
| Canna-Tsu-lbc1_C1_P2_NumReads146_Q33DQ2.1_A430T | 146 | CBDA | | | Q33DQ2.1_A430T |
| Canna-Tsu-lbc1_C1_P1_NumReads81_Q33DQ2.1_S221F | 81 | CBDA | | | Q33DQ2.1_S221F |
| Canna-Tsu-lbc1_C1_P0_NumReads55_STOP | 55 | CBDA | | | Stop_Frameshift |
| C4-lbc2_C0_P0_NumReads500_Q8GTB6.1 | 500 | THCA | Salmon Creek, CA | | Q8GTB6.1 |
| C4-lbc2_C1_P1_NumReads96_Q33DQ2.1_N22S_P81S | 96 | THCA | | | Q33DQ2.1_N22S, P81S |
| C4-lbc2_C1_P2_NumReads92_Q33DQ2.1 | 92 | THCA | | | Q33DQ2.1 |
| C4-lbc2_C1_P0_NumReads91_Q33DQ2.1_A375V | 91 | THCA | | | Q33DQ2.1_A375V, |
| C4xCanna_tsu-lbc3_C0_P1_NumReads270_Q8GTB6.1_P333R | 270 | Unknown | Salmon Creek, CA | Not direct cross | Q8GTB6.1_P333R |
| C4xCanna_tsu-lbc3_C0_P0_NumReads230_Q8GTB6.1 | 230 | Unknown | | from C4 | Q8GTB6.1 |
| C4xCanna_tsu-FS-lbc3_C1_P0_NumReads47_STOP | 47 | Unknown | qual flag | from CannaTsu | Stop_Frameshift |
| C4xCanna_tsu-lbc3_C2_P0_NumReads20_Q33DQ2.1 | 20 | Unknown | | from C4 | Q33DQ2.1 |
| AK47-lbc4_C0_P0_NumReads500_Q8GTB6.1_A250D | 500 | THCA | Salmon Creek, CA | Deduced Active | Q8GTB6.1_A250D, |
| AK47-lbc4_C1_P1_NumReads400_Q33DQ2.1 | 400 | THCA | | | Q33DQ2.1 |
| AK47-lbc4_C1_P0_NumReads100_STOP | | THCA | | | Stop_Frameshift |
| Alaskan Ice-lbc6_C0_P0_NumReads184_Q8GTB6.1 | 184 | THCA | Mass | | Q8GTB6.1 |
| Alaskan Ice-lbc6_C0_P1_NumReads183_Q8GTB6.1_V125L_G410E | 183 | THCA | | | Q8GTB6.1_V125L,G410E |
| Alaskan Ice-lbc6_C0_P4_NumReads82_Q33DQ2.1 | 82 | THCA | | | Q33DQ2.1 |

| | |
|---|---|
| Q33DQ2.1= Inactive | |
| Somewhere in between | |
| Q8GTB6.1 = Active | |

**Table 3.** List of all haplogroups discovered from each cultivar other than Chemdog91 (Table 2). Red highlighting are Inactive haplogroups. Green highlightings are unknown haplogroups and blue highlightings are active haplogroups.

| Haplogroup Clusters | SubRead Coverage (ZMW) | Chemotype | PredictedAccuracy | ConsensusConverged | Activity | Reference_AAVariant |
|---|---|---|---|---|---|---|
| Black84_2130-lbc11_C0_P0_NumReads500 | >500 | THCA | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1_A250D |
| Blueberry Essence_2130-lbc12_C0_P0_NumReads500 | >500 | 8.11%THCA:10.8%CBDA | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1_P333R |
| C4_SSC_2130-lbc14_C0_P0_NumReads500 | >500 | THCA | 0.9999369 | TRUE | Active | Q8GTB6.1 |
| C4 x CannaTsu_2130-lbc15_C0_P1_NumReads265 | 265 | Unknown | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1_P333R |
| C4 x CannaTsu_2130-lbc15_C0_P0_NumReads235 | 235 | Unknown | 0.9999369 | TRUE | Active | Q8GTB6.1 |
| WIFI_2130-lbc16_C0_P0_NumReads500 | >500 | THCA | 0.9999369 | TRUE | Active | Q8GTB6.1 |
| ChemDog91-2130_lbc17_C0_P1_NumReads334 | 334 | THCA | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1_A250D |
| Chemdog91-2130_lbc17_C0_P0_NumReads166 | 166 | THCA | 0.9999369 | TRUE | Active | Q8GTB6.1 |
| AK47_SCC-lbc18_C0_P0_NumReads500 | >500 | THCA | 0.9999369 | TRUE | Active by deduction | Q8GTB6.1_A250D |
| Alaskan Ice_2130-lbc19_C0_P0_NumReads500 | >500 | THCA | 0.9999369 | TRUE | Active | Q8GTB6.1 |
| Otto_2130-lbc20_C0_P0_NumReads500 | >500 | 11.4% CBDA:0.4%THCA | 0.9999369 | TRUE | Low Activity by deduct | Q8GTB6.1_A411V |

**Table4**- MGC-2130 primers only amplify active alleles. Active transcript Q8GTB6.1 is the most common haplogroup sequenced. A250D is the second most common haplogroup and is the only haplogroup present in a high THCA expressing cultivar AK-47 and is thus labeled as "Active by deduction". P333R is also present in Blueberry Essence with 8.11% THCA suggesting Active by deduction Classification. A411V may be low activity or in a cultivar with Bd:Bd genotype.
