Open access • Posted Content • DOI:10.1101/2021.07.05.451212

# The Phytochemical Diversity of Commercial Cannabis in the United States
— **Source link**  ⧉

Smith Cj, Daniela Vergara, Brian Keegan, Jikomes N

**Institutions:** University of Colorado Boulder

Related papers:

- Cannabinomics: Application of Metabolomics in Cannabis (Cannabis sativa L.) Research and Development.

- Cannabis sativa research trends, challenges and new-age perspectives

- Medicinal cannabis added in food

- Professionals or Amateurs? Revisiting the Notion of Professional Crime in the Context of Cannabis Cultivation

- The Cream of the Crop: Biology, Breeding, and Applications of Cannabis Sativa

1    **The Phytochemical Diversity of Commercial *Cannabis* in the United States**

2    Christiana J. Smith[1], Daniela Vergara[2], Brian Keegan[3], Nick Jikomes[4*]

3

4    1. Seattle, Washington USA 98177

5    2. Department of Ecology and Evolutionary Biology, University of Colorado, Boulder, Colorado USA

6    80309

7    3. Department of Information Science, University of Colorado Boulder, Boulder, Colorado, USA 80309

8    4. Department of Science & Innovation, Leafly Holdings Inc, 600 1st Ave (Ste. LL20), Seattle, WA 98104

9    *Correspondence: njikomes@gmail.com

10

11   **Abstract**

12   The legal status of *Cannabis* is changing, fueling an increased diversity of *Cannabis*-derived products.
13   Because *Cannabis* contains dozens of chemical compounds with potential psychoactive or medicinal
14   effects, understanding its phytochemical diversity is crucial. The legal *Cannabis* industry heavily markets
15   products to consumers based on widely used labelling systems purported to predict the effects of different
16   *Cannabis* "strains." We analyzed the cannabinoid and terpene content of tens of thousands of commercial
17   *Cannabis* samples across six US states, finding distinct chemical phenotypes (chemotypes) which are
18   reliably present. After careful descriptive analysis of the phytochemical diversity and comparison to the
19   commercial labels commonly attached to *Cannabis* samples, we show that commercial labels do not
20   consistently align with the observed chemical diversity. However, certain labels are statistically
21   overrepresented for specific chemotypes. These results have important implications for the classification
22   of commercial *Cannabis*, the design of animal and human research, and the regulation of legal *Cannabis*
23   marketing.

24

25

26   **Introduction**

27       *Cannabis sativa* L., a flowering plant from the family Cannabacea (Clarke and Merlin 2013;

28   Clarke and Merlin 2016), is one of the oldest domesticated plants (Russo 2007). The plant has been used

29   by humans for more than 10,000 years (Abel 2013) and has spread throughout the globe such that, today,

30   distinct varieties exist, which have been cultivated for multiple purposes. This versatile and phenotypically

31   diverse plant has been used for a wide variety of commercial and medicinal purposes (Clarke and Merlin

32   2013). The *Cannabis* genus is considered to have a single species, *Cannabis sativa* L (Watts 2006),

33   inclusive of all forms of hemp and marijuana, with high genomic and phenotypic variation (Vergara et al.

34   2016; Kovalchuk et al. 2020) across multiple lineages (Sawler et al. 2015; Lynch et al. 2016; Vergara et

35   al. 2016). 'Marijuana-type' lineages are used for human consumption (recreational and medical), while

36   the 'hemp' lineage is used in industry settings for fiber or oil extraction.

37       For human consumption, the mature female inflorescences are grown, harvested and processed

38   into dried plant material commonly called "marijuana", "weed," "flower," or other informal names. New

39  laws leading to decriminalization and legalization have given rise to a global, multibillion dollar industry

40  that is projected to continue to grow aggressively (Hutchison et al. 2019). The cannabis industry has

41  innovated across genetics, cultivation, extraction, distribution, and compliance to keep pace with the

42  demands of consumers, competitors, and regulators. Beyond dried flowers, there are concentrated oils,

43  confections and beverages, topicals, suppositories, and many other delivery mechanisms (Steigerwald et

44  al. 2018; Goodman et al. 2020). To avoid confusion with the confounding terminology (Riboulet-Zemouli

45  2020), we will use "*Cannabis*" in reference to the plant genus including its different varieties, and

46  "cannabis" as a generic term encompassing processed *Cannabis* in all forms or in reference to the cannabis

47  industry generally.

48      *Cannabis* is renowned for the production of secondary metabolites, including cannabinoids and

49  terpenes. Cannabinoids are a class of compounds that can interact with the endocannabinoid system

50  (Gertsch et al. 2008) and many have medicinal (Russo 2011; Swift et al. 2013) or psychoactive (ElSohly

51  and Slade 2005; Russo 2007) properties. Two of the most abundant cannabinoids are Δ-9-

52  tetrahydrocannabinolic acid (THCA) and cannabidiolic acid (CBDA), which are converted to the neutral

53  forms Δ-9-tetrahydrocannabinol (THC) and cannabidiol (CBD) once heated (Hart et al. 2001). The

54  enzymes responsible for the production of these cannabinoids are highly similar at the biochemical

55  structure and genetic sequence levels (Onofri et al. 2015; Vergara et al. 2019) and accept the same

56  substrate, Cannabigerolic Acid (CBGA) (Franco 2011; Chakraborty et al. 2013).

57      Beyond THC and CBD, there are various "minor cannabinoids," typically present at much lower

58  levels. This includes CBGA, the aforementioned precursor molecule to both THCA and CBDA. A third

59  compound, CBCA (cannabichrommenic acid), is also part of the same biochemical pathway that gives

60  rise to CBDA and THCA (Page and Stout 2017). Other minor cannabinoids include cannabinol (CBN), a

61  byproduct that accumulates with the breakdown of THC (Turner and Elsohly 1979; Ross and ElSohly

62  1997; Trofin et al. 2012), Δ-9-tetrahydrocannabivarin carboxylic acid (THCVA), and others. Similar to

63  THCA and CBDA, decarboxylation is responsible for the formation of cannabigerol (CBG), Δ-9-

64  tetrahydrocannabivarin (THCV), and other neutral cannabinoids (Valliere et al. 2019). Due to their low

65  abundance, these have generally been less well-studied than THC and CBD, although they display a range

66  of interesting pharmacological properties with potential medicinal value (Izzo et al. 2012; Borrelli et al.

67  2014; McPartland et al. 2015).

68      Cannabinoid levels have been used both for setting legal definitions for different categories of

69  cannabis products and for 'chemotaxonomic' purposes to classify different *Cannabis* varieties based on

70  THC:CBD ratios (Hillig and Mahlberg 2004). For example, the legal definition of hemp in the United

71  States is any *Cannabis* plant containing up to 0.3% THC. This arbitrary number intends to distinguish

72    *Cannabis* with low intoxication potential from varieties containing high THC levels. Commercial

73    marijuana-type *Cannabis* usually falls within discrete groups based on THC:CBD ratios (Hillig and

74    Mahlberg 2004), and has been categorized as either "THC-dominant" (low CBD levels), "CBD-

75    dominant," (low THC levels and high CBD levels), or "Balanced THC/CBD" (comparable levels of THC

76    and CBD), although the vast majority is THC-dominant (Jikomes and Zoorob 2018). The level of other

77    minor cannabinoids has additionally been measured in a limited number of studies (Orser et al. 2017;

78    Henry et al. 2018). However, a more comprehensive quantification of both major and minor cannabinoids

79    from a large sample representative of commercial *Cannabis*, across multiple legal markets in the United

80    States, is needed.

81        In addition to cannabinoids, *Cannabis* harbors a diverse class of related compounds known as

82    terpenes (Potter 2004, 2009). These are a type of secondary metabolite which often play defensive roles

83    for the plant (Langenheim 1994; Sirikantaramas et al. 2005). They are responsible for its odors, can be

84    pharmacologically active (McPartland and Russo 2001; ElSohly and Slade 2005), and may serve as

85    reliable chemotaxonomic markers for classifying *Cannabis* beyond THC:CBD ratios *(Orser et al. 2017;*

86    *Reimann-Philipp et al. 2019)*. It has been shown that the chemical phenotype ("chemotype") of plants can

87    be used to classify *Cannabis* into chemical varieties ("chemovars") (Hazekamp and Fischedick 2012;

88    Lewis et al. 2018). Distinct chemovars, each with different ratios of cannabinoids and terpenes, are

89    hypothesized to cause distinct effects for human consumers (Lewis et al. 2018).

90        A variety of studies have looked at the chemical composition of *Cannabis* samples limited to a

91    single geographic location (Hazekamp and Fischedick 2012; Orser et al. 2017; Henry et al. 2018;

92    Reimann-Philipp et al. 2019), included measurements of a limited number of cannabinoids (Hillig and

93    Mahlberg 2004; Elzinga et al. 2015; Hazekamp et al. 2016; Vergara et al. 2017; Jikomes and Zoorob 2018;

94    Vergara et al. 2020), or included measurements of terpenes without cannabinoid content (Hillig 2004).

95    Few studies have investigated the major and minor cannabinoids together with the terpenes (Mudge et al.

96    2019) and none have performed a thorough chemotaxonomic analysis on a dataset with tens of thousands

97    of samples across several legal cannabis markets in the United States. Mapping the chemical diversity of

98    the *Cannabis* consumed by millions of people has important implications for consumer health and safety,

99    such as identifying how many chemically distinct types of *Cannabis* are currently consumed in legal

100   markets. This may be consequential if distinct chemotypes are later determined to cause reliably different

101   effects.

102        It has been suggested that the multiple compounds produced by *Cannabis* may act in combination

103   to produce specific medicinal and psychoactive effects, the so-called 'entourage effect' (Russo 2011).

104 There is limited suggestive evidence for such an effect (McPartland and Russo 2001; Adams and Taylor
105 2010), including improved patient outcomes in those who use whole-plant extracts (containing THC and
106 unknown quantities of other compounds) versus synthetic THC (Venderová et al. 2004). For example,
107 synthetic THC alone in manufactured products such as 'Marinol' may produce unpleasant effects
108 (Calhoun et al. 1998; Carter et al. 2011). Whether or not distinct ratios of cannabinoids and terpenes are
109 able to consistently yield different subjective effects or therapeutic outcomes is unknown, and a topic of
110 debate (Russo 2019).

111 Combinatorial effects, when the ingestion of two or more compounds yields different effects from
112 either compound in isolation, may be more likely when a drug acts on multiple target systems
113 (polypharmacology, (Proschak et al. 2018; Bolognesi 2019)), as CBD is known to do (Zlebnik and Cheer
114 2016). Two compounds can also act directly on the same target, either by augmenting or antagonizing
115 each other's effect. CBD appears to ameliorate THC-elicited side-effects (Laprairie et al. 2015; Boggs et
116 al. 2018); it acts as a negative allosteric modulator of the CB1 receptor (Laprairie et al. 2015), whereas
117 THC is a partial agonist (Pertwee 2008). Randomized control trials observed different effects from both
118 compounds consumed alone versus in combination (Solowij et al. 2019). These effects depend both on
119 dose and consumers' past experience, suggesting that future studies looking for possible THC-CBD
120 combinatorial effects must control for these factors, which may be why previous studies have had
121 conflicting results (Boggs et al. 2018). Carefully controlled *in vivo* studies are needed to determine
122 whether distinct ratios of compounds have combinatorial effects. A first step toward defining possible
123 chemical ratios to be used for *vivo* studies is to quantify the ratios present in commercial *Cannabis*. Doing
124 so will also be important for informing the design of human clinical studies aimed at investigating the
125 purported therapeutic effects of cannabis products. Ideally, such studies will test formulations with
126 comparable cannabinoid and terpene ratios to those widely encountered by millions of consumers.

127 Another important reason to quantitatively map the chemotaxonomy of commercial *Cannabis* is
128 that products are commonly labelled with distinct "strain names" or categories with alleged effects,
129 implying that distinct chemical combinations are consistently linked to those labels. For example,
130 consumers believe that *Cannabis* flower labelled "Indica" are reliably sedating, while flower labelled as
131 "Sativa" provide energizing effects (Clarke and Merlin 2013; Lynch et al. 2016; Vergara et al. 2016).
132 Cannabis products are aggressively marketed using these labels. Thus, a better understanding of whether
133 these labels have any reliable association with distinct chemical profiles may have implications for
134 consumer health and safety as well as the regulation of cannabis product marketing.

135       The lack of a standardized, regulated naming system for commercial *Cannabis* varieties has been
136    discussed previously (Sawler et al. 2015; Vergara et al. 2016; Vergara et al. 2020). Various studies, each
137    limited in different ways, have investigated whether these labels capture real chemical variation. For
138    example, cannabinoid and terpene measurements from California samples found limited differences
139    between "Indica" and "Sativa," with some strain names more consistently associated with specific
140    chemical compositions than others (Elzinga et al. 2015). Flower samples from the Netherlands were found
141    to contain specific terpenes more often associated with "Indica" than to "Sativa" samples (Hazekamp et
142    al. 2016). Samples from Washington state limited to total THC and CBD content found no differences
143    between "Indica" and "Sativa," with potency variation between certain strain names (Jikomes and Zoorob
144    2018). Cannabinoid samples across the US did not find a clear relationship between strain name and
145    chemotype, although terpene measurements were not included (Vergara et al. 2020).
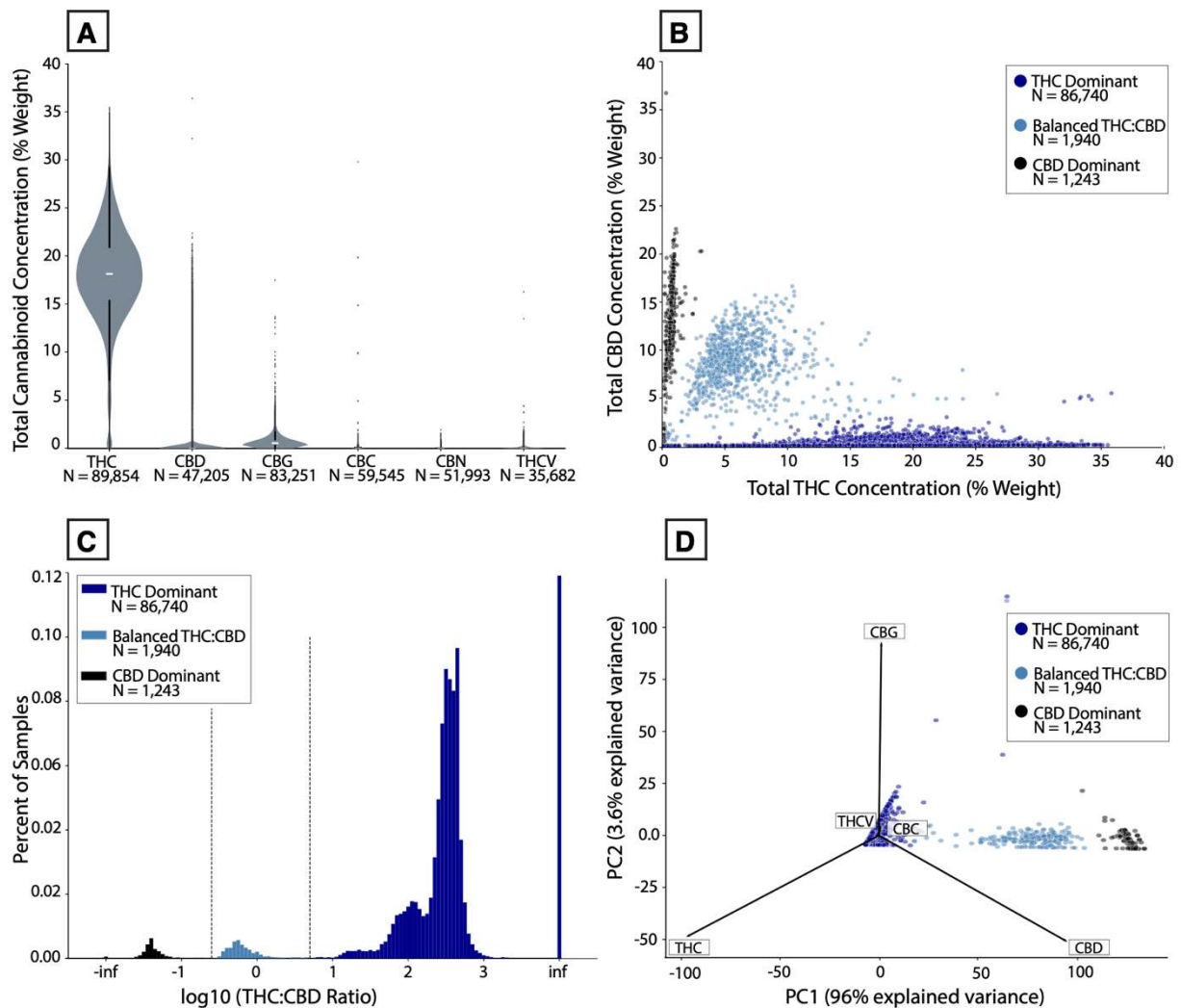
146       In this study, we conducted the largest chemotaxonomic analysis of commercial *Cannabis* flower
147    to date (N = 89,923), using samples from cannabis testing labs in six US states. We analyzed both the
148    cannabinoid and terpene content available for these samples, together with common industry labels and
149    popularity metrics associated with them by the consumer-facing cannabis platform, Leafly. We defined
150    distinct chemotypes that reliably show up across US states and quantified how well the industry labels
151    "Indica," "Hybrid," and "Sativa" map to these chemotypes. We also examined the consistency of "strain
152    names" across samples from different producers. These results provide new possibilities for systematically
153    categorizing commercial *Cannabis* based on chemistry, the design of preclinical and clinical research
154    experiments, and the regulation of consumer marketing in the legal cannabis industry.