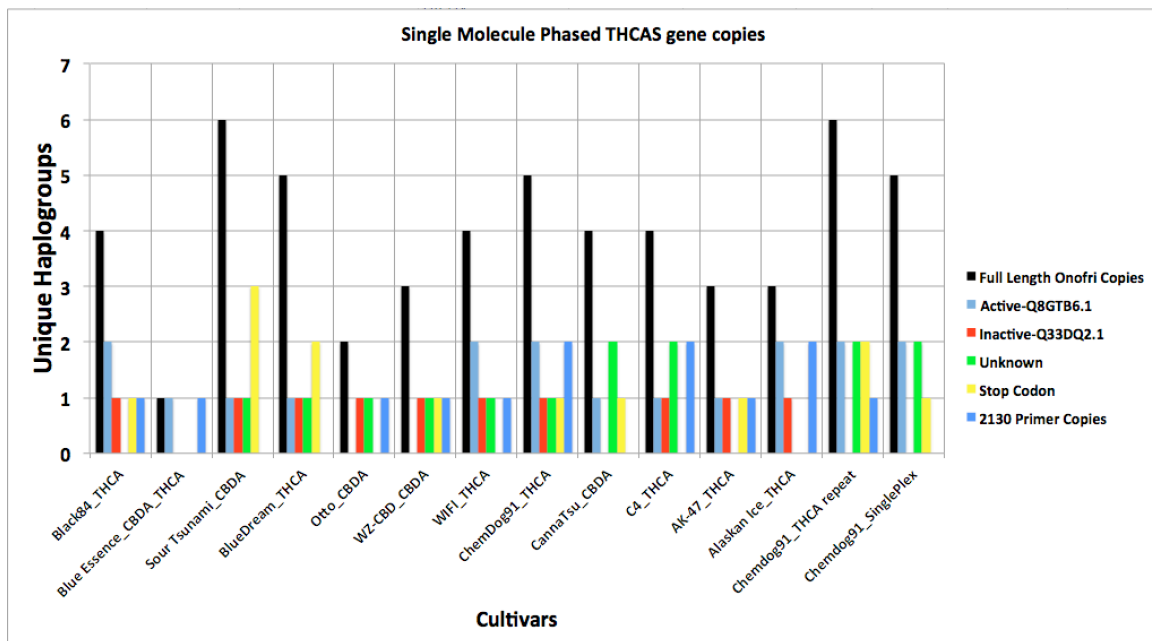
**Figure 1**. Unique full length THCA Synthase sequences for each cultivar demonstrate diploidy is a rare event in Drug Type cultivars selected for high cannabinoid content. All THCA positive strains have an active form of THCAS with the exception of AK-47 that has an A250D variant of the active form Q8GTB6.1. A250D variants are also present in Black84 as the only possible functional THCAS and as the 2nd THCAS haplogroup in Chemdog91. While it is also present in a CannaTsu high CBDA cultivar, this cultivar does not amplify with the MGC-2130 primers suggesting an MGC-2130 A250D haplogroup is in fact an active THCAS when full length.

**Figure 2**. Mapping 2x250 Illumina reads to all of the Phased Pacific Bioscience generated haplogroups compared to mapping them to a single THCAS haplogroup sequence as you might find in CanSat3. 5 haplogroups of Chemdog91 shown in IGV with Illumina reads mapped to a refernce containing all 5 references. The center bottom view "All Mapped to Haplogroup3" represents all of the reads mapped to just Haplogroup 3 demonstrating collapsed polymorphic reads when mapped to the unphased reference sequence of a single haplogroup. Most reads in IGV are labeled with poor mapping scores due to homology between the haplogroups. IGV also filters out alleles under 20% frequency inducing a more detailed dissection of this across more samples in Figure 5.

## Figure 3. THCA Synthase Annotation.

MGC-2130 Primers bracket the displayed region of interest. Onofri and Weiblen primers are in green. THCAS ORF is in yellow. Various nonsynonmous SNPs (I63L, V125L, A250D, V287M, M289T, P333R, A411V) found in Active THCAS haplogroups are annotated near respective functional groups like Glycosylation sites, FAD binding domains and mPIF signals. Cascini primers and BBE Pfam domains are labeled accordingly.
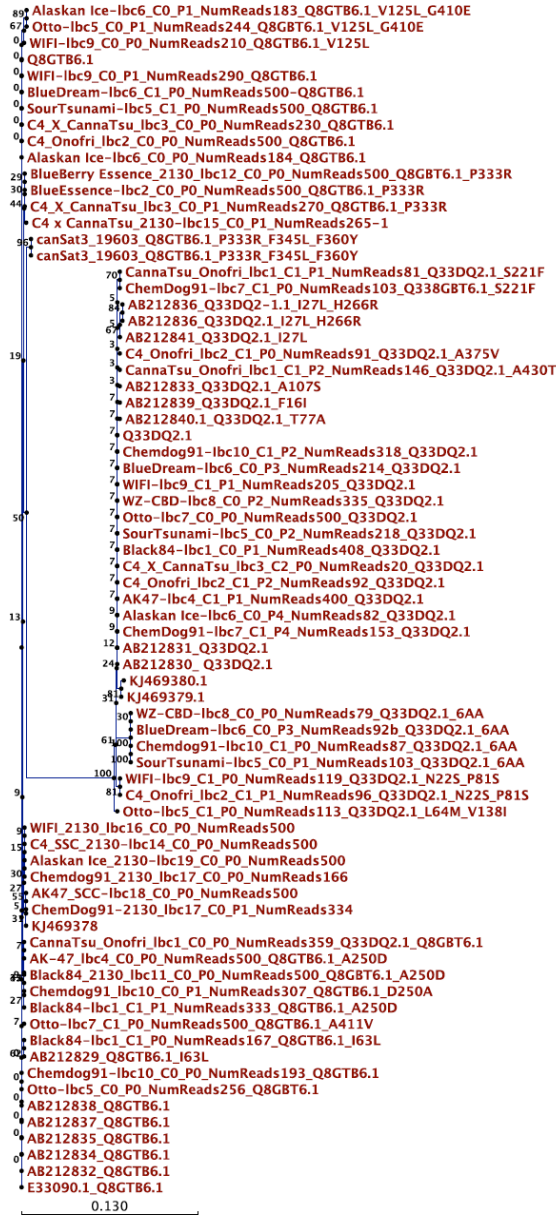
**Figure 4.** Phylogenetic tree of amino acid sequences of THCAS using E33090 as the root. Long pseudogenes were omitted from the tree since they had excessive divergence. The NumReads reflects the number of subreads that support a haplogroup after circular consensus has been achieved. Most haplogroups cluster around Active and Inactive like accessions (Q8GTB6.1 & Q33DQ2.1 accessions). Only Q8GTB6.1 cluster amplify with MGC-2130 primers. Other GenBank Accessions are label as AB#.