155
156    **RESULTS**
157    **Cannabinoid Composition of U.S. Commercial *Cannabis***

158       To assess total cannabinoid levels across samples, we plotted the distribution for each cannabinoid
159    that was consistently measured across regions (Figure 1A) and for every cannabinoid measured within
160    each region (Figure S1). In all regions, total THC levels were much higher compared to levels of all other
161    cannabinoids. Total CBD and CBG were present at modest levels in some samples, while other minor
162    cannabinoids were usually present at very low levels (Figure 1A; Figure S1). Following past work (Hillig
163    2004; Jikomes and Zoorob 2018), we established the presence of three distinct chemotypes based on
164    THC:CBD ratios by plotting total THC against total CBD levels (Figure 1B; see Methods). Most samples
165    belonged to the THC-dominant chemotype (96.5%) in the aggregate dataset (Figure 1B-C) and in each
166    individual region (Figure S2). A much smaller proportion of samples were classified as CBD-dominant
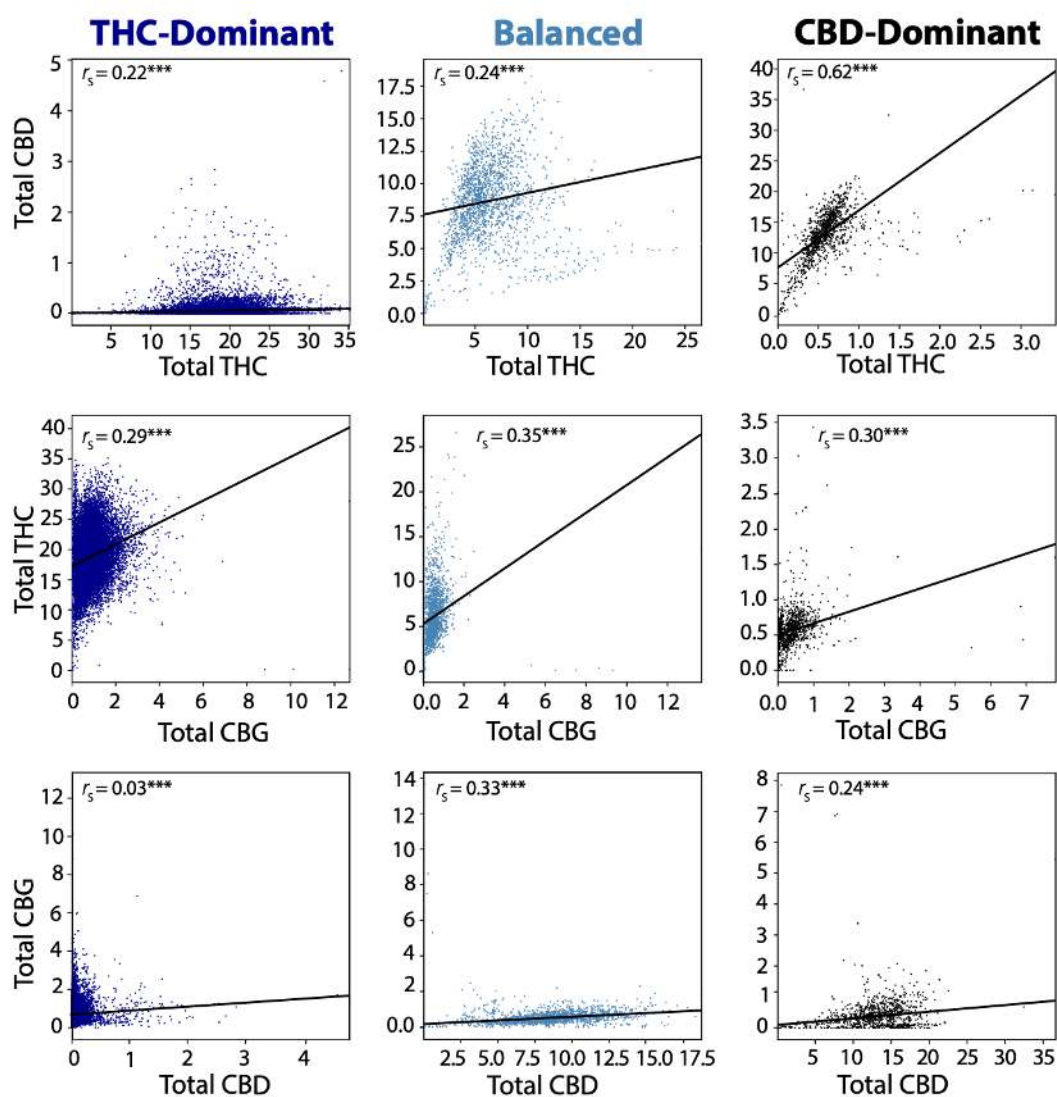167    (1.4%) or Balanced THC:CBD (2.2%; Figure 1; Figure S1).

**Figure 1: Cannabinoid variation among commercial *Cannabis* samples in the US. (A)** Violin plot of distribution of the set of common cannabinoids measured across all regions **(B)** Total THC vs. Total CBD levels, color-coded by THC:CBD chemotype. **(C)** Histogram showing THC:CBD distribution on a $\log_{10}$ scale. "Inf" stands for "infinite" (any samples with 0 total THC or CBD). **(D)** Principal Component Analysis of all cannabinoids shown in panel A, color-coded by THC:CBD chemotype.

Although most samples contained low levels of cannabinoids beyond THC, we observed that 3.9% and 23.1% of samples, respectively, had total CBD or total CBG of 1% by weight or higher. To further understand any systematic patterns of variation in cannabinoid profiles beyond THC and CBD levels, we performed Principal Component Analysis (PCA) on all samples that contained measurements for total THC, CBD, CBG, CBC, CBN, and THCV content. Most of the variance in this dataset (96%) was explained by the first principal component (Figure 1D), which was highly correlated with samples' THC:CBD ratios ($r_s = -0.51$, $P < 0.0001$). Most of the remaining variation (3.6%) was explained by the second principal component, which was highly correlated with total CBG levels ($r_s = 0.95$, $P < 0.0001$).

6

183    Thus, the vast majority of variance in cannabinoid profiles is explained by variation among the three most

184    abundant cannabinoids (THC, CBD, CBG) in commercial *Cannabis* in the US.

185          To further understand the relationship between levels of each pair of these three cannabinoids, we

186    plotted total levels of THC, CBD, and CBG against each other, separately for each THC:CBD chemotype.

187    Given that CBGA is the precursor molecule to both THCA and CBDA, we expected to see positive

188    correlations between each cannabinoid pair. This is what we observed, with the strength of each

189    correlation varying across THC:CBD chemotypes (Figure 2). One notable finding with potential

190    regulatory consequences is the substantial correlation between total THC and CBD levels in CBD-

191    dominant samples ($r_s = 0.65$, $P < 0.0001$). 84.5% of CBD-dominant samples had total THC levels above

192    0.3%, the threshold used to legally define hemp in the US. This indicates that a substantial fraction of

193    CBD-dominant *Cannabis* would not meet the legal definition of hemp in the US.



194
195  **Figure 2: Correlations among total THC, CBD, and CBG levels in each THC:CBD chemotype.** Scatterplots
196  showing the linear correlation between total THC, CBD, and CBG levels in each of the main THC:CBD

197  chemotypes. Top Row: Total THC vs. Total CBD; middle row: Total CBD vs. Total THC. Bottom row: Total CBD
198  vs. Total CBG. ***$P < 0.0001$

199

## Terpene Composition of U.S. Commercial *Cannabis*

201          We next assessed which terpene compounds were most prominent in samples by plotting the
202  distribution of each terpene that was consistently measured in each region. On average, the terpenes
203  myrcene, β-caryophyllene, and limonene were present at the highest levels (Figure 3A). In most cases,
204  individual terpenes were rarely present at more than 0.5% weight and most were present at low levels (<
205  0.2%) in the majority of samples. Overall, total terpene content averaged 2% by weight and displayed a
206  modest but robust positive correlation with total cannabinoid content ($r_s = 0.37$, $P < 0.0001$), suggesting
207  that the production of one type of compound doesn't come at the expense of the other.

208          To validate that patterns expected from previous studies were observed in the terpene data, we first
209  looked for correlations between specific terpene pairs. We chose pairs that have been previously observed
210  to display robust positive correlations, likely stemming from constraints on their biochemical synthesis
211  (Booth et al. 2017; Allen et al. 2019; Booth and Bohlmann 2019). Strong positive correlations were seen
212  between α- and β-pinene (Figure 3B; $r_s = 0.78$, $P < 0.0001$), as well as β-caryophyllene and humulene
213  (Figure 3C; $r_s = 0.88$, $P < 0.0001$). These correlations held for both the aggregate dataset (Figure 3) and
214  for each individual US state (Figures S3 and S4), demonstrating their robustness across regions.
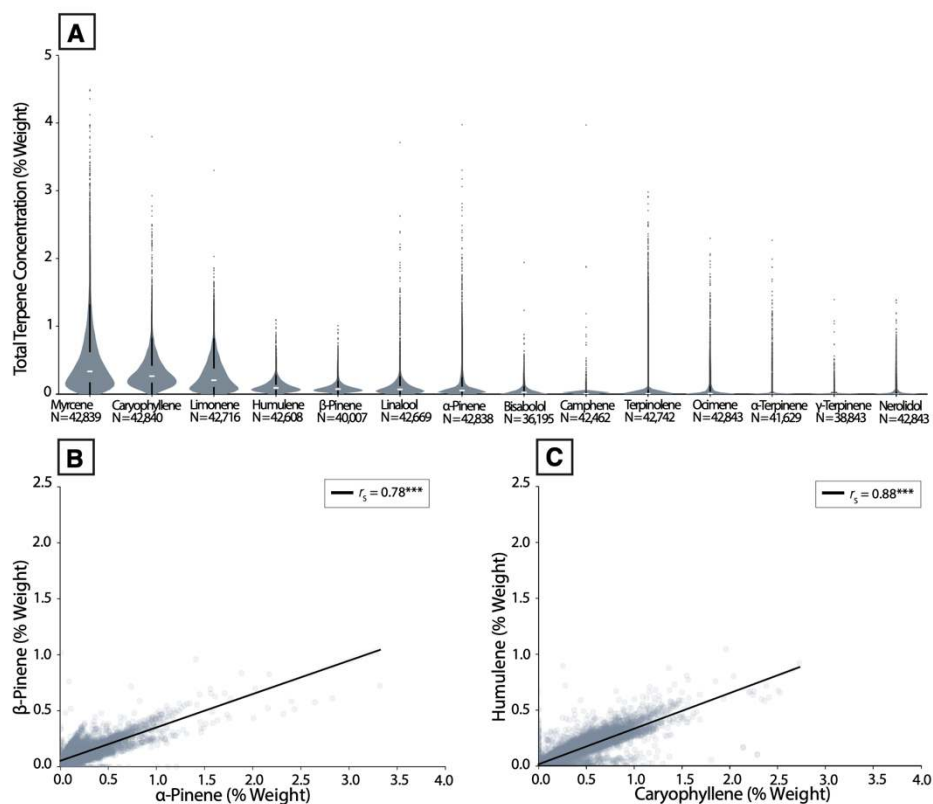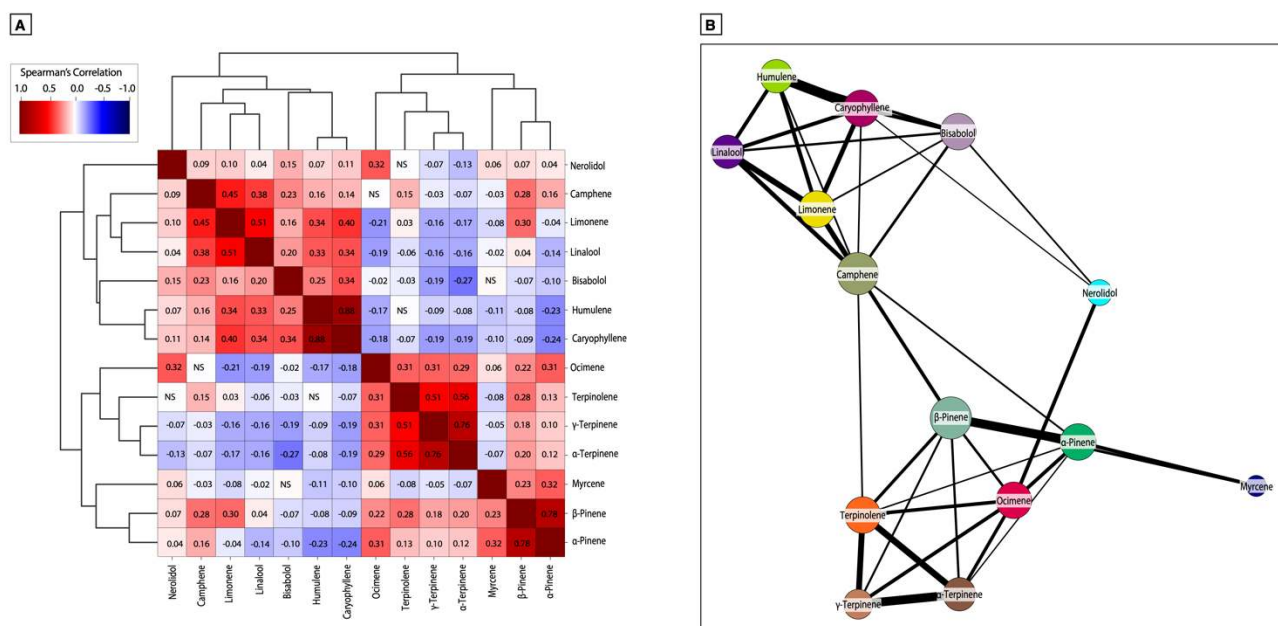


**Figure 3: Terpene abundance across commercial *Cannabis* samples in the US. (A)** Violin plots showing distributions of the set of common terpenes measured across all regions **(B)** Scatterplot showing the correlation between α- and β-pinene, two common pinene isomers. $r_s = 0.78$, ***$P < 0.0001$ **(C)** Scatterplot showing the correlation between β-caryophyllene and humulene, two *Cannabis* terpenes co-produced by common enzymes. $r_s = 0.88$, ***$P < 0.0001$

237       In order to systematically understand relationships between all terpene pairs, we performed

238    hierarchical clustering on all pairwise correlations among terpenes (Figure 4A; see Methods). This

239    revealed distinct clusters of co-occurring terpenes. After controlling for multiple comparisons, we

240    observed many robust correlations between terpenes (see Methods). We also plotted this data in the form

241    of a network diagram configured to display connections between terpenes with the strongest correlations

242    (Figure 4B). This diagram provides a more compact picture of terpene co-occurrence and likely reflects

243    the underlying biosynthesis pathways that give rise to these correlations (Booth et al. 2017; Allen et al.

244    2019; Booth and Bohlmann 2019).



245
246 **Figure 4: Patterns of terpene co-occurrence among commercial *Cannabis* samples in the US. (A)**
247 Hierarchically clustered correlation matrix showing pairwise correlations between all terpenes consistently
248 measured across regions. **(B)** Network diagram where nodes are terpenes and edges are thresholded to the strongest
249 observed correlations and their widths correspond to the strength of the correlation. [explanation of circle sizes and
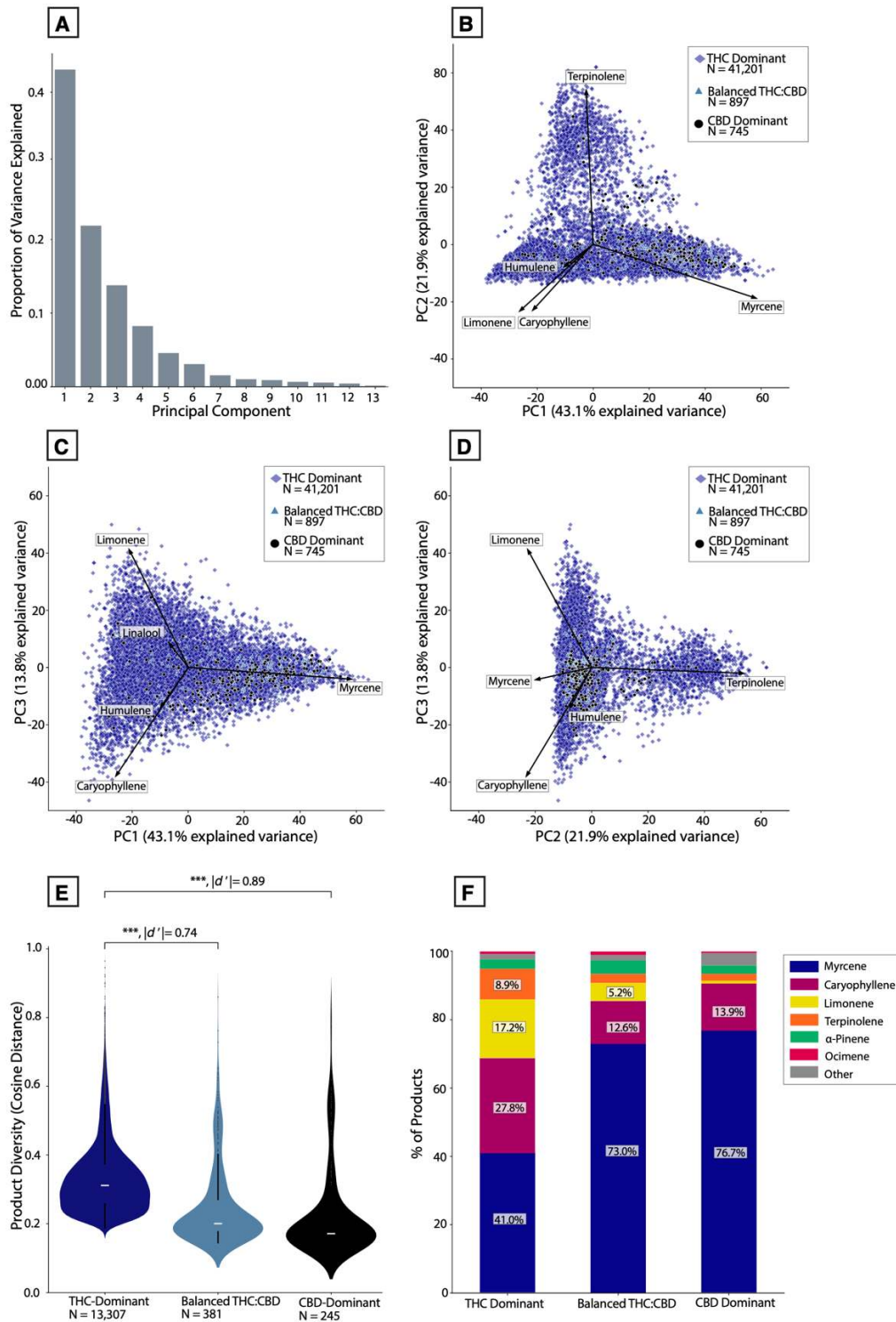250 line widths]

251

252 **THC-Dominant And High-CBD *Cannabis* Display Distinct Levels of Terpene Diversity**

253       Historically, the major focus of both clandestine and legal *Cannabis* breeding in the US has been

254    on THC-dominant varieties, which is why they predominate in the commercial marketplace (Figure 1)

255    (Clarke and Merlin 2016). It is therefore expected that THC-dominant cultivars will display a more diverse

256    array of terpene profiles than CBD-dominant and balanced THC:CBD cultivars. To visualize patterns of

257    variation among terpene profiles, we performed a Principal Component Analysis (PCA) on the terpene

258    data (see Methods). The first three principal components explained 78.7% of the variance in the data

9

259    (Figure 5A), indicating that most of the variance in terpene profiles can be explained with just a few
260    components.

261         To visualize how patterns of terpene profile variation map to the major THC:CBD chemotypes
262    shown in Figure 1, we plotted PCA scores for all samples along the first three principal components, with
263    each sample color-coded by its THC:CBD chemotype (Figure 5 B-D). The superimposed vectors encoding
264    the five terpenes with the strongest loadings onto each principal component help clarify the terpene
265    composition of different points on the graph. Most CBD-dominant and balanced THC:CBD samples
266    cluster within a smaller subsection of the plots compared to THC-dominant samples. To quantify terpene
267    profile variation across each THC:CBD chemotype, we computed the mean pairwise cosine distance in
268    terpene profiles within each THC:CBD chemotype and used this as a measure of diversity. We conducted
269    this analysis at the product level rather than sample level, as individual samples of the same product tend
270    to be highly similar (see Methods). THC-dominant products displayed significantly higher levels of
271    diversity than both balanced THC:CBD (Figure 5E; $P < 0.0001$, $|d'| = 0.74$) and CBD-dominant products
272    (Figure 5E; $P < 0.0001$, $|d'| = 0.89$). In particular, a higher proportion of CBD-dominant and balanced
273    THC:CBD products displayed myrcene-dominant terpene profiles compared to THC-dominant samples
274    (Figure 5F).

**Figure 5: Patterns of terpene profile diversity across THC:CBD chemotypes. (A)** Histogram showing the proportion of variation explained by each principal component after performing Principal Component Analysis on the terpene dataset. **(B)** PCA scores plotted along PC1 and PC2, color-coded by major THC:CBD chemotype. Vectors depict the loadings of the five individual terpenes onto these principal axes. **(C)** PCA scores plotted along PC1 and PC3. **(D)** PCA scores plotted along PC2 and PC3. **(E)** Violin plot showing distribution of 'product diversity' values (cosine distances) for each THC:CBD chemotype. Product values are calculated by averaging samples with the same strain name linked to a given producer ID. ***$P < 0.0001$, Welch's t-test and Cohen's d'. **(F)** Stacked bar chart showing the percent products with a given dominant terpene for each THC:CBD chemotype.
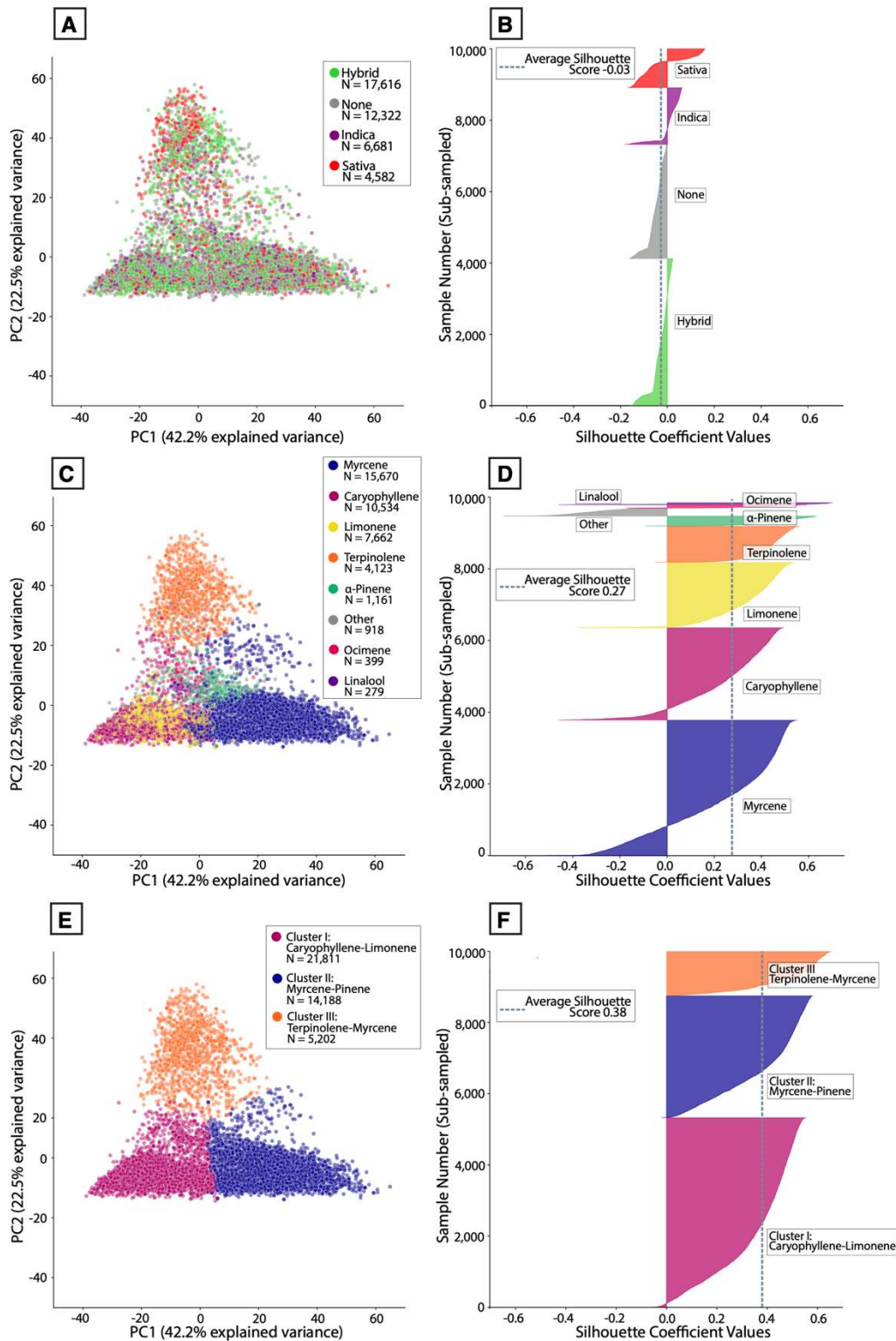
11

**284**  **Cluster Analysis Reveals Distinct Terpene Chemotypes And Poor Validity of Common Commercial**
**285**  **Labels**

286      Given the observed diversity of terpene profiles displayed by THC-dominant samples, we wanted
287  to establish how this diversity is captured by the categorization system most commonly used for
288  commercial THC-dominant *Cannabis*. Commercial products are routinely labelled "Indica," "Hybrid," or
289  "Sativa." Prevailing folk theories assert that "Indica" products provide sedating effects, "Sativa"
290  energizing effects, and "Hybrids" intermediate effects (McPartland and Small 2020). If this were true, we
291  would expect to see a reliable difference between the chemical composition of samples attached to each
292  label. To test this, we devised an approach using silhouette analysis to quantify how well these industry
293  labels capture the observed chemical diversity (see Methods). We compared this commercial labelling
294  system to labelling the data with simplified chemical designations (each samples' dominant terpene), as
295  well as an unbiased approach using k-means clustering.

296      Figure 6A displays THC-dominant samples plotted along the first two principal components,
297  color-coded by their Indica/Hybrid/Sativa label. The samples are highly intermingled, with no obvious
298  segregation of data points by commercial label. This is reflected in the corresponding silhouette plot,
299  which displays a low mean silhouette score (Figure 6B). The majority of samples have a negative score,
300  indicating that many samples with one label could be easily confused with samples of a different label in
301  terms of terpene profile. In other words, it is likely that a sample with the label 'Indica' will have an
302  indistinguishable terpene composition as samples labelled "Sativa" or "Hybrid." By comparison, when
303  samples are labelled by their dominant terpene, there is better visual separation of data points by their
304  label (Figure 6C) and a higher mean silhouette score (Figure 6D). These results indicate that even a
305  simplistic labeling system, in which THC-dominant samples are labelled by their dominant terpene, is
306  better at discriminating samples than the industry-standard labelling system.

307      To segment samples in an unbiased fashion based on terpene profile, we applied the k-means
308  clustering algorithm to define clusters of samples in the data. This approach allowed us to cluster the data
309  using a standard method for determining a number of clusters that fits this dataset well (Figure 6E; Figure
310  S6-8; see Methods). Three major clusters were defined. As expected, this algorithmic partitioning of the
311  data is better at assigning points to distinct groups, especially compared to the Indica/Sativa labels. This
312  is reflected in the higher mean silhouette score and low proportion of samples with negative silhouette
313  values (Figure 6F). This data can be clustered in different ways, such as defining additional sub-clusters
314  within the clusters displayed here (Figure S5). Ideally, this type of analysis would be further constrained
315  by other data sources, such as sample genotypes and other classes of metabolites. For the purposes of this

12

316 study, we focused on the three large clusters depicted in Figure 6 and conducted further analysis on their

317 relationship to common commercial categories.

318



319 **Figure 6: Commercial 'strain category' labels poorly align to patterns of phytochemistry. (A)** PCA scores for
320 all THC-dominant samples plotted along PC1 and PC2, color-coded by Indica/Hybrid/Sativa label attached to each
321 sample. **(B)** Silhouette coefficients for each sample with a given Indica/Hybrid/Sativa label. **(C)** PCA scores for all
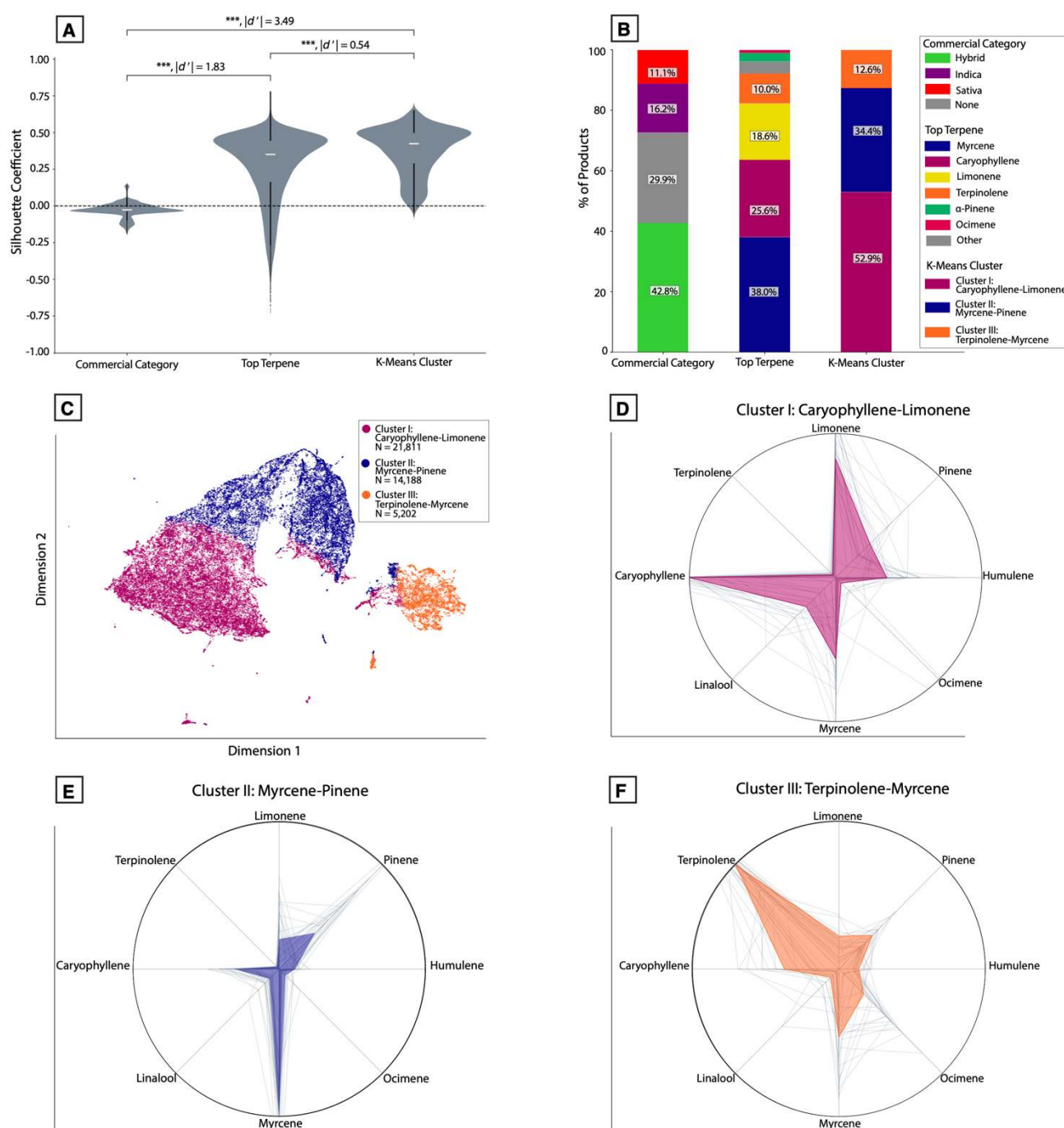
322    THC-dominant samples plotted along PC1 and PC2, color-coded by the dominant terpene of each sample. **(D)**
323    Silhouette coefficients for each sample with a given dominant terpene. **(E)** PCA scores for all THC-dominant
324    samples plotted along PC1 and PC2, color-coded by k-means cluster labels attached to each sample. **(F)** Silhouette
325    coefficients for each sample with a given k-means cluster label. Each silhouette plot depicts a random subset of
326    10,000 samples from the full dataset (n=41,201).

327

328        The distribution of silhouette scores across each of the three labelling systems allows us to compare

329    the results depicted in Figure 6. Labelling data either by dominant terpene or by k-means cluster was

330    significantly better at capturing the terpene diversity seen in THC-dominant samples compared to the

331    commercial labels (Figure 7A; $P < 0.0001$, $|d'| = 3.49$, k-means vs. commercial labels). Regardless of the

332    labelling system, samples are not evenly distributed among groups (Figure 7B). To further visualize the

333    clusters defined in Figure 6E-F, we used Uniform Manifold Approximation and Projection (UMAP) to

334    visualize the data (Figure 7B). UMAP is a dimensionality reduction technique like PCA but without

335    linearity assumptions. The dimensions returned by UMAP lack the interpretability (e.g. factor loadings)

336    associated with PCA but are superior at recovering latent clustered structure within high-dimensional data

337    (Dorrity et al. 2020). More of the individual data points are visible in this plot compared to the PCA plots

338    shown in Figure 6.

339        Averaging the full cannabinoid and terpene profile of all products within each cluster allowed us

340    to depict the average chemical composition of each cluster. We plotted mean terpene profiles as

341    normalized polar plots together with the total THC, CBD, and CBG distributions of each cluster (Figure

342    7C-F). In relative terms, a simplified description for the terpene profiles characterizing each cluster is:

343    "high caryophyllene-limonene" (Cluster I), "high myrcene-pinene" (Cluster II), and "high terpinolene-

344    myrcene" (Cluster III; Figure 4 B-D). Similar groups are seen across regional datasets (Figure S6). We

345    also observed that one cluster (Cluster III: "high terpinolene-myrcene") had somewhat higher total CBG

346    levels compared to the other clusters (median CBG 0.98% vs 0.65%; $P < 0.0001$, $|d'| = 0.57$). This

347    appeared to be due to a modest but significant correlation between total CBG and terpinolene levels ($r_s =$

348    0.17, $P < 0.0001$).

349

350

**Figure 7: Cluster analysis reveals distinct chemotypes of THC-dominant commercial *Cannabis* commonly present in US states. (A)** Violin plot showing the distribution of silhouette coefficients for each labelling method. ***$P < 0.0001$, Welch's t-test and Cohen's d'. Absolute effect sizes are given as Cohen's d' values. ***p<0.0001, **p<0.001; *p<0.01 **(B)** Stacked bar chart showing the percent of samples falling within each group for each labelling system. **(C)** UMAP embedding in two dimensions showing samples classified into each k-means cluster. **(D)** Polar plot showing the mean, normalized levels of eight of the terpenes most commonly observed for Cluster I (high caryophyllene-limonene) products. **(E)** Similar polar plot for Cluster II (high myrcene-pinene) products. **(F)** Similarly polar plot for Cluster III (high terpinolene-myrcene) products. Gray lines represent the top 25 products from each cluster with the most samples per product.

360

**Commercial "Strain Names" Display Differential Levels of Chemical Consistency**

The cannabis industry also uses colloquial "strain names" to label and market products. Distinct "strains" of THC-dominant *Cannabis* are purported to offer distinct psychoactive effects, such as "sleepy," "energizing," or "creative." While the commercial use of nomenclature is not accepted by the scientific community, it is conceivable that distinct chemovars of THC-dominant *Cannabis* could cause different psychoactive effects, on average. In principle, if commercial "strain names" are indicative of different psychoactive effects in a discernible way, then different strain names should reliably map to distinct chemotypes. Alternatively, because there are few regulatory constraints on the nomenclature of commercial *Cannabis*, it is possible that *Cannabis* producers attach strain names to their products in arbitrary or inconsistent ways. If this were true, we would not expect to see strain names consistently map to specific chemotypes above chance levels.

To quantify chemical consistency among THC-dominant products, we compared each product's chemical composition in terms of the 14 major terpenes depicted in Figures 3-4. We did this for all strain names where the underlying data was attached to at least five product IDs each having five or more samples with that particular strain name. To validate whether the strain names attached to more testing data are representative of those encountered by consumers, we plotted the number of products attached to each strain name vs. consumer popularity, measured in terms of unique online pageviews to the consumer *Cannabis* database, Leafly.com. We observed a strong positive correlation ($r_s$ = 0.59, $P$ < 0.0001), indicating that the strain names in our analysis are representative of the names encountered by consumers in commercial settings.