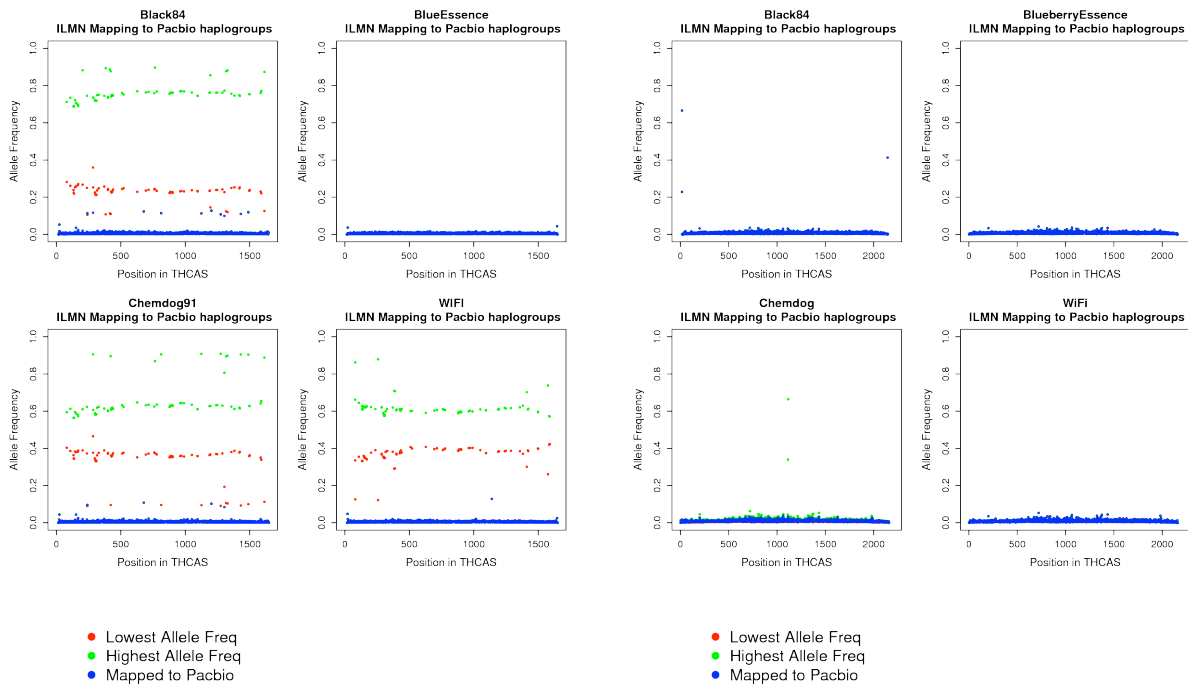
**Figure 5.** Over 100,000X coverage with Illumina 2x250bp reads derived from Onofri and MGC-2130 primer based amplification of THCAS. The Illumina data is mapped to the Pacbio FASTA files derived from Pacbio LAA analysis. The strain sequence is respectively mapped- for example- *Black84* Illumina Data is mapped to the *Black84* Pacbio FASTA. As an example, if there are 5 haplogroups, the data is mapped 6 times; once to each separate haplogroup and once to all 5 together in one FASTA reference. An allele count is performed on each position of THCAS. A clustal alignment is performed between all the haplogroups so we are able to compare the same position across the multiple reference/bam files. The minimum and maximum non-reference ratio at each position (~1600) is calculated across all the haplogroups (i.e. 5). This is the green and red dots shown. The non-reference ratio for all haplogroups mapped together is also calculated and shown in blue. Left two columns of the chart is Onofri amplification allelic coverage charts. Right 2 columns are MGC-2130 amplicon allelic coverage charts. On the X-axis is THCAS bases 1-1654bp or 1-2130bp respectively. On the Y-axis is Heterozygosity where homozygous genotypes are either 1 or 0 and heterozygous genotypes would display at 0.5. Samples Blueberry essence has been fully sampled with the PacBio depth on both primer sets (all blue homozygous alignments). Samples with red and green dots at equal allele frequency across the amplicon represent PacBio data that may have a remaining haplogroups detected at 100,000 X Illumina sequence but not yet sampled at the PacBio read depth. All 2130-MGC amplicons look fully sampled with phased active alleles.