As a measure of consistency, we computed the pairwise cosine similarity of all products attached to each strain name and visualized this in a similarity matrix (Figure 8B, ten most abundant strain names shown). Next, we quantified the average pairwise similarity of all products sharing a common strain name. For each strain name, we plotted the distribution of product similarity scores, sorted from highest to lowest mean similarity, for the 41 strain names used in this analysis (Figure 8C). We compared these values to the average similarity score computed after randomly shuffling strain names across all product IDs (Figure 8C, dashed line). This allowed us to model the situation where each producer has arbitrarily labelled their product with a given strain name. The mean between-product similarity was significantly higher compared to the shuffled dataset for the majority strain names (Figure 8C, $P$ < 0.0001, $|d'|$ = 1.44). For some strain names, product similarity did not significantly differ from the shuffled distribution or was even below this, and there was a large amount of variability in mean consistency scores across all strain names. To illustrate this variability further, we overlaid the individual profiles of all products with a given name,

16

393     separately for two strain names: one with a relatively high level of between-product similarity ("Purple

394     Punch") and one with a low level ("Tangie"; Figure 8D).

395        To assess between-product similarity in terms of the major clusters defined previously, we applied

396     the same clustering approach from Figures 6-7 to the product averages analyzed in Figure 8. These data

397     were visualized in a UMAP embedding, with all products attached to the two example strain names (Figure

398     8D highlighted Figure 8E). This illustrates how a relatively consistent (Purple Punch) vs. inconsistent

399     (Tangie) strain name maps to this space. 96% of product averages attached to Purple Punch fall within

400     Cluster I (high caryophyllene-limonene), while only 62.5% of product averages for Tangie fall into a
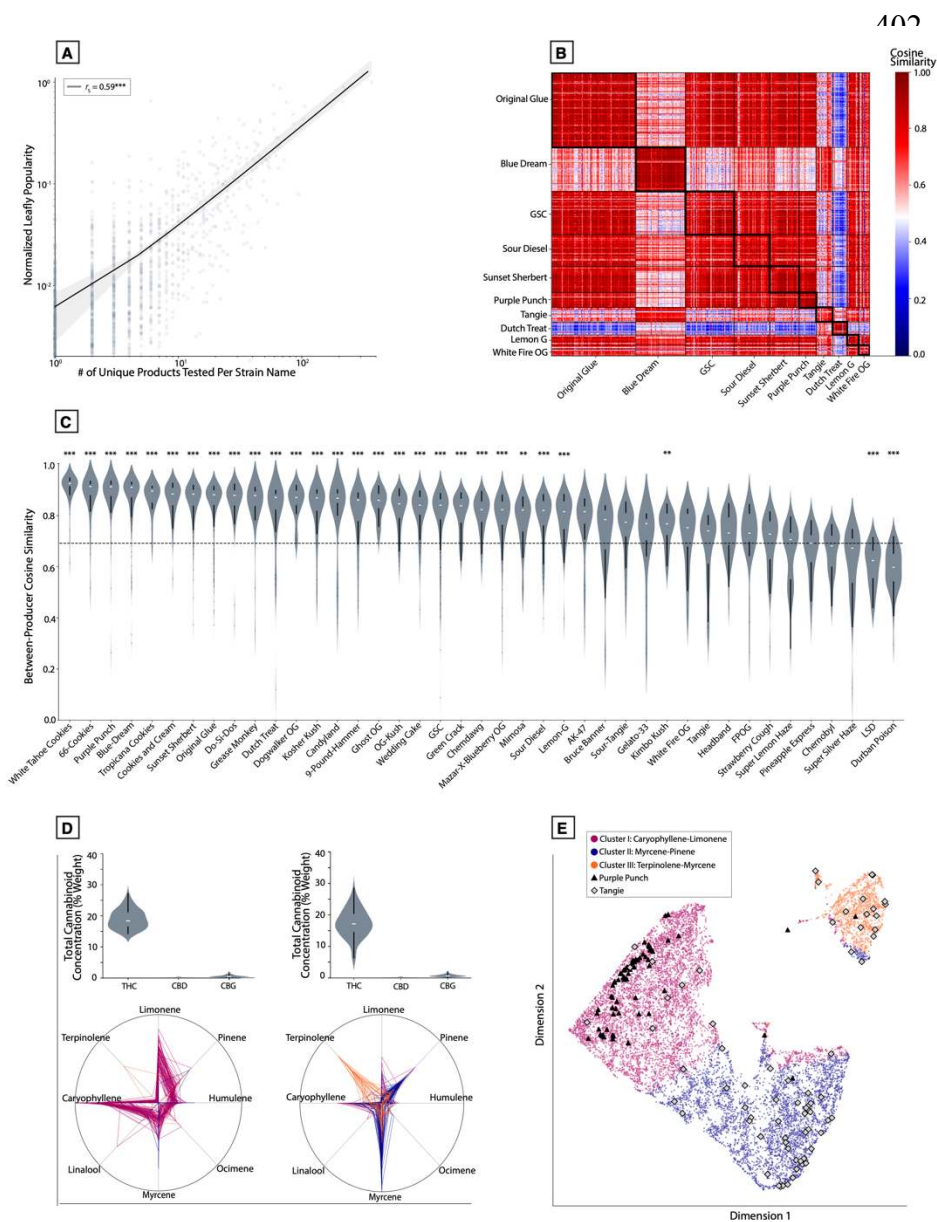
401     single cluster.



**Figure 8: Strain names are associated with variable levels of chemical consistency across *Cannabis* products. (A)** Scatterplot of the number of products tested vs. normalized Leafly popularity for all product-level data attached to strain names ($\log_{10}$ scale). $r_s = 0.59$, ***$P < 0.0001$ **(B)** Similarity matrix depicting pairwise cosine similarities between all product-level data attached to the ten most common strain names by abundance. **(C)** Violin plot depicting the distribution of cosine similarity scores between products attached to the same strain name. Dashed line represents the average similarity level after randomly shuffling strain names. **$P < 0.001$, ***$P < 0.0001$, Welch's t-test.***p<0.0001; **p<0.00024; *p<0.0012 Welch's t-test. **(D)** Violin plots representing total cannabinoid distributions and polar plots representing terpene profiles for all products attached to the strain names "Purple Punch" (left) and "Tangie" (right); **(E)** UMAP embedding showing where each of the product

439     samples for Purple Punch and Tangie from panel D show up in this representation.
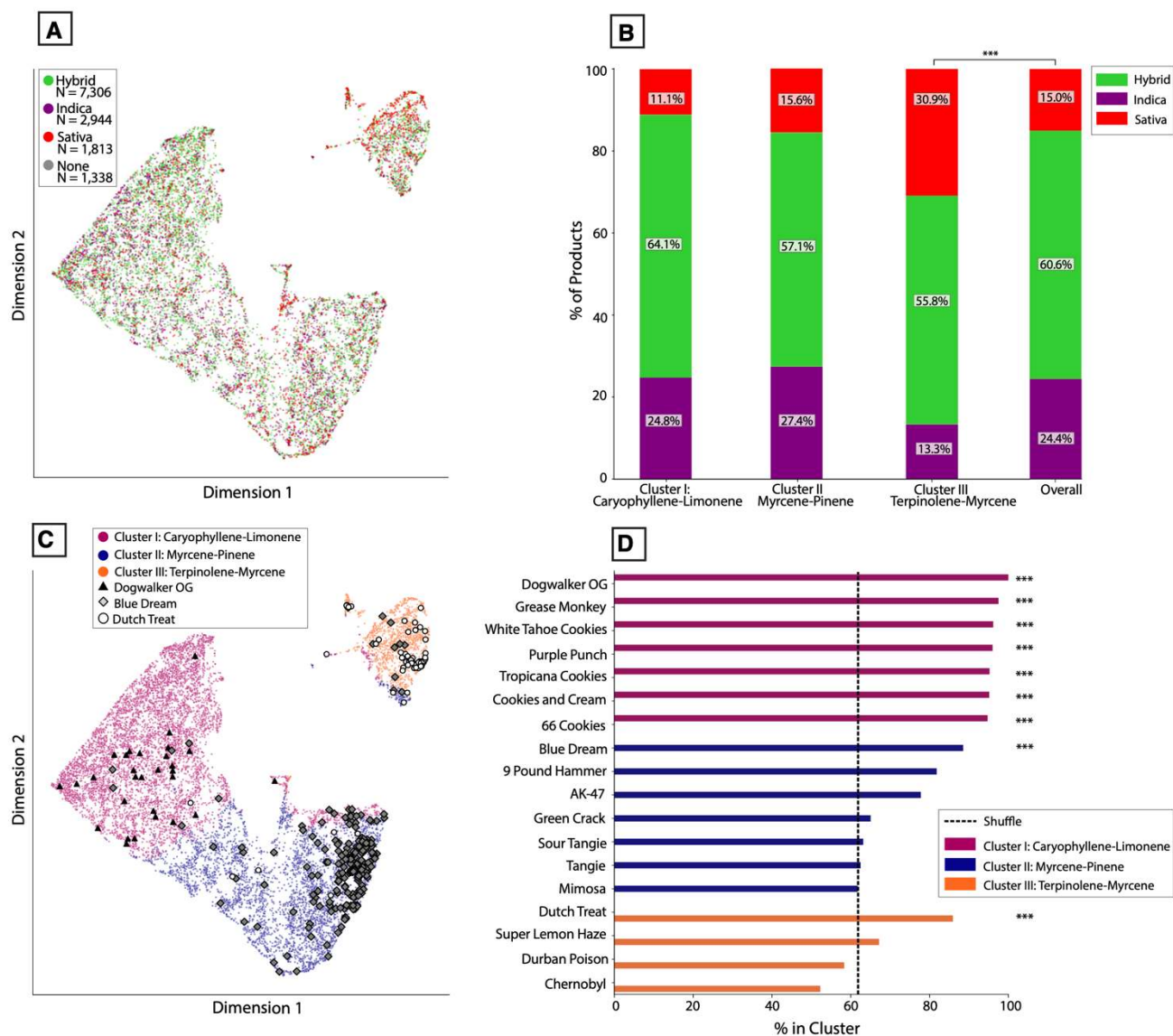
**Some Commercial Labels Are Over-Represented in Specific Chemically Defined Clusters**

440

441    To further understand whether any strain names were overrepresented in our algorithmically

442    defined clusters, as appeared true for Purple Punch (Figure 8E), we calculated the proportion of all

443    products with a given strain name that belonged to each cluster. For each strain name displayed in Figure

444    8C, we calculated that proportion for whichever cluster contained the highest count of products with that

445    name. For example, 96% of products attached to the name "Purple Punch" were found in Cluster I, much

446    higher than the 61.8% expected if product strain names are randomly shuffled ($P < 0.0001$, $|d'| = 2.47$).

447    We plotted this proportion for the 18 most overrepresented strain names, grouped by their primary cluster

448    and compared these to the average cluster frequency expected from shuffling strain names across products

449    (Figure 9A). For each cluster, there are strain names that are highly overrepresented. 100% of "Dogwalker

450    OG" products are found within Cluster I ("high caryophyllene-limonene"; $P < 0.0001$, $|d'| = 1110.4$),

451    88.5% of "Blue Dream" products are found within Cluster II ("high myrcene-pinene"; $P < 0.0001$, $|d'| =$

452    1.2), and 85.9% of "Dutch Treat" products are found within Cluster III ("high terpinolene"; $P < 0.0001$,

453    $|d'| = 1.0$).

454    Similar to Figure 8E, we plotted the single most over-represented strain name associated with each

455    cluster in a UMAP embedding of all the product-level data (Figure 9B). These strain names represent

456    those that are the most consistently associated with a given chemotype. Notably, even these strain names

457    are not perfectly associated with a single chemotype, and products attached to each name display

458    variability within each cluster. This indicates that even the strain names with the highest levels of

459    consistency across products still display a non-trivial amount of variation. An interactive 3-D version of

460    this product-level UMAP (including high-CBD products) is also included (see Methods).

461    In doing this analysis, we noticed that one cluster (Cluster III, characterized by high terpinolene

462    levels) contained a paucity of products attached to strain names labelled as "Indica." To understand

463    whether any of the Indica/Hybrid/Sativa industry labels were over- or under-represented within any of

464    these clusters, we performed a similar analysis for commercial categories as we did for strain names: for

465    each of the three clusters, we calculated the proportion of products attached to Indica/Hybrid/Sativa labels.

466    For each of these, we compared it to the population frequency of each category. For Cluster I and Cluster

467    II, the frequency of products attached to Indica/Hybrid/Sativa labels did not significantly differ from those

468    observed in the full set of products with Indica/Hybrid/Sativa labels. In contrast, Cluster III (high

469    terpinolene) did show a significant difference, with approximately twice as many Sativa-labelled products

470    and half as many Indica-labelled products as expected from the full population (Figure 9B; $X^2 = 22.2$, $P$

471    $< 0.0001$, Chi-squared test). This over-representation of Sativa-labelled products can also be seen in the

472  UMAP embedding (Figure 9D), which displays product-level data color-coded by Indica/Hybrid/Sativa

473  label.



474
475  **Figure 9: Some commercial *Cannabis* labels are overrepresented for specific chemotypes. (A)** UMAP
476  embedding of product-level data as in Figure 8E, color-coded by Indica/Hybrid/Sativa label. **(B)** Stacked bar chart
477  showing the proportion of products labelled as Indica, Hybrid, or Sativa within each k-means cluster, compared to
478  the overall distribution. ***$P < 0.0001$, Chi-squared test. **(C)** UMAP embedding of product-level data as in Figure
479  8D, color-coded by k-means cluster label, showing where all products attached to either "Blue Dream" or "Dutch
480  Treat" are found. **(D)** Bar charts showing the percent of products attached to each strain name that are found in a
481  given k-means cluster, color-coded by its most prominent cluster. Dashed line represents expected percent after
482  randomly shuffling strain names. ***$P < 0.0001$, Welch's t-test.

483

484  **Discussion**

485           To our knowledge, this study represents the largest quantitative chemical mapping of commercial

486  *Cannabis* to date. It builds on a literature examining the chemotaxonomy of *Cannabis* samples taken from

487  individual regions of the US (Elzinga et al. 2015; Henry et al. 2018; Vergara et al. 2020), Canada (Mudge

488  et al. 2019), and Europe (Hazekamp and Fischedick 2012; Hazekamp et al. 2016), as well as classic studies

19

489    of the chemotaxonomy of non-commercial *Cannabis* (Hillig 2004; Hillig and Mahlberg 2004). We

490    mapped and analyzed the cannabinoid and terpene diversity of almost 90,000 samples from six US states

491    and found distinct chemotypes of *Cannabis* that are reliably present across regions.

492         Because *Cannabis* remains federally illegal in the US, the laboratory-derived data from each state

493    represent distinct pools of *Cannabis* found within those states. Even with clones, environmental factors

494    such as variation in growing conditions and preparation procedures can cause differences in morphology

495    and chemotype expressions that are measured by testing labs (Magagnini et al. 2018). Moreover, the

496    measurements themselves are made by different labs, using methodologies that may not be standardized

497    (See Methods, Data Collection). Nonetheless, we observed similar patterns across regions. In all states,

498    the sample population is comprised mostly of THC-dominant samples, each with a similar distribution of

499    major terpenes (Figures S2, S6) and displaying the terpene-terpene correlations expected based on the

500    constraints of terpene biosynthesis (Booth et al. 2017; Booth and Bohlmann 2019; Booth et al. 2020), as

501    has been observed elsewhere (Allen et al. 2019; Mudge et al. 2019). The pooled dataset also displays

502    features seen in sample populations from US states not represented here (Henry et al. 2018). Collectively,

503    these results suggest that, while some regional variation may exist, the major patterns of cannabinoid and

504    terpenes profiles are similar throughout the US.

505         We used cluster analysis to define at least three major chemotypes of THC-dominant *Cannabis*

506    prevalent in the US (Figures 6-7; Figure S5). In simplified terms, samples from each cluster tend to be

507    characterized by relatively high levels of β-caryophyllene and limonene (Cluster I), myrcene and pinene

508    (Cluster II), or terpinolene and myrcene (Cluster III). Samples across these clusters display similar total

509    THC distributions, while Cluster III is associated with modestly higher CBG levels (Figure 7). The

510    chemotype landscape of commercial *Cannabis* is highly uneven, with less than 96.5% of samples

511    classified as THC-dominant, and 87.4% of these samples belonging to either the Cluster I (high

512    caryophyllene-limonene) or Cluster II (high myrcene-pinene). Breeding new *Cannabis* chemotypes not

513    represented in the current commercial landscape will be a key area of future innovation.

514         We observed that the diversity of cannabinoid profiles displayed by commercial *Cannabis* in the

515    US is explained almost entirely by variation in total THC, CBD, and CBG content, with the majority of

516    variation coming from THC content (Figure 1). Similar to classic work on non-commercial *Cannabis*

517    *(Hillig and Mahlberg 2004)*, our results show distinct THC:CBD chemotypes: THC-dominant, balanced

518    THC:CBD, and CBD-dominant. These likely arise from distinct genotypes. The genes giving rise to the

519    cannabinoid synthases responsible for producing the major cannabinoid acids are highly similar (Vergara

520    et al. 2019; van Velzen and Schranz 2020; Vergara et al. 2021b). Copy number variation (Vergara et al.

521    2019; Vergara et al. 2021b) or allelic variation (Onofri et al. 2015) in the genes encoding these enzymes
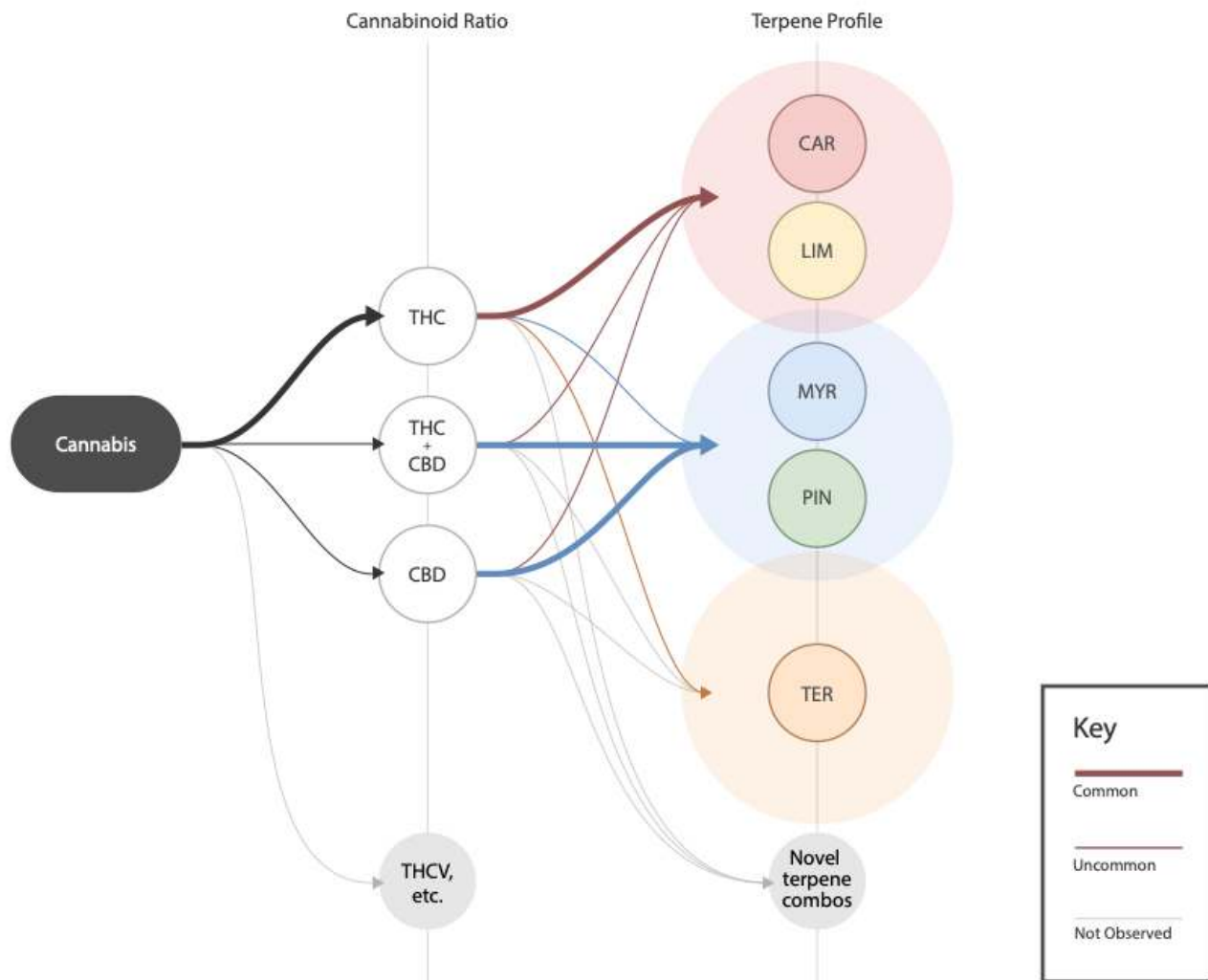
522  may explain the observed variation in cannabinoid ratios. Interesting areas of future study will be to

523  correlate chemotype and genotype directly and determine why other minor cannabinoids have such low

524  abundance in commercial *Cannabis*. For example, there are numerous CBC-related genes (van Velzen

525  and Schranz 2020) but we observe very low levels of CBC (Figures 1-2), supporting previous claims that

526  CBCA synthase may not be selective for CBC production (Vergara et al. 2020).

527        The observed variation in terpene profiles is also likely related to underlying genotypic variation.

528  While environmental and developmental modulation of terpene profiles is possible (Aizpurua-Olaizola et

529  al. 2016), the fact that we observe a similar set of major profiles across US states (Figure S6) suggests that

530  these profiles have a strong genetic component. *Cannabis* terpenes are synthesized from enzymes encoded

531  by multiple genes (Booth et al. 2017; Allen et al. 2019; Booth and Bohlmann 2019; Booth et al. 2020).