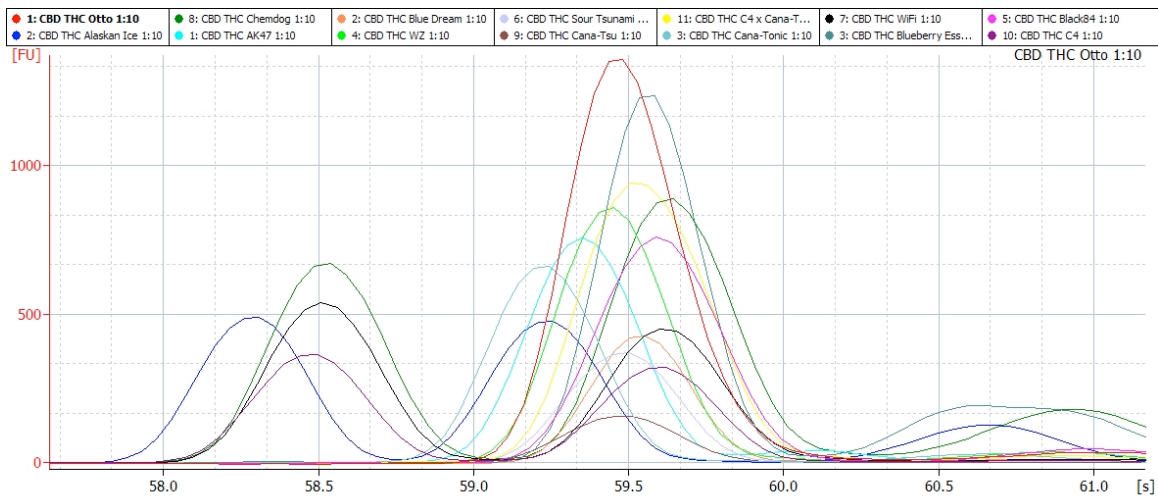
**Figure 6**. Bt:Bd allele amplified and run on an Agilent High Sensitivity chip. High THCA cultivars (WIFI, Chemdog, C4, Alaskan Ice) replicate a small ~190bp band known as Bt:Bt. High CBDA alleles generate a major 200 bp and a minor 215bp band known as Bd:Bd. Hybrid Bt:Bd alleles generate similar bands as Bd:Bd but with more 215bp product. This later peak is a bit subjective to measure.
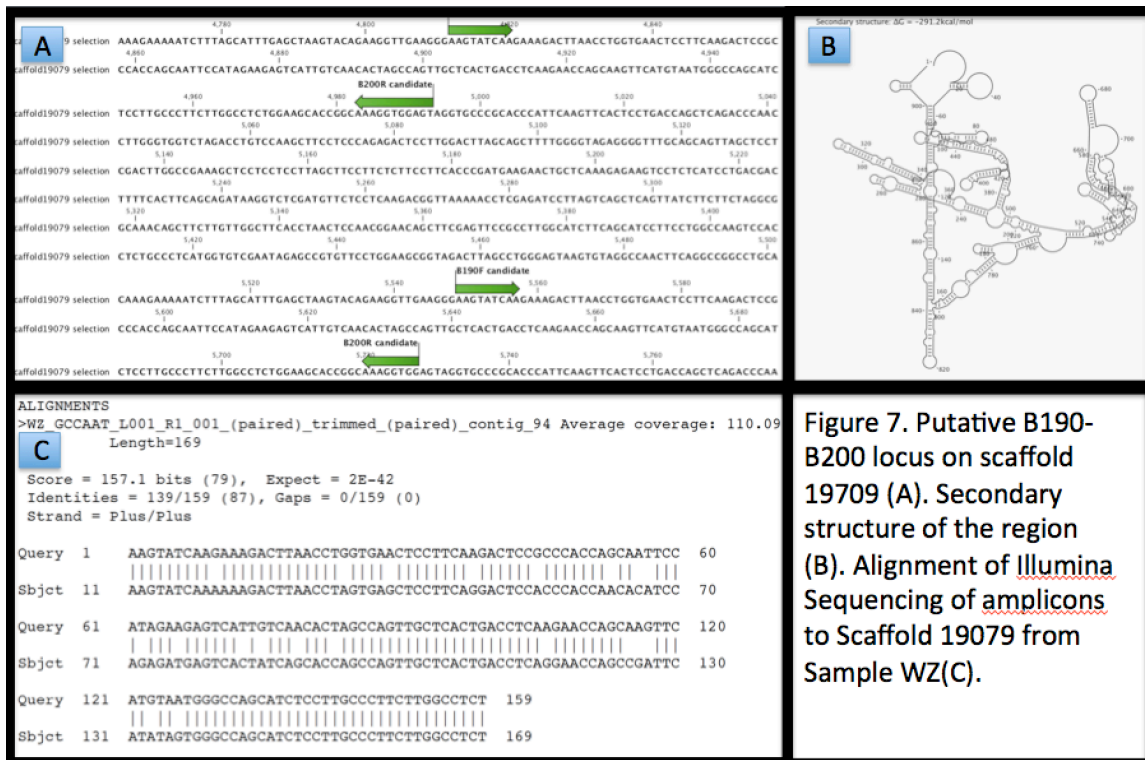
Figure 7. Putative B190-B200 locus on scaffold 19709 (A). Secondary structure of the region (B). Alignment of Illumina Sequencing of amplicons to Scaffold 19079 from Sample WZ(C).

**Figure 7**. Assembly of Illumina reads from Bt:Bd allele in Cultivar WZ with homology to the B190/B200 partial primer sequences highlighted in green arrows.