532  The robust correlation patterns we observed among many of the most abundant *Cannabis* terpenes likely

533  arise from variation in biosynthetic enzymes. The underlying genetic networks regulating these

534  biochemical pathways are complex (Booth et al. 2017; Allen et al. 2019; Booth and Bohlmann 2019;

535  Booth et al. 2020) and more research may be needed to inform efficient breeding programs to generate

536  novel chemotypes.

537        Despite the chemotypic diversity we observed for THC-dominant *Cannabis*, this likely represents

538  a fraction of the diversity the plant is capable of expressing. For example, although one of the clusters we

539  defined is characterized by especially high myrcene levels, each of the three clusters contain samples

540  where myrcene is more abundant than most other terpenes. This pattern is stronger for CBD-dominant

541  and balanced THC:CBD chemotypes, where the majority of samples are myrcene-dominant. This may

542  reflect a historical genetic bottleneck, whereby most *Cannabis* grown in the US is descended from a subset

543  of the worldwide lineages (McPartland and Small 2020). The relative lack of diversity among high-CBD

544  cultivars is likely due to the historical focus on breeding high potency THC-dominant *Cannabis* in the

545  US. In principle, there is no biological limitation preventing the breeding of high-CBD cultivars with

546  similar terpene diversity to what is seen in THC-dominant cultivars. Many of the genes encoding the

547  synthetic enzymes for terpene production are located on different chromosomes from those involved in

548  cannabinoid acid synthesis (Booth et al. 2020) or are found far apart from each other in the same genomic

549  region (Allen et al. 2019), and therefore could be assorted through recombination. These two aspects of

550  chemical phenotype may therefore be independently inherited, similar to other phenotypic traits (Vergara

551  et al. 2021a).

552        While not observed in this commercial dataset, chemovars that predominate in other cannabinoids,

553  such as CBG, have been bred and may offer distinct psychoactive or medicinal effects compared with the

554  high-THC chemovars that predominate commercially (Hutchison et al. 2019). There were few samples

21

555 that contained an abundance of minor cannabinoids, suggesting that commercial *Cannabis* in the US is

556 much more homogenous than it could be. An exciting area for academic research and product innovation

557 lies in the breeding of new varieties with higher levels of other cannabinoids. For example, cannabinoids

558 like THCV have interesting pharmacological properties suggesting they may be dose-dependently

559 psychoactive (Pertwee 2008), with potential medicinal benefits (Bolognini et al. 2010). Chemotypes

560 expressing distinct ratios of minor cannabinoids and terpenes, with and without significant THC levels,

561 will likely elicit effects of interest to consumers and clinical researchers. Our results are consistent with

562 the notion that the full chemotype landscape of *Cannabis* has yet to be filled in (Figure 10).



563

**Figure 10: Potential scheme for classifying commercial *Cannabis* based on cannabinoid and terpene profiles.**
Flow chart showing a potential classification framework for commercial *Cannabis*. Level 1 represents cannabinoid
ratios and displays the three common THC:CBD chemotypes as well as novel cannabinoids that could be bred.
Level 2 represents terpene profiles and displays the three clusters we identified as well as other terpene combinations
which could come to exist. Terpene clusters overlap slightly to illustrate that terpenes in each cluster are not

22

569 mutually exclusive. Grey lines demonstrate a chemotype that may be possible (e.g., CBD-dominant and terpinolene-
570 dominant) but has not yet been observed.
571

572      In addition to mapping the chemical landscape of commercial *Cannabis* in the US, we also
573 quantified how well commonly used industry labels align with the chemical composition of samples. In
574 general, we found that industry labels are poorly or inconsistently aligned with the underlying chemistry.
575 In particular, the Indica/Hybrid/Sativa nomenclature does not reliably distinguish samples based on their
576 chemical content, making it highly unlikely that this widely used commercial labeling system is a reliable
577 indicator of systematically different effects. Marketing emphasizing Indica-labelled products as sedating
578 and Sativa-labelled products as energizing are not borne out by our analysis of the underlying chemistry.

579      We also examined the popular "strain names" commonly attached to products, which are used
580 commercially to reference cultivars purported to offer distinct effects. In particular, we quantified the
581 terpene profile consistency of THC-dominant products sharing the same strain name across different
582 producers. We modeled the situation where strain names are randomly applied to products, finding that
583 many strain names are more consistent from product-to-product, on average, than would be expected by
584 chance. However, we also observed a wide range of consistencies for all strain names, suggesting that
585 some are more homogeneous than others (Schwabe and McGlaughlin 2019), perhaps because these names
586 are more often attached to cultivars that are clonally propagated. These results indicate that while strain
587 names may be a better marker of product chemistry than the Indica/Sativa/Hybrid category labels, they
588 are far from ideal (Figure 8).

589      While commercial labels tended to have poor validity overall, we found evidence that certain strain
590 names and categories were statistically overrepresented within specific chemically defined clusters. In
591 particular, Cluster III samples (high terpinolene-myrcene) displayed an over-representation of Sativa-
592 labelled products. While certain strain names were over-represented in Clusters I and II, neither of these
593 Clusters displayed an over-representation of Indica or Sativa labels. Although the origins of this pattern
594 are unclear, one hypothesis is that it echoes patterns of phytochemistry that may have been more
595 distinctive prior to the long history of *Cannabis* hybridization in the US. It is conceivable, for example,
596 that certain cultivars commonly associated with "Sativa" lineages may have historically displayed a
597 chemotype reliably distinct from those in other lineages. Over time, hybridization and a lack of
598 standardized naming conventions may have decorrelated chemotaxonomic markers from the linguistic
599 labels used by cultivators. Thoroughly tracing which chemotypes tend to map to different lineages will
600 require datasets that combine both genotype and chemotype data for modern commercial cultivars and,
601 ideally, the landrace cultivars from which they descended (Clarke and Merlin 2016).

602 Medical *Cannabis* has been described as a "pharmacological treasure trove" (Mechoulam 2005)

603 due to the diversity of pharmacologically active compounds it harbors. *Cannabis*-derived formulations

604 and specific cannabinoids (namely THC and CBD) have demonstrated efficacy for conditions ranging

605 from chronic pain (Haroutounian et al. 2016) to childhood epilepsy (Lattanzi et al. 2018). Medical

606 *Cannabis* patients report an even wider array of conditions they believe *Cannabis* is efficacious for,

607 including mental health outcomes (Lucas et al. 2019). It has also been hypothesized that distinct

608 chemotypes of *Cannabis*, each with different ratios of cannabinoids and terpenes, may offer distinct

609 medical benefits and psychoactive effects (Russo 2019; Koltai and Namdar 2020). This hypothesized

610 "entourage effect" has been difficult to confirm experimentally due to onerous regulations that make it

611 challenging to execute *in vivo* studies with controlled administration of the myriad compounds found in

612 *Cannabis*.

613 The results of this study can serve as a guide for future research, including *in vitro* assays, animal

614 studies, and human trials. Studies seeking to falsify claims about the psychoactive and medical effects of

615 different *Cannabis* types should test chemical ratios that match those found commercially. If it is true that

616 different chemotypes of THC-dominant *Cannabis* reliably produce distinct psychoactive or medicinal

617 effects, then a sensible starting point is to design studies comparing the effects of common, distinctive

618 commercial chemotypes, such as those described by our cluster analysis (Figures 6-7). Likewise, if there

619 is any modulatory effect of specific cannabinoids or terpenes on the effects of THC, then this should be

620 tested using formulations designed to match the ratios that people choose to consume under 'ecological'

621 conditions.

622 While the present study represents the largest chemotaxonomic analysis of commercial *Cannabis*

623 to-date, there are important caveats. One is that the dataset we analyzed was an aggregation of lab data

624 from different states. We had no access to the genotype or the growing conditions for any of these samples

625 and important outstanding questions remain for how these factors relate to chemotype in *Cannabis*. It is

626 also possible one or more compounds that were not consistently measured in each region is an important

627 chemotaxonomic marker. State-level markets have different regulations which may influence the expertise

628 of commercial growers or the choice and development of *Cannabis* products. Finally, this dataset did not

629 include the variation found in hemp. An exciting area of future research will be to investigate these

630 questions using datasets that combine sample-level features about genotype, chemotype, and

631 environmental conditions.

632 Our results also have regulatory implications. For example, we observed a robust correlation

633 between total THC and total CBD levels for CBD-dominant *Cannabis* samples. Because the legal

634 definition of hemp in the US is based on an arbitrary threshold of total THC levels, the majority of CBD-

24

635    dominant samples would not be legally classified as hemp within the US, despite such samples being

636    characterized by low THC:CBD ratios distinct from those seen in high-THC samples (Figure 1-2).

637          Legal THC-dominant *Cannabis* products are marketed to consumers as if there are clear-cut

638    associations between a product's label and its psychoactive effects. This is deceptive, as there is currently

639    no clear scientific evidence for these claims and our results show that these labels have a tenuous

640    relationship to the underlying chemistry. In contrast to other widely used but federally regulated plants

641    (e.g., corn and other crops regulated by the Federal Seed Act), there are no enforced rules for the naming

642    of *Cannabis* varieties. This stems from the fact that *Cannabis* is not federally legal in the US, which

643    prevents an overarching, enforceable naming standard from emerging. As a consequence, legacy

644    classification systems inherited from the illicit market have persisted with unwarranted trust in the

645    provenance and predictability of products' effects.

646          We have shown that in the US, multiple, distinct chemotypes of commercial *Cannabis* are reliably

647    present across regions. Due to the chemical complexity of these products, which may contain dozens of

648    pharmacologically active compounds with potentially psychoactive or medicinal effects, we believe it is

649    in the public interest to devise a classification system and naming conventions that reflect the true

650    chemotaxonomic diversity of this plant. The general approach we have used in this study can serve as a

651    basic guide for cannabis product segmentation and classification rooted in product chemistry. Consumer-

652    facing labelling systems should be grounded in such an approach so that consumers can be guided to

653    products with reliably different sensory and psychoactive attributes.

654

655    **MATERIAL & METHODS**

656    **Data Collection**

657          The data analyzed in this paper was shared by Leafly, a technology company in the legal cannabis

658    industry. Leafly made a variety of data available as part of a data sharing program where university-

659    affiliated researchers can access data for research purposes with the intent to publish results in peer-

660    reviewed scientific journals. The data Leafly made available included laboratory testing data (cannabinoid

661    and terpene profiles; see below) as well as metrics related to consumer behavior and preferences,

662    including: normalized values of the number of unique views to each of the web pages within its online,

663    consumer-facing strain database; consumer ratings and common categorical designations associated with

664    commercial strain names (Indica, Hybrid, or Sativa); crowd-sourced metrics related to the perceived

665    flavors and effects of associated with popular strain names, derived from online consumer reviews. For

666    the purposes of this study, we focused mainly on analyzing the laboratory testing data and its relationship

667    with popular commercial labelling systems (i.e. strain names and Indica/Hybrid/Sativa designations).

668    Laboratory testing data came from Leafly via partnerships they have with cannabis testing labs
669    across the US. Each lab consented to allowing researchers to analyze its data for academic research
670    purposes. Each laboratory dataset consisted of the complete set of cannabinoid and terpene compounds
671    measured by each lab within a given time period between December 2013 and January 2021. The name
672    of each lab is listed below, together with the US state their data was measured in and a link to their
673    websites, which contain more detailed information on their specific testing methodologies. Each lab used
674    different variations of High Performance Liquid Chromatography to measure cannabinoid levels and Gas
675    Chromatography (GC-FID or GC-MS) to measure terpene levels.

676    - CannTest, Alaska, http://www.canntest.com/
677    - Confidence Analytics, Washington, https://www.conflabs.com/
678    - ChemHistory, Oregon, https://chemhistory.com/
679    - Modern Canna Labs, Florida, https://www.moderncanna.com/
680    - PSI Labs, Michigan, https://psilabs.org/
681    - SC Labs, California, https://www.sclabs.com/

682

683    Leafly shared a single, standardized lab dataset composed of *Cannabis* flower samples that had
684    been tested for cannabinoid, or for both cannabinoid and terpene content. Raw cannabinoid acid,
685    cannabinoid, and terpene measurements had been converted to common units (% weight) together with
686    additional information for each sample: anonymized producer ID, test date, and the producer-given sample
687    name.

688    For each lab testing sample, Leafly included the strain name associated with each web page in its
689    online *Cannabis* strain database together with the popular industry category ("Indica," "Hybrid," or
690    "Sativa") associated with each strain name. The strain names from Leafly's database were matched to the
691    producer-given strain name of each flower sample (e.g. "blue-dream"), wherever such a match was found,
692    using a similar string-matching algorithm as described in Jikomes & Zoorob (2018), supplemented with a
693    human expert-supplied dictionary used to standardize names with common variations (e.g. "SLH" =
694    "super-lemon-haze," "GDP" = "granddaddy-purple," and so on). In total, 81.5% of samples were attached
695    to popular strain names and 73.4% additionally attached to a Indica/Hybrid/Sativa label, with the
696    remainder labelled as "Unknown."

697

698    **Technologies Used**

699    All data cleaning and analysis for this paper was performed using the Python programming
700    language (Python Software Foundation, https://www.python.org) and utilized the following libraries:

26

701 NumPy, pandas, SciPy, and scikit-learn. All data visualizations were made using the Python libraries

702 Seaborn and Matplotlib.

703

**Data Processing: Raw Data Filtering & Outlier Removal**

705 The standardized dataset consisting of rows of lab data was cleaned and processed using custom

706 code in Python. A small number of duplicate rows were removed from the dataset ($n = 11$). We also

707 removed any samples with biologically implausible values (i.e. very high or low) for dried *Cannabis*,

708 which likely represent rare measurement anomalies or come from samples which do not truly represent

709 dried *Cannabis* flower (e.g. "shake" or other plant material different from the dried female inflorescence).

710 We used the following, conservative criteria: any single cannabinoid measured at over 40% (percent

711 weight; $n = 80$), or samples which had summed total cannabinoid measurements over 50% ($n = 2$); samples

712 which had null or 0.0 measurements for both total THC and total CBD ($n = 591$). The total number of

713 samples dropped from the dataset was 684, or 0.75% of the raw dataset. The final number of samples was

714 89,923.

715 Terpene data was also removed for samples which had a terpene measurement variance less than

716 0.001 ($n = 2,048$), samples which had any single terpene measurement over 5% ($n = 8$), or for samples

717 which had over 10 measurements equalling zero among the 14 most common terpenes ($n = 2,178$). The

718 total number of samples which had terpene data removed was 4,234, or 9% of samples having any terpene

719 data. The final number of samples with terpene data was 42,843, or 47.6% of the final dataset. The reason

720 that many laboratory testing samples contain only cannabinoid measurements is that terpene levels are

721 generally not legally required to be measured. Nonetheless, we were still left with 42,843 samples with

722 terpene measurements attached, which to our knowledge is the largest such dataset of commercial

723 *Cannabis* analyzed to date.

724

**Data Processing: Total Cannabinoid Levels**

726 Total cannabinoid levels were calculated from the raw cannabinoid and cannabinoid acid values

727 attached to each flower sample. This widely used convention calculates the total levels of a cannabinoid

728 found in a *Cannabis* product assuming complete decarboxylation of a cannabinoid acid to its

729 corresponding cannabinoid. For total THC, the formula is:

730 Total THC = (0.877 * THCA) + THC

731

732 0.877 is a scaling factor which accounts for the difference in molecular weight between raw cannabinoid

733 and cannabinoid acid values for THC, CBD, CBG, CBC, CBN, CBT, and delta-8 THC. The equivalent

734  formula, with the scaling factor of 0.8668, was used to calculate total cannabinoid levels for THCV and
735  CBDV.
736

**Data Processing: THC:CBD Chemotypes**

738  Following past work (Hillig and Mahlberg 2004; Jikomes and Zoorob 2018), we classified all
739  flower samples as THC-dominant, CBD-dominant, or Balanced THC:CBD based on the THC:CBD ratio
740  of the sample. THC-dominant samples are those with a 5:1 THC:CBD or higher, CBD-dominant samples
741  are those with a 1:5 THC:CBD or lower, and Balanced THC:CBD are in between.
742

**Data Analysis: Cannabinoid and Terpene Analysis**

744  Given that cannabis testing is not standardized nationally, each lab had a unique set of
745  cannabinoids and terpenes that they measured. Because of this, we established a list of compounds
746  common across every lab and used these in our main analyses. These compounds were:

747  ● Common Cannabinoids:
748  ○ Tetrahydrocannabinol (THC)
749  ○ Cannabidiol (CBD)
750  ○ Cannabigerol (CBG)
751  ○ Cannabichromene (CBC)
752  ○ Cannabinol (CBN)
753  ○ Tetrahydrocannabivarin (THCV)
754

755  ● Common Terpenes:
756  ○ Bisabolol
757  ○ Camphene
758  ○ β-Caryophyllene (Caryophyllene)
759  ○ α-Humulene (Humulene)
760  ○ Limonene
761  ○ Linalool
762  ○ β-Myrcene (Myrcene)
763  ○ cis- and trans-Nerolidol (Nerolidol)
764  ○ α-, β-, cis-, and trans-Ocimene (Ocimene)
765  ○ α-Pinene
766  ○ β-Pinene

767            ○   α-Terpinene

768            ○   γ-Terpinene

769            ○   Terpinolene

770

771       In the case of polar plots used to describe basic terpene profiles, α-pinene and β-pinene were

772 summed together and shown as "pinene" (see figures 7D-F and 8D). For certain terpenes (ocimene and

773 nerolidol), some labs measured individual isomers, and some reported a single total sum. In our main

774 analyses using data aggregated across labs, we summed across cis- and trans-nerolidol, and across α-, β-,

775 cis-, and trans-ocimene.

776

777 **Data Analysis: Sample- vs. Product-level Analysis**

778       Most of the analysis was conducted on the sample-level, meaning the data analyzed were the

779 individual *Cannabis* flower samples labs received and measured. We conducted some analyses at the

780 product-level. A product represents the average cannabinoid and terpene measurements for all strain

781 name-anonymized producer combinations. For example, Producer 101 might have 15 separate samples

782 attached to the name "blue-dream" that were submitted over some period of time. For product-level

783 analyses (Figures 5E-F, 7D-F, 8AB-E, and 9A-D), we averaged across such samples for each unique

784 combination of Producer IDs and strain names. THC:CBD chemotype was assigned to products based on

785 the average total THC and CBD values.

786

787 **Data Analysis: Statistics**

788       When performing statistical tests, we opted for statistical tests that do not depend on assumptions

789 about the distribution of the underlying data. For comparing groups, we used the Welch's t-test, which

790 does not assume equal population variances. For correlations, we computed Spearman's rank correlation

791 coefficient by default, as it provides a nonparametric measure of correlation. Any samples with null values

792 among the variables being analyzed were excluded in the calculation. Significance levels were corrected

793 using the most conservative Bonferroni correction to adjust for multiple comparisons, when applicable.

794 All p-values reported in the figures and text as significant are significant at the particular corrected alpha

795 level. Stars in figures (*, **, ***) correspond to the alpha levels 0.01, 0.001, and 0.0001 (with Bonferroni

796 correction), respectively. Due to the large sample sizes in our dataset, we tended to obtain very small p-

797 values that vary by many orders of magnitude. In these cases, p-values are reported as $< 0.0001$ (with

798 Bonferroni correction).

799       With sufficiently large sample sizes, statistically significant p-values can be found even when

800   differences are negligible. For this reason, we report effect sizes in addition to the p-values obtained from

801   Welch's t-test. We used an adjusted version of Cohen's d ("d-prime") in order to estimate the effect size

802   for independent samples without the assumption of equal variances (Navarro 2020).

803       This version averages the two population variances:

804

805
$$d' = \frac{X_1 - X_2}{\sqrt{\dfrac{\sigma_1{}^2 + \sigma_2{}^2}{2}}}$$

806

807 **Data Analysis: Figure 1**

808       The total levels for the six common cannabinoids were visualized as combination violin and box

809   plots. A scatter plot and a histogram of the relationship between total THC and total CBD were visualized

810   with the THC:CBD chemotypes color-coded. Principal component analysis (PCA) was run on the

811   normalized values of the six common cannabinoids (i.e., the % of measured common cannabinoids). Null

812   values were filled with zeros. A PCA biplot was created to visualize the PCA scores of the samples and

813   the weight of each cannabinoid on the first two principal components.

814

815 **Data Analysis: Figure 2**

816       The data was filtered by each of the three chemotype classes identified in Figure 1 (THC-dominant,

817   CBD-dominant, and balanced THC:CBD). Pairwise scatterplots for each permutation of the three most

818   abundant cannabinoids (THC, CBD, CBG) were made for the three THC:CBD chemotype classes. No

819   additional filtering or outlier removal was performed. The resulting nine plots are visualized in Figure 2.

820   The Spearman rank correlation for each cannabinoid relationship in each class was computed to measure

821   the strength of the relationship. Statistical significance was evaluated after using the Bonferroni correction

822   for 9 multiple comparisons. All observed relationships were significant at the (corrected) $P < 0.0001$ level.

823

824 **Data Analysis: Figure 3**

825       The fourteen common terpenes were visualized for samples with terpene data in a combination

826   violin/box plot, ordered by median value, descending. The linear relationships between two pairs of

827   terpenes (α- and β-pinene, and β-caryophyllene and humulene) were quantified with a linear regression

828   and Spearman rank correlation. Statistical significance was evaluated after using the Bonferroni correction

829   for two multiple comparisons.

830

**Data Analysis: Figure 4**

The fourteen terpene levels were correlated with each other using a Spearman rank correlation. A clustermap visualization in Figure 4 combining a heatmap and hierarchical clustering visualizations was made. Because of the multiple pairwise comparisons (14 x 13 / 2 = 91), statistical significance was evaluated after using the Bonferroni correction for 91 multiple comparisons. Cells were colored by the strength of the relationship (bluer are stronger negative correlations, redder are stronger positive correlations) and annotated with the correlation value only if the relationship was significant at the (corrected) $p < 0.05$ level. Only four compound combinations had non-significant corrected relationships: (1) terpinolene-nerolidol, (2) terpinolene-humulene, (3) myrcene-bisabolol, and (4) ocimene-camphene. The distances between clusters were evaluated using the "average" method in the "hierarchy.linkage" function and the "euclidean" function was used as a distance metric.

The clusters recovered by the clustermap visualization can also be represented as a network where the nodes are the terpenes and the (weighted) edges are the correlations. Because nearly all compound combinations have statistically significant correlations (even after Bonferroni correction), the resulting network would be (nearly) completely connected. To sparsify the network for visualization purposes, the correlation values were thresholded to greater than or equal to 0.10 to show the strongest relationships. There were 38 remaining edges after this thresholding procedure. This threshold value was chosen through qualitative iteration to generate a network that preserves all 14 compounds but is sufficiently sparse to visually recover the clusters identified in Figure 4A. The network was visualized using a spring-embedding layout algorithm and visualized using the "networkx" library in Python.

851

**Data Analysis: Figure 5**

Principal component analysis (PCA) was run on the normalized values of the fourteen common terpenes (i.e., the % of measured common terpenes) on all samples with terpene data. Null values were filled with zeros. A bar plot was created to visualize how much variation each principal component captured in the data. PCA biplots were created to visualize the PCA scores of the samples and the weight of each terpene on the first three principal components (Figure 5X-Y).

Sample level data was averaged across strain name/producer ID pairs to create a product level dataset. Pairwise cosine distances of terpene profiles were calculated for products in each chemotype. We then averaged the cosine distances across each product, so each product had an associated average cosine distance. These values were plotted in a violin/box plot (Figure 5E). Welch's t-tests and effect sizes were calculated between each chemotype. Statistical significance was evaluated after using the Bonferroni

863    correction for three multiple comparisons. The top terpene among the 14 common terpenes was found for

864    each product. If the most abundant terpene was not either myrcene, caryophyllene, limonene, terpinolene,

865    alpha pinene, or ocimene, the top terpene was listed as "other" (Figure 5F).

866

**Data Analysis: Figure 6**

867

868    For figures 6A-F, the sample level data was filtered to include only THC-dominant samples with

869    terpene data. Terpene data were normalized to be % of measured common terpenes. Null values were

870    filled with zeros. PCA was run on these normalized values and then plotted.

871    Silhouette coefficients for each sample were calculated using the mean nearest-cluster Euclidean

872    distance (b) minus the mean intra-cluster Euclidean distance (a), divided by max (a,b). This value

873    measures how similar a sample is to its labeled cluster compared to other clusters. The individual

874    silhouette sample scores plotted were obtained from a random subsample of the data (n=10,000) due to

875    graphic memory limitations, however the average silhouette score displayed on the figure was obtained

876    using the full filtered dataset.

877    We used the k-means clustering algorithm to segment THC-dominant samples based on terpene

878    profiles. To determine the optimal number of clusters we created an 'elbow plot', which plots a range of

879    number of clusters versus within-cluster sum of squared errors (Figure S5A). This revealed that the

880    optimal number of clusters to use was k = 3. K-means clustering was applied to the normalized dataset. A

881    color palette was created using the color of the most abundant terpene for each cluster's average terpene

882    profile. The correct choice of k can be ambiguous, so we also explored our cluster analysis for k=2 and

883    k=4 clusters (Figure S5B-C).

884

**Data Analysis: Figure 7**

885

886    To evaluate the difference between the labeling methods described above, silhouette scores

887    (described above) were calculated on the full dataset for the three different methods. Welch's t-tests and

888    effect sizes were calculated between these methods. Statistical significance was evaluated after using the

889    Bonferroni correction for three multiple comparisons.

890    A UMAP embedding (McInnes and Healy 2018) was run on the terpene data of THC-dominant

891    samples and color coded by k-means cluster label. The parameters for number of components and number

892    of neighbors were specified as 2 and 15, respectively. An interactive 3-D version of a similar product-

893    level UMAP can be found here: https://plotly.com/~cj.smith015/5/. Each data point can be hovered over

894    to reveal the following information: strain name, Indica/Hybrid/Sativa label, THC and CBD concentration,

895    dominant terpene, and k-means cluster label information

32

896      To illustrate a simple terpene profile, we ran k-means clustering (k = 3) on the product-level

897      dataset. α- and β-pinene were summed together. The normalized terpene values and total THC, CBD, and

898      CBG values from the THC-dominant product dataset were grouped by k-means cluster label and averaged.

899      Polar plots were constructed based on the average terpene profiles and limited to eight terpenes to help

900      with visual legibility. The terpene profiles of the top 25 products in each cluster with the most samples

901      were drawn in grey behind the cluster-level average.

902

### Data Analysis: Figure 8

904      To quantify consistency between products attached with the same name we needed to ensure that

905      the underlying data contained multiple samples per producer ID and several unique producer IDs each.

906      We used the following thresholds: to be included, a strain name must be linked to at least five producers

907      with at least five samples from each producer. If the strain met this threshold, we included all samples of

908      that strain in our examination, averaging all samples linked to each unique producer ID to create product

909      averages. 41 strain names met this threshold. Due to the predominance of THC-dominant samples in the

910      dataset, all strain names in the list happened to be THC-dominant. Measures of strain name popularity

911      were supplied by Leafly in the form of normalized values for how many unique views each page of its

912      public strain database received.

913      In figure 8B, a correlation matrix was constructed on the terpene values of THC-dominant samples

914      for the ten strain names attached to the most samples. The samples were put in descending order based on

915      the number of samples, and within each strain name, ordered by producer ID. Pairwise cosine similarity

916      scores were calculated on the samples and plotted as a heat map with a Gaussian filter for visualization

917      purposes.

918      Cosine similarities were calculated for the terpene profiles of products for each strain name, then

919      averaged to assign a mean similarity score to each product (identity values of 1 were replaced with nulls

920      so as to not artificially increase the average). A violin/box plot was created with these similarity scores,

921      ordered by median value. The dashed line in figure 8C represents the average similarity score one would

922      expect if strain names were randomly assigned, obtained by running a bootstrap simulation where strain

923      names were shuffled across the product IDs. Average similarity scores for products were calculated based

924      on these randomized strain names. Those scores were then averaged to give each (randomized) strain

925      name a similarity score. A weighted average was created by taking the randomized strain-level similarity

926      scores and weighing them by the number of products associated with each randomized strain name. This

927      process was repeated 200 times and the mean of this distribution was calculated and displayed as the

928      dashed line. Welch's t-tests and effect sizes were calculated comparing the similarity scores for each strain

929 to the bootstrapped distribution of average randomized strain-level similarity scores. Statistical
930 significance was evaluated after using the Bonferroni correction for 41 multiple comparisons.

931        A UMAP embedding was run on the normalized terpene data of the entire THC-dominant product
932 dataset and color coded by k-means cluster label, k = 3. The parameters for number of components and
933 number of neighbors were specified as 2 and 15, respectively.

934

**Data Analysis: Figure 9**

936        Using the THC-dominant product dataset with k-means clustering (k = 3), a UMAP embedding
937 was run on the normalized terpene data and color coded by Indica/Sativa/Hybrid labels.

938        Excluding products without an associated Indica/Sativa/Hybrid label, the percentage of
939 Indica/Sativa/Hybrid labels for products was found for each k-means cluster label. Chi-squared tests were
940 calculated comparing these percentages with the overall percentages. Statistical significance was
941 evaluated after using the Bonferroni correction for three multiple comparisons.

942        Using the list of 41 strains obtained by the thresholds described for figure 8, the most frequent k-
943 means cluster label was identified for each strain name. The number of products with that cluster label
944 divided by the total number of products for that strain multiplied by 100 gave the percentage of products
945 in the top cluster. Up to seven strains in each cluster were displayed in the bar chart in figure 9D, ordered
946 by k-means cluster label and then by the percentage of products in the top cluster. The dashed line in
947 figure 9D represents the average percentage of products one would expect if strain names were randomly
948 assigned, obtained by running a bootstrap simulation where strain names were shuffled across the product
949 dataset, as described above for Figure 8. Welch's t-tests and effect sizes were calculated by comparing the
950 distribution of products in the top cluster for each strain to the bootstrapped distribution of average
951 percentage of randomized products in the top cluster. Statistical significance was evaluated after using the
952 Bonferroni correction for 41 multiple comparisons.

953

954

957

**Author contributions:** C.S. and B.K performed all data analysis and visualization. N.J., C.S., and B.K conceived all the analysis; D.V produced final figures; All authors contributed to manuscript preparation.
960

**Competing interests:** D.V. is the founder and president of the non-profit organization Agricultural Genomics Foundation, and the sole owner of CGRI, LLC. N.J. is employed by Leafly Holdings, Inc.

963 Leafly allowed N.J. to use some professional time to oversee this research project and work on the
964 manuscript.
965

966 **Data and materials availability:** All code used to conduct analysis and generate figures can be made
967 available upon request. Lab data analyzed in the study can be made available with written consent from
968 each testing lab.
969

970
971 **References**
972

973 Abel, E. L. 2013. Marihuana: the first twelve thousand years. Springer Science & Business Media.
974 Adams, T. B. and S. V. Taylor. 2010. Safety evaluation of essential oils: a constituent-based approach.
975     CRC Press: London, UK.
976 Aizpurua-Olaizola, O., U. Soydaner, E. Öztürk, D. Schibano, Y. Simsir, P. Navarro, N. Etxebarria, and
977     A. Usobiaga. 2016. Evolution of the cannabinoid and terpene content during the growth of
978     Cannabis sativa plants from different chemotypes. Journal of natural products 79:324-331.
979 Allen, K. D., K. McKernan, C. Pauli, J. Roe, A. Torres, and R. Gaudino. 2019. Genomic
980     characterization of the complete terpene synthase gene family from Cannabis sativa. PloS one
981     14:e0222363.
982 Boggs, D. L., J. D. Nguyen, D. Morgenson, M. A. Taffe, and M. Ranganathan. 2018. Clinical and
983     preclinical evidence for functional interactions of cannabidiol and Δ 9-tetrahydrocannabinol.
984     Neuropsychopharmacology 43:142-154.
985 Bolognesi, M. L. 2019. Harnessing polypharmacology with medicinal chemistry. ACS Publications.
986 Bolognini, D., B. Costa, S. Maione, F. Comelli, P. Marini, V. Di Marzo, D. Parolaro, R. A. Ross, L. A.
987     Gauson, and M. G. Cascio. 2010. The plant cannabinoid Δ9-tetrahydrocannabivarin can decrease
988     signs of inflammation and inflammatory pain in mice. British journal of pharmacology 160:677-
989     687.
990 Booth, J. K. and J. Bohlmann. 2019. Terpenes in Cannabis sativa–From plant genome to humans. Plant
991     Science 284:67-72.
992 Booth, J. K., J. E. Page, and J. Bohlmann. 2017. Terpene synthases from Cannabis sativa. Plos one
993     12:e0173911.
994 Booth, J. K., M. M. Yuen, S. Jancsik, L. L. Madilao, J. E. Page, and J. Bohlmann. 2020. Terpene
995     synthases and terpene variation in Cannabis sativa. Plant physiology 184:130-147.
996 Borrelli, F., E. Pagano, B. Romano, S. Panzera, F. Maiello, D. Coppola, L. De Petrocellis, L. Buono, P.
997     Orlando, and A. A. Izzo. 2014. Colon carcinogenesis is inhibited by the TRPM8 antagonist
998     cannabigerol, a Cannabis-derived non-psychotropic cannabinoid. Carcinogenesis:bgu205.
999 Calhoun, S. R., G. P. Galloway, and D. E. Smith. 1998. Abuse potential of dronabinol (Marinol®).
000     Journal of psychoactive drugs 30:187-196.
001 Carter, G. T., A. M. Flanagan, M. Earleywine, D. I. Abrams, S. K. Aggarwal, and L. Grinspoon. 2011.
002     Cannabis in palliative medicine: improving care and reducing opioid-related morbidity.
003     American Journal of Hospice and Palliative Medicine:1049909111402318.
004 Chakraborty, S., R. Minda, L. Salaye, A. M. Dandekar, S. K. Bhattacharjee, and B. J. Rao. 2013.
005     Promiscuity-based enzyme selection for rational directed evolution experiments. Enzyme
006     Engineering: Methods and Protocols:205-216.
007 Clarke, R. and M. Merlin. 2013. Cannabis: evolution and ethnobotany. Univ of California Press.
008 Clarke, R. C. and M. D. Merlin. 2016. Cannabis domestication, breeding history, present-day genetic
009     diversity, and future prospects. Critical reviews in plant sciences 35:293-327.
010 Dorrity, M. W., L. M. Saunders, C. Queitsch, S. Fields, and C. Trapnell. 2020. Dimensionality reduction
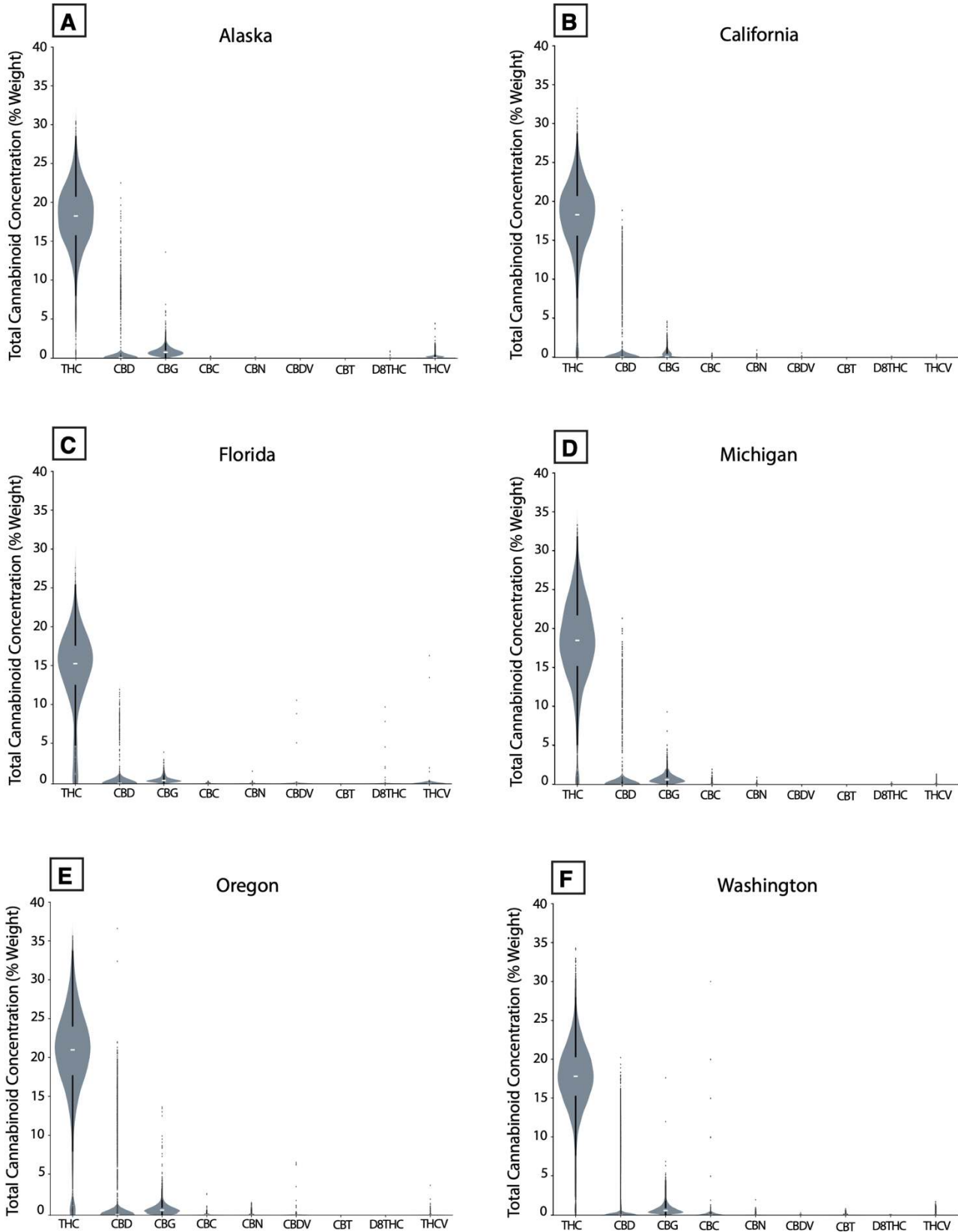011     by UMAP to visualize physical and genetic interactions. Nature communications 11:1-6.

ElSohly, M. A. and D. Slade. 2005. Chemical constituents of marijuana: the complex mixture of natural cannabinoids. Life sciences 78:539-548.

Elzinga, S., J. Fischedick, R. Podkolinski, and J. Raber. 2015. Cannabinoids and terpenes as chemotaxonomic markers in cannabis. Nat Prod Chem Res 3:2.

Franco, O. L. 2011. Peptide promiscuity: an evolutionary concept for plant defense. FEBS letters 585:995-1000.

Gertsch, J., M. Leonti, S. Raduner, I. Racz, J.-Z. Chen, X.-Q. Xie, K.-H. Altmann, M. Karsak, and A. Zimmer. 2008. Beta-caryophyllene is a dietary cannabinoid. Proceedings of the National Academy of Sciences 105:9099-9104.

Goodman, S., E. Wadsworth, C. Leos-Toro, and D. Hammond. 2020. Prevalence and forms of cannabis use in legal vs. illegal recreational cannabis markets. International Journal of Drug Policy 76:102658.

Haroutounian, S., Y. Ratz, Y. Ginosar, K. Furmanov, F. Saifi, R. Meidan, and E. Davidson. 2016. The effect of medicinal cannabis on pain and quality-of-life outcomes in chronic pain. The Clinical journal of pain 32:1036-1043.

Hart, C. L., W. Van Gorp, M. Haney, R. W. Foltin, and M. W. Fischman. 2001. Effects of acute smoked marijuana on complex cognitive performance. Neuropsychopharmacology 25:757-765.

Hazekamp, A. and J. T. Fischedick. 2012. Cannabis-from cultivar to chemovar. Drug testing and analysis 4:660-667.

Hazekamp, A., K. Tejkalová, and S. Papadimitriou. 2016. Cannabis: from cultivar to chemovar II—a metabolomics approach to Cannabis classification. Cannabis and Cannabinoid Research 1:202-215.

Henry, P., A. Hilyard, S. Johnson, and C. Orser. 2018. Predicting chemovar cluster and variety verification in vegetative cannabis accessions using targeted single nucleotide polymorphisms. PeerJ Preprints 6:e27442v27441.

Hillig, K. W. 2004. A chemotaxonomic analysis of terpenoid variation in Cannabis. Biochemical Systematics and Ecology 32:875-891.

Hillig, K. W. and P. G. Mahlberg. 2004. A chemotaxonomic analysis of cannabinoid variation in Cannabis (Cannabaceae). American Journal of Botany 91:966-975.

Hutchison, K. E., L. C. Bidwell, J. M. Ellingson, and A. D. Bryan. 2019. Cannabis and Health Research: Rapid Progress Requires Innovative Research Designs. Value in Health.

Izzo, A. A., R. Capasso, G. Aviello, F. Borrelli, B. Romano, F. Piscitelli, L. Gallo, F. Capasso, P. Orlando, and V. Di Marzo. 2012. Inhibitory effect of cannabichromene, a major non-psychotropic cannabinoid extracted from Cannabis sativa, on inflammation-induced hypermotility in mice. British journal of pharmacology 166:1444-1460.

Jikomes, N. and M. Zoorob. 2018. The cannabinoid content of legal cannabis in Washington state varies systematically across testing facilities and popular consumer products. Scientific reports 8:4519.

Koltai, H. and D. Namdar. 2020. Cannabis Phytomolecule'Entourage': From Domestication to Medical Use. Trends in Plant Science.

Kovalchuk, I., M. Pellino, P. Rigault, R. van Velzen, J. Ebersbach, J. R. Ashnest, M. Mau, M. Schranz, J. Alcorn, and R. Laprairie. 2020. The Genomics of Cannabis and Its Close Relatives. Annual Review of Plant Biology 71.

Langenheim, J. H. 1994. Higher plant terpenoids: a phytocentric overview of their ecological roles. Journal of chemical ecology 20:1223-1280.

Laprairie, R., A. Bagher, M. Kelly, and E. Denovan-Wright. 2015. Cannabidiol is a negative allosteric modulator of the cannabinoid CB1 receptor. British journal of pharmacology 172:4790-4805.

Lattanzi, S., F. Brigo, E. Trinka, G. Zaccara, C. Cagnetti, C. Del Giovane, and M. Silvestrini. 2018. Efficacy and safety of cannabidiol in epilepsy: a systematic review and meta-analysis. Drugs 78:1791-1804.

061  Lewis, M. A., E. B. Russo, and K. M. Smith. 2018. Pharmacological foundations of cannabis
062       chemovars. Planta medica 84:225-233.
063  Lucas, P., E. P. Baron, and N. Jikomes. 2019. Medical cannabis patterns of use and substitution for
064       opioids & other pharmaceutical drugs, alcohol, tobacco, and illicit substances; results from a
065       cross-sectional survey of authorized patients. Harm reduction journal 16:1-11.
066  Lynch, R. C., D. Vergara, S. Tittes, K. White, C. J. Schwartz, M. J. Gibbs, T. C. Ruthenburg, K.
067       deCesare, D. P. Land, and N. C. Kane. 2016. Genomic and Chemical Diversity in Cannabis.
068       Critical Reviews in Plant Sciences 35:349-363.
069  Magagnini, G., G. Grassi, and S. Kotiranta. 2018. The effect of light spectrum on the morphology and
070       cannabinoid content of Cannabis sativa L. Medical Cannabis and Cannabinoids 1:19-27.
071  McInnes, L. and J. Healy. 2018. Umap: Uniform manifold approximation and projection for dimension
072       reduction. arXiv preprint arXiv:1802.03426.
073  McPartland, J. M., M. Duncan, V. Di Marzo, and R. G. Pertwee. 2015. Are cannabidiol and Δ9-
074       tetrahydrocannabivarin negative modulators of the endocannabinoid system? A systematic
075       review. British journal of pharmacology 172:737-753.
076  McPartland, J. M. and E. B. Russo. 2001. Cannabis and cannabis extracts: greater than the sum of their
077       parts? Journal of Cannabis Therapeutics 1:103-132.
078  McPartland, J. M. and E. Small. 2020. A classification of endangered high-THC cannabis (Cannabis
079       sativa subsp. indica) domesticates and their wild relatives. PhytoKeys 144:81.
080  Mechoulam, R. 2005. Plant cannabinoids: a neglected pharmacological treasure trove. British journal of
081       pharmacology 146:913-915.
082  Mudge, E. M., P. N. Brown, and S. J. Murch. 2019. The terroir of Cannabis: terpene metabolomics as a
083       tool to understand Cannabis sativa selections. Planta medica 85:781-796.
084  Navarro, D. 2020. Effect Size [Internet]. University of New South Wales. University of New South
085       Wales, Available from: https://stats.libretexts.org/@go/page/8266.
086  Onofri, C., E. P. M. de Meijer, and G. Mandolino. 2015. Sequence heterogeneity of cannabidiolic-and
087       tetrahydrocannabinolic acid-synthase in Cannabis sativa L. and its relationship with chemical
088       phenotype. Phytochemistry.
089  Orser, C., S. Johnson, M. Speck, A. Hilyard, and I. Afia. 2017. Terpenoid Chemoprofiles Distinguish
090       Drug-type Cannabis sativa L. Cultivars in Nevada. Natural Products Chemistry and Research 6.
091  Page, J. E. and J. M. Stout. 2017. Cannabichromenic acid synthase from Cannabis sativa. Google
092       Patents.
093  Pertwee, R. 2008. The diverse CB1 and CB2 receptor pharmacology of three plant cannabinoids: Δ9-
094       tetrahydrocannabinol, cannabidiol and Δ9-tetrahydrocannabivarin. British journal of
095       pharmacology 153:199-215.
096  Potter, D. 2004. Growth and morphology of medicinal cannabis. The Medicinal Uses of Cannabis and
097       Cannabinoids. Pharmaceutical Press, London.
098  Potter, D. 2009. The Propagation, Characterisation and Optimisation of Cannabis Sativa L as a
099       Phytopharmaceutical. King's College London.
100  Proschak, E., H. Stark, and D. Merk. 2018. Polypharmacology by design: a medicinal chemist's
101       perspective on multitargeting compounds. Journal of medicinal chemistry 62:420-444.
102  Reimann-Philipp, U., M. Speck, C. Orser, S. Johnson, A. Hilyard, H. Turner, A. J. Stokes, and A. L.
103       Small-Howard. 2019. Cannabis Chemovar Nomenclature Misrepresents Chemical and Genetic
104       Diversity; Survey of Variations in Chemical Profiles and Genetic Markers in Nevada Medical
105       Cannabis Samples. Cannabis and Cannabinoid Research.
106  Riboulet-Zemouli, K. 2020. 'Cannabis' ontologies I: Conceptual issues with Cannabis and cannabinoids
107       terminology. Drug Science, Policy and Law 6:2050324520945797.
108  Ross, S. A. and M. A. ElSohly. 1997. CBN and Δ 9-THC concentration ratio as an indicator of the age of
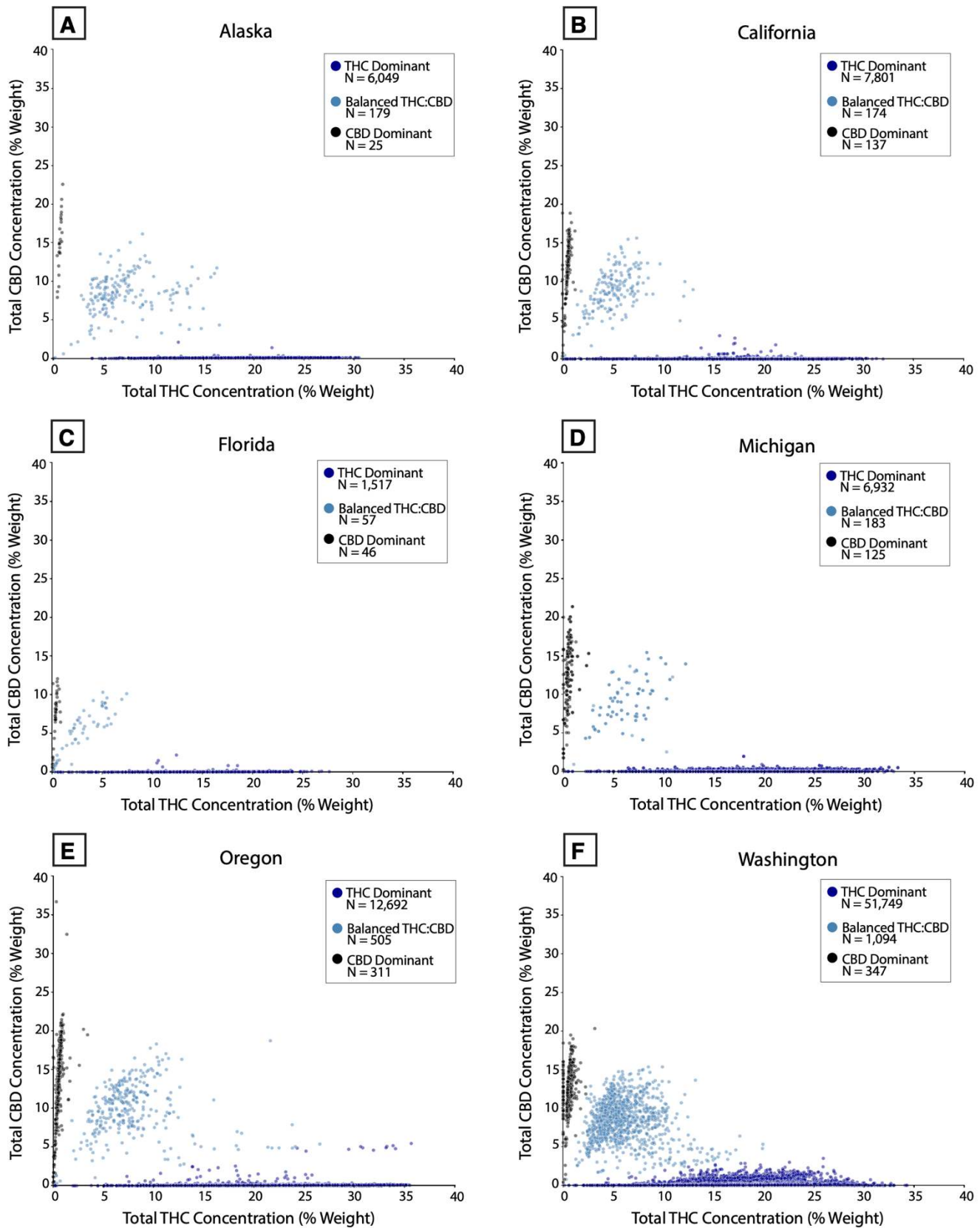109       stored marijuana samples. Bulletin on Narcotics 49:139-139.

110  Russo, E. B. 2007. History of cannabis and its preparations in saga, science, and sobriquet. Chemistry &
111       Biodiversity 4:1614-1648.
112  Russo, E. B. 2011. Taming THC: potential cannabis synergy and phytocannabinoid-terpenoid entourage
113       effects. British Journal of Pharmacology 163:1344-1364.
114  Russo, E. B. 2019. The Case for the Entourage Effect and Conventional Breeding of Clinical Cannabis:
115       No "Strain," No Gain. Frontiers in Plant Science 9.
116  Sawler, J., J. M. Stout, K. M. Gardner, D. Hudson, J. Vidmar, L. Butler, J. E. Page, and S. Myles. 2015.
117       The Genetic Structure of Marijuana and Hemp. PloS one 10:e0133292.
118  Schwabe, A. L. and M. E. McGlaughlin. 2019. Genetic tools weed out misconceptions of strain
119       reliability in Cannabis sativa: implications for a budding industry. Journal of Cannabis Research
120       1:3.
121  Sirikantaramas, S., F. Taura, Y. Tanaka, Y. Ishikawa, S. Morimoto, and Y. Shoyama. 2005.
122       Tetrahydrocannabinolic acid synthase, the enzyme controlling marijuana psychoactivity, is
123       secreted into the storage cavity of the glandular trichomes. Plant and Cell Physiology 46:1578-
124       1582.
125  Solowij, N., S. Broyd, L.-m. Greenwood, H. van Hell, D. Martelozzo, K. Rueb, J. Todd, Z. Liu, P.
126       Galettis, and J. Martin. 2019. A randomised controlled trial of vaporised Δ 9-
127       tetrahydrocannabinol and cannabidiol alone and in combination in frequent and infrequent
128       cannabis users: acute intoxication effects. European archives of psychiatry and clinical
129       neuroscience 269:17-35.
130  Steigerwald, S., P. O. Wong, A. Khorasani, and S. Keyhani. 2018. The Form and Content of Cannabis
131       Products in the United States. Journal of General Internal Medicine 33:1426-1428.
132  Swift, W., A. Wong, K. M. Li, J. C. Arnold, and I. S. McGregor. 2013. Analysis of cannabis seizures in
133       NSW, Australia: cannabis potency and cannabinoid profile. PloS one 8:e70052.
134  Trofin, I. G., G. Dabija, D. I. Vaireanu, and L. Filipescu. 2012. The influence of long-term storage
135       conditions on the stability of cannabinoids derived from cannabis resin. Rev Chim Bucharest
136       63:422-427.
137  Turner, C. E. and M. A. Elsohly. 1979. Constituents of cannabis sativa L. XVI. A possible
138       decomposition pathway of Δ9-tetrahydrocannabinol to cannabinol. Journal of heterocyclic
139       chemistry 16:1667-1668.
140  Valliere, M. A., T. P. Korman, N. B. Woodall, G. A. Khitrov, R. E. Taylor, D. Baker, and J. U. Bowie.
141       2019. A cell-free platform for the prenylation of natural products and application to cannabinoid
142       production. Nature communications 10:565.
143  van Velzen, R. and M. E. Schranz. 2020. Origin and evolution of the cannabinoid oxidocyclase gene
144       family. bioRxiv.
145  Venderová, K., E. Růžička, V. Voříšek, and P. Višňovský. 2004. Survey on cannabis use in Parkinson's
146       disease: subjective improvement of motor symptoms. Movement Disorders 19:1102-1106.
147  Vergara, D., H. Baker, K. Clancy, K. G. Keepers, J. P. Mendieta, C. S. Pauli, S. B. Tittes, K. H. White,
148       and N. C. Kane. 2016. Genetic and Genomic Tools for Cannabis sativa. Critical Reviews in Plant
149       Sciences 35:364-377.
150  Vergara, D., L. C. Bidwell, R. Gaudino, A. Torres, G. Du, T. C. Ruthenburg, K. deCesare, D. P. Land,
151       K. E. Hutchison, and N. C. Kane. 2017. Compromised External Validity: Federally Produced
152       Cannabis Does Not Reflect Legal Markets. Scientific Reports 7:46528.
153  Vergara, D., C. Feathers, E. L. Huscher, B. Holmes, J. A. Haas, and N. C. Kane. 2021a. Widely assumed
154       phenotypic associations in Cannabis sativa lack a shared genetic basis. PeerJ 9:e10672.
155  Vergara, D., R. Gaudino, T. Blank, and B. Keegan. 2020. Modeling cannabinoids from a large-scale
156       sample of Cannabis sativa chemotypes. PloS one 15:e0236878.

157  Vergara, D., E. L. Huscher, K. G. Keepers, R. M. Givens, C. G. Cizek, A. Torres, R. Gaudino, and N. C.
158      Kane. 2019. Gene copy number is associated with phytochemistry in Cannabis sativa. AoB
159      PLANTS 11:plz074.
160  Vergara, D., E. L. Huscher, K. G. Keepers, R. Pisupati, A. L. Schwabe, M. E. McGlaughlin, and N. C.
161      Kane. 2021b. Genomic evidence that governmentally produced Cannabis sativa poorly
162      represents genetic variation available in state markets. bioRxiv:431041.
163  Watts, G. 2006. Science commentary: Cannabis confusions. BMJ: British Medical Journal 332:175.
164  Zlebnik, N. E. and J. F. Cheer. 2016. Beyond the CB1 receptor: is cannabidiol the answer for disorders
165      of motivation? Annual review of neuroscience 39:1-17.
166
167
168

169 **Supplementary Materials**



170

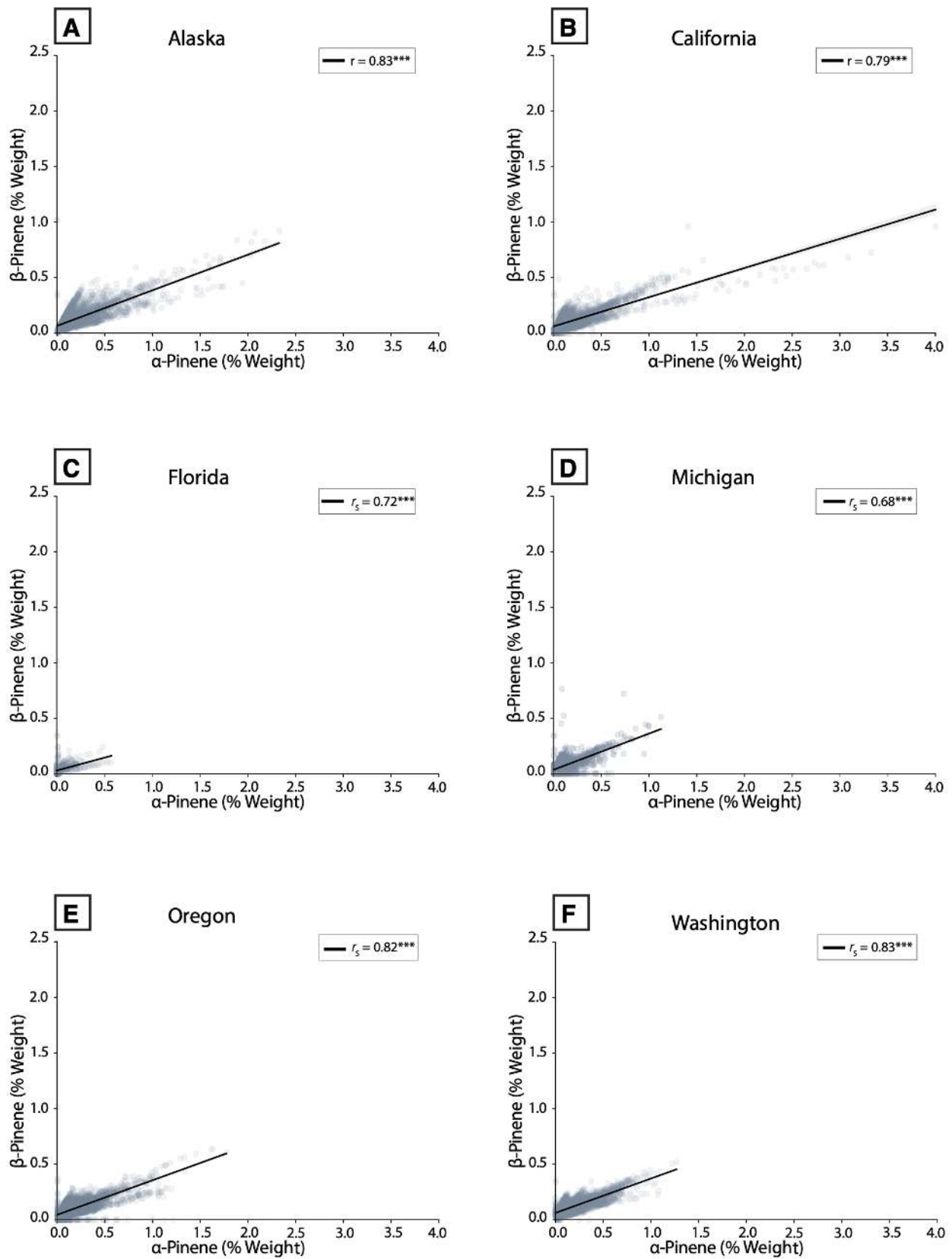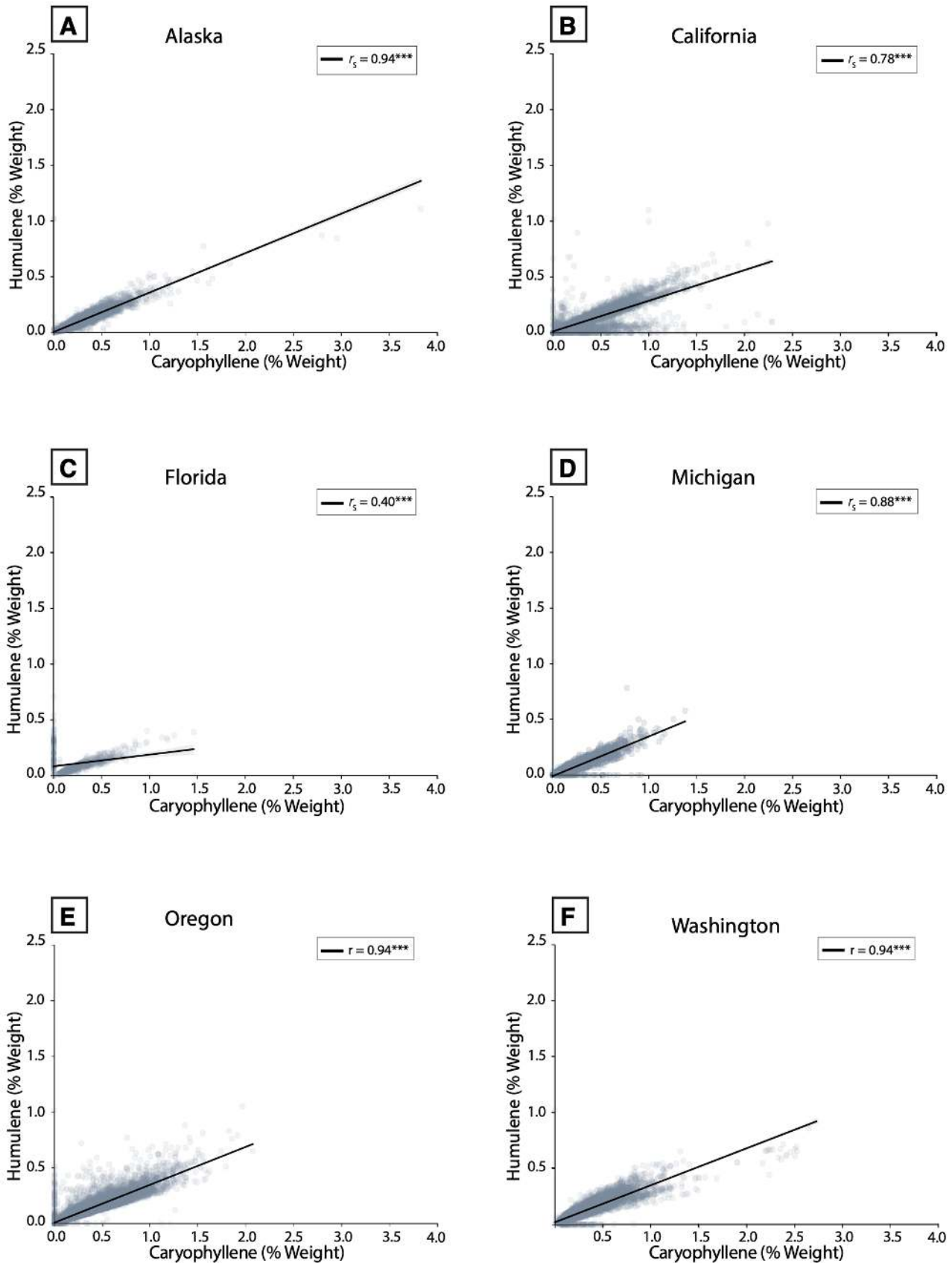171 **Figure S1:** Violin plot of distribution of all cannabinoids measured, by region.

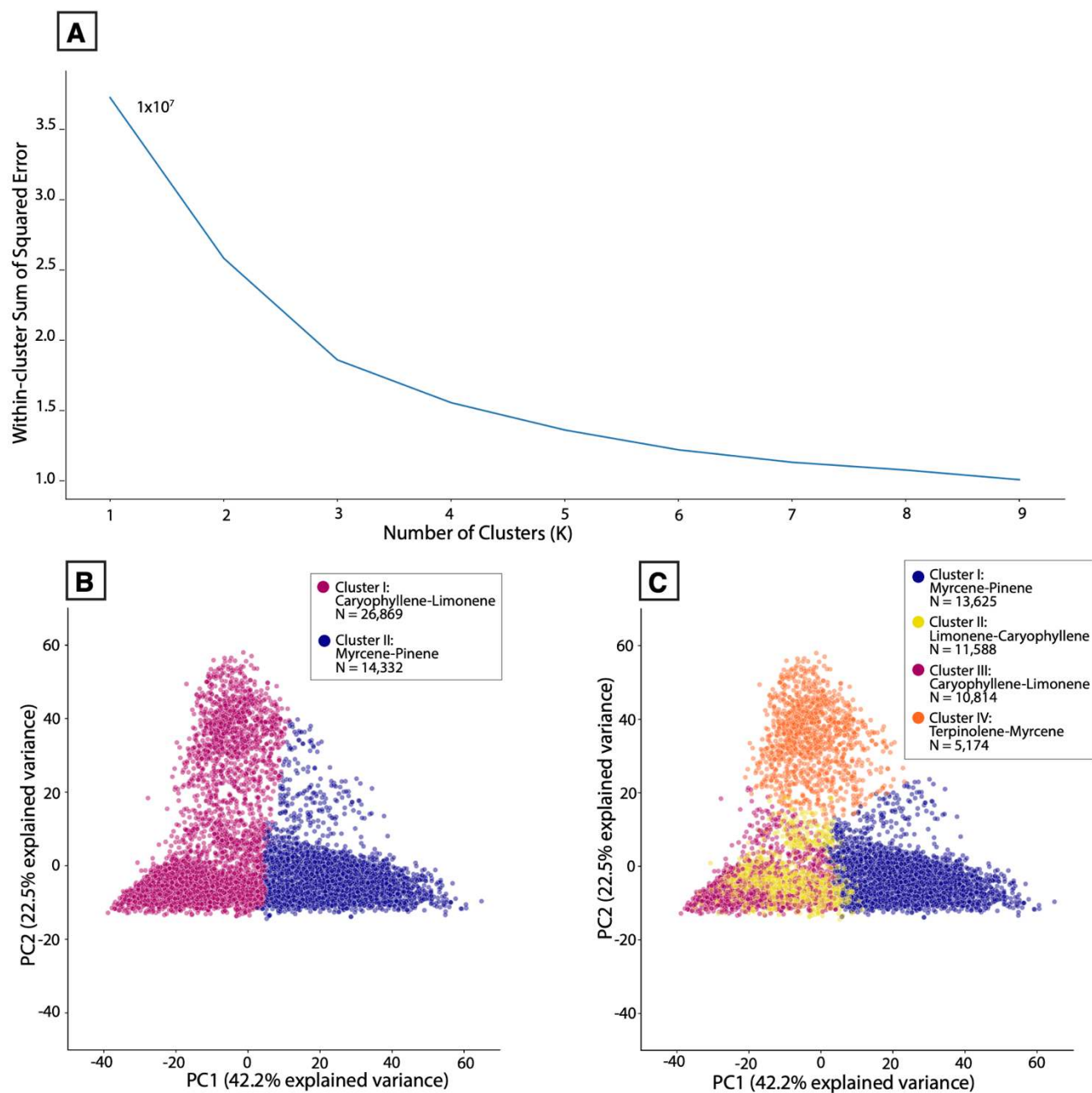**Figure S2:** Total THC vs. Total CBD levels, by region.

175

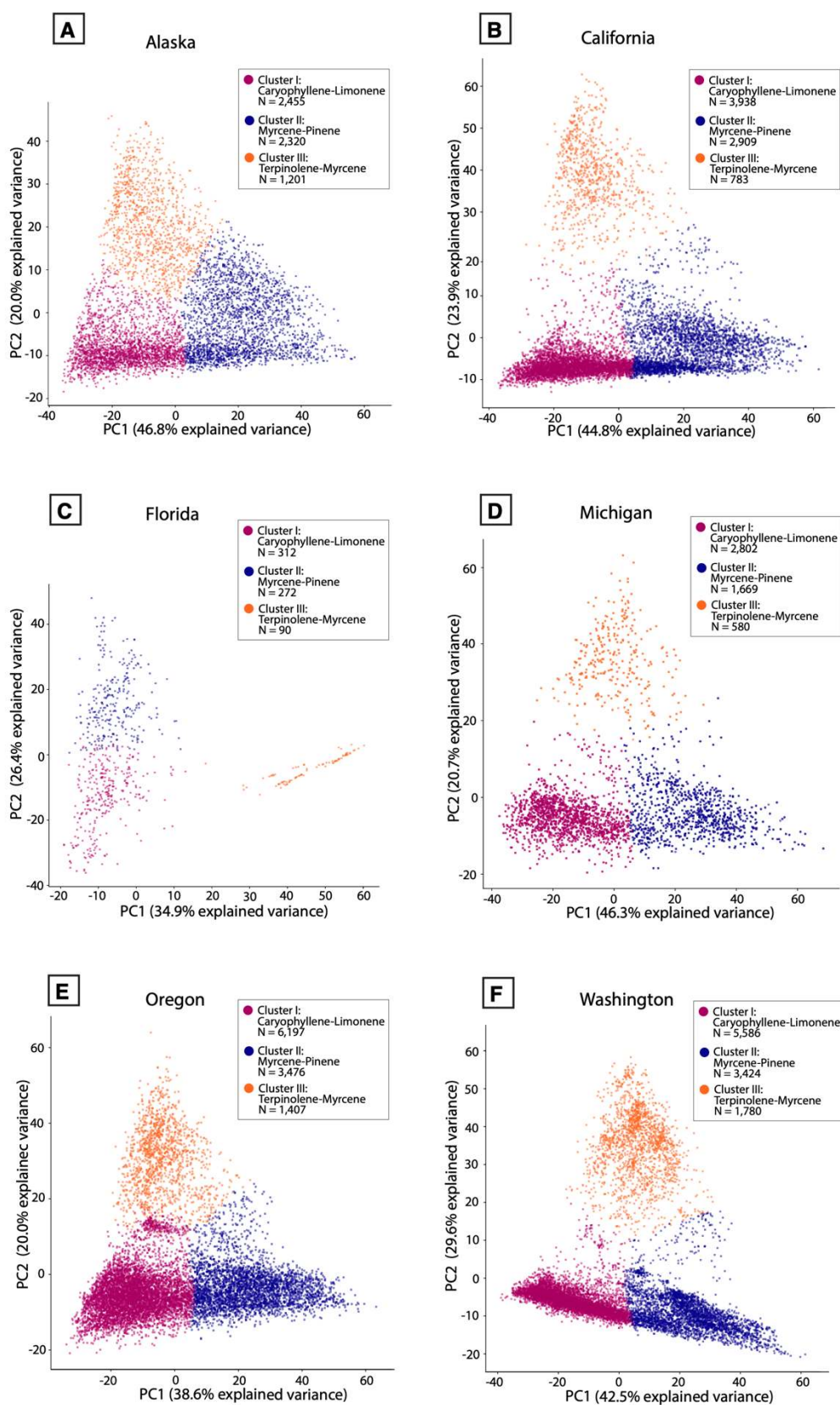176    **Figure S3:** Scatterplots showing the correlation between α- and β-pinene, by region. ***P < 0.0001

177

178   **Figure S4:** Scatterplots showing the correlation between β-caryophyllene and humulene, by region. ***P < 0.0001

43

**Figure S5: (A)** Line plot showing the relationship between number of clusters in k-means clustering and within-cluster sum of squared errors, using THC-dominant sample terpene data. "Elbow point" was determined to be at k=3. **(B)** PCA scores for all THC-dominant samples plotted along PC1 and PC2, color-coded by k-means cluster labels, k=2. **(C)** PCA scores for all THC-dominant samples plotted along PC1 and PC2, color-coded by k-means cluster labels, k=4.

**Figure S6:** PCA scores for THC-dominant samples plotted along PC1 and PC2, color-coded by k-means cluster labels attached to each sample, by region